

# Enforcing Privacy in Distributed Learning with Performance Guarantees

Elsa Rizk\*, *Member, IEEE*, Stefan Vlaski†, *Member, IEEE*, and Ali H. Sayed\*, *Fellow, IEEE*.

**Abstract**—We study the privatization of distributed learning and optimization strategies. We focus on differential privacy schemes and study their effect on performance. We show that the popular additive random perturbation scheme degrades performance because it is not well-tuned to the graph structure. For this reason, we exploit two alternative graph-homomorphic constructions and show that they improve performance while guaranteeing privacy. Moreover, contrary to most earlier studies, the gradient of the risks is not assumed to be bounded (a condition that rarely holds in practice; e.g., quadratic risk). We avoid this condition and still devise a differentially private scheme with high probability. We examine optimization and learning scenarios and illustrate the theoretical findings through simulations.

**Index Terms**—distributed learning, privatized learning, differential privacy, distributed optimization

## I. INTRODUCTION

**D**ISTRIBUTED learning and optimization strategies are relevant in many contexts in real-world problems, such as in the design of robotic swarms for rescue missions, or the design of cloud computing services, or the exchange of information over social networks. Even in scenarios where a centralized solution is possible, it is often preferable to rely on a distributed implementation for various reasons. For instance, the centralized solution tends to have high maintenance costs and is sensitive to the failure of the central processor. Agents may also be reluctant to share their data with a remote central processor due to privacy and safety considerations. The amount of data available at each agent may be significant in size, which makes it difficult to regularly transmit large amounts of data between the dispersed agents and the central processor. Distributed implementations offer an attractive and robust alternative. The architecture can tolerate the failure of individual agents since processing can continue to occur among the remaining agents. Also, agents are only required to share minimal processed information with their neighbours.

There exist several schemes for distributed optimization, which have been studied extensively in the literature. Among these schemes we list the incremental strategy [1]–[8], consensus strategy [9]–[18], and diffusion strategy [19]–[24]. The incremental algorithm requires a renumbering of the agents over a cyclic path to cover the entire graph. This is usually a challenging task since the determination of an appropriate cycle is an NP-hard problem and, moreover, the failure of any edge along the path turns the solution moot.

The consensus and diffusion strategies avoid the need for a circular path over the graph. They rely on the local sharing of information among neighbouring agents. One main difference between both classes of strategies is that consensus updates are asymmetrical, where the starting point for the gradient-descent step is different from the location where the gradient vector is evaluated — see expression (16). It was shown in several earlier works (see, e.g., [20], [25], [26]) that this asymmetry reduces the stability range of consensus implementations in comparison to diffusion solutions, especially in scenarios involving the need for continuous learning and adaptation.

Now, one key aspect of distributed architectures is that they require agents to share information with their neighbours. This aspect raises an important privacy question about whether the information that is being shared over the edges in the graph can be intercepted. For instance, it is known that in algorithms that rely on gradient-descent updates, information leakage can occur through the sharing of the local gradients or the models that they estimate [27]–[30]. This can be problematic when the network is dealing with classified or sensitive data such as healthcare or financial data. In such cases, attackers may be able to recover certain elements of an individual’s personal information. There is no question that it is useful to pursue distributed strategies that guarantee a certain level of privacy.

There exists several useful works in the literature that address privacy questions for distributed algorithms. These contributions rely mainly on two types of tools: differential privacy or cryptography. Cryptographic methods range from using secure aggregation to multiparty computation and homomorphic encryption [31]–[35]. Although these methods do not hinder the performance of the learned model, they add significant computational and communication overhead.

On the other hand, differentially private methods mask the messages by adding some random noise [36]–[46]. They are simple to implement, but they introduce errors into the learned model and reduce the overall utility of the network. One main reason for this degradation is that the noise is often added at will *without* accounting for the graph topology.

In this work, we focus on differential privacy since it is simpler to apply and more scalable. We explain how to adjust its application to match the graph topology, while ensuring privacy and performance guarantees. In particular, we examine the effect of two differentially private schemes: the traditional random perturbations scheme and a graph-homomorphic scheme. We establish the superiority of the latter over the former in the mean-squared-error (MSE) sense. We also devise a third scheme, called *local* graph-homomorphic processing, which fully removes the degrading effect of the

\*The authors are with the School of Engineering, École Polytechnique Fédérale de Lausanne (e-mail: {elsa.rizk; ali.sayed}@epfl.ch). †The author is with the Department of Electrical and Electronic Engineering, Imperial College London (e-mail: s.vlaski@imperial.ac.uk).

noise on performance. These results apply to a broad class of distributed learning and optimization formulations.

## II. PROBLEM SETUP

We consider a graph topology with  $P$  agents, labelled  $p = 1, 2, \dots, P$ , as illustrated in Fig. 1. The objective is for the agents to approach the minimizer of an aggregate convex optimization problem of the form:

$$w^o \triangleq \underset{w \in \mathbb{R}^M}{\operatorname{argmin}} \frac{1}{P} \sum_{p=1}^P \left\{ J_p(w) \triangleq \frac{1}{N_p} \sum_{n=1}^{N_p} Q_p(w; x_{p,n}) \right\}, \quad (1)$$

where the risk function  $J_p(\cdot)$  is associated with the  $p$ th agent and is defined as an empirical average of the corresponding loss function  $Q_p(\cdot; \cdot)$  evaluated at the local data  $\{x_{p,n}\}_{n=1}^{N_p}$ . We associate two non-negative weights  $a_{mp}$  and  $a_{pm}$  with the edge linking neighbouring agents  $m$  and  $p$ . In this notation,  $a_{mp}$  is the weight used by agent  $p$  to scale information arriving from  $m$ , and similarly for  $a_{pm}$ ; it scales information from  $p$  toward  $m$ . The neighbourhood of an agent  $p$  is denoted by  $\mathcal{N}_p$  and consists of all agents that are connected to  $p$  by an edge. We assume that  $\mathcal{N}_p$  includes agent  $p$  as well.

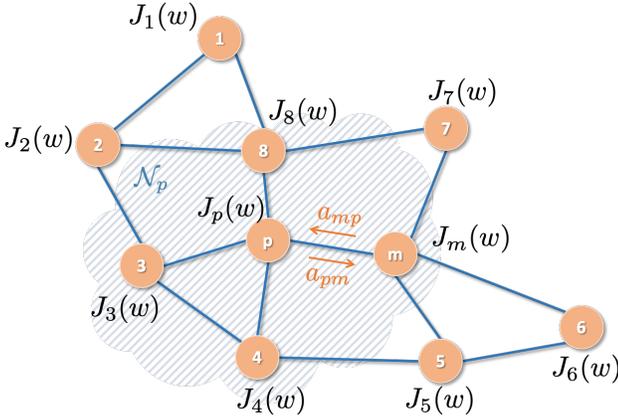


Fig. 1: Illustration of a network of agents.

### A. Modeling Conditions

We assume the individual risk functions  $J_p(w)$  are strongly convex and the loss functions  $Q_p(w; \cdot)$  have Lipschitz continuous gradients and are twice differentiable. These conditions are common in the study of distributed methods. Although the conditions can be relaxed and the results extended to broader scenarios (see, e.g., [25], [26], [47]–[49]), it is sufficient for the purposes of this work to illustrate the main ideas under these assumptions.

**Assumption 1 (Convexity and smoothness).** *The risks  $J_p(\cdot)$  are  $\nu$ -strongly convex, and the losses  $Q_p(\cdot; \cdot)$  are convex and twice differentiable, namely for some  $\nu > 0$ :*

$$J_p(w_2) \geq J_p(w_1) + \nabla_{w^\top} J_p(w_1)(w_2 - w_1) + \frac{\nu}{2} \|w_2 - w_1\|^2, \quad (2)$$

$$Q_p(w_2; \cdot) \geq Q_p(w_1; \cdot) + \nabla_{w^\top} Q_p(w_1; \cdot)(w_2 - w_1). \quad (3)$$

*The loss functions have  $\delta$ -Lipschitz continuous gradients, meaning there exists  $\delta > 0$  such that for any data point  $x_{p,n}$ :*

$$\|\nabla_{w^\top} Q_p(w_2; x_{p,n}) - \nabla_{w^\top} Q_p(w_1; x_{p,n})\| \leq \delta \|w_2 - w_1\|. \quad (4)$$

□

Since we assume the loss functions are twice differentiable, then the above strong-convexity and Lipschitz continuity conditions are equivalent to (see [25], [26], [49]):

$$0 < \nu I \leq \nabla_{w^\top}^2 J_p(w) \leq \delta I. \quad (5)$$

We further assume that the graph topology is strongly connected. This means that there exist paths linking any arbitrary pair of agents  $(m, p)$  in both directions and, moreover, at least one agent  $p$  in the network has a self-loop with  $a_{pp} > 0$ . In other words, at least one agent has some trust in its local information. The combination matrix  $A = [a_{mp}]$  is usually left-stochastic meaning that its entries satisfy:

$$a_{mp} \geq 0, \quad \sum_{m \in \mathcal{N}_p} a_{mp} = 1. \quad (6)$$

That is, the weights on edges connecting agents are nonnegative, and the entries on each column of  $A$  add up to one. The strong connectedness of the graph translates into guaranteeing that  $A$  is a primitive matrix. As a result, it follows from the Peron-Frobenius theorem [26] that  $A$  will have a single eigenvalue at one, while all other eigenvalues are strictly inside the unit circle. Moreover, an eigenvector  $q$  will exist with positive entries  $\{q_p\}$  adding up to one and satisfying:

$$Aq = q, \quad q_p > 0, \quad \mathbf{1}^\top q = 1. \quad (7)$$

We refer to  $q$  as the Peron eigenvector of  $A$ . Furthermore, it holds that  $\rho(A - q\mathbf{1}^\top) < 1$ , where  $\rho(\cdot)$  denotes the spectral radius of its matrix argument.

Next, let  $w^o$  denote the global minimizer for (1) and let  $w_p^o$  denote the local minimizer for  $J_p(\cdot)$ . We assume that the difference between these global and local models is bounded since, otherwise, collaboration would not be beneficial and one would instead follow a different optimization approach such as multi-task learning [50].

To clarify this point further, we consider a simple example involving a quadratic loss. Assume the data arriving at node  $p$ , denoted by  $\mathbf{d}_p(n)$ , is generated by some linear regression model under additive noise of the form:

$$\mathbf{d}_p(n) = \mathbf{u}_{p,n}^\top w^* + \mathbf{o}_p(n), \quad (8)$$

where  $\mathbf{u}_{p,n}$  is the feature vector and  $w^*$  is the model. We can seek to estimate  $w^*$  by solving:

$$\min_w \frac{1}{P} \sum_{p=1}^P \frac{1}{N_p} \sum_{n=1}^{N_p} (\mathbf{d}_p(n) - \mathbf{u}_{p,n}^\top w)^2. \quad (9)$$

The global minimizer in this case is given by:

$$w^o = w^* + \widehat{R}_u^{-1} \widehat{r}_{u,o}, \quad (10)$$

where:

$$\widehat{R}_u \triangleq \frac{1}{P} \sum_{p=1}^P \left\{ \widehat{R}_{p,u} \triangleq \frac{1}{N_p} \sum_{n=1}^{N_p} \mathbf{u}_{p,n} \mathbf{u}_{p,n}^\top \right\}, \quad (11)$$

$$\widehat{r}_{uo} \triangleq \frac{1}{P} \sum_{p=1}^P \left\{ \widehat{r}_{p,uo} \triangleq \frac{1}{N_p} \sum_{n=1}^{N_p} \mathbf{o}_p(n) \mathbf{u}_{p,n} \right\}, \quad (12)$$

while the local minimizers of  $J_p(w)$  are given by:

$$w_p^o = w^* + \widehat{R}_{p,u}^{-1} \widehat{r}_{p,uo}. \quad (13)$$

Thus, the global model (10) can be written as a weighted average of the local models (13):

$$w^o = \frac{1}{P} \sum_{p=1}^P \widehat{R}_u^{-1} \widehat{R}_{p,u} w_p^o. \quad (14)$$

This implies that the global model is a mixture of the local models. Therefore, the bound imposed below on the model difference amounts to an assumption on how different the distributions of the data across the agents are. This condition is weaker than a uniform bound on the difference between the gradients of the cost functions, which is more commonly assumed in the literature (see [51], [52]).

**Assumption 2 (Model drifts).** *The distance of each local model  $w_p^o$  to the global model  $w^o$  is uniformly bounded, i.e., there exists  $\xi \geq 0$  such that  $\|w^o - w_p^o\| \leq \xi$ .*  $\square$

### III. DISTRIBUTED LEARNING

#### A. Generalized Distributed Learning

We focus on two main strategies: consensus and diffusion. The consensus strategy for solving (1) takes the form:

$$\psi_{p,i-1} = \sum_{m \in \mathcal{N}_p} a_{mp} \mathbf{w}_{m,i-1}, \quad (15)$$

$$\mathbf{w}_{p,i} = \psi_{p,i-1} - \mu \widehat{\nabla}_{w^\top} J_p(\mathbf{w}_{p,i-1}), \quad (16)$$

where  $\widehat{\nabla}_{w^\top} J_p(\cdot)$  denotes a stochastic gradient *approximation* for the true gradient of  $J_p(\cdot)$ . Usually, the approximation is taken as the gradient of the loss function, namely,  $\nabla_{w^\top} Q_p(\mathbf{w}_{p,i-1}, \mathbf{x}_{p,i})$ . Here, the quantities  $\{\psi_{p,i}, \mathbf{w}_{p,i}\}$  denote estimates for  $w^o$  at node  $p$  at time  $i$ . Observe that the gradient vector in (16) is evaluated at the prior local model  $\mathbf{w}_{p,i-1}$  and not at the intermediate model  $\psi_{p,i-1}$ . The diffusion strategy, in turn, admits two related implementations known as combine-then-adapt (CTA) and adapt-then-combine (ATC). They differ by the order in which the calculations are performed with combination coming before adaptation in one case, and with the order reversed in the other case. The CTA diffusion strategy is described by:

$$\psi_{p,i-1} = \sum_{m \in \mathcal{N}_p} a_{mp} \mathbf{w}_{m,i-1}, \quad (17)$$

$$\mathbf{w}_{p,i} = \psi_{p,i-1} - \mu \widehat{\nabla}_{w^\top} J_p(\psi_{p,i-1}). \quad (18)$$

Comparing with (15)–(16), observe now that the starting point in (18) for the gradient-descent step is the same as the point

where the gradient vector is evaluated. Similarly, the ATC diffusion strategy is given by:

$$\psi_{p,i} = \mathbf{w}_{p,i-1} - \mu \widehat{\nabla}_{w^\top} J_p(\mathbf{w}_{p,i-1}), \quad (19)$$

$$\mathbf{w}_{p,i} = \sum_{m \in \mathcal{N}_p} a_{mp} \psi_{m,i}. \quad (20)$$

The above three algorithms can be combined into a single general description as follows [26]:

$$\phi_{p,i-1} = \sum_{m \in \mathcal{N}_p} a_{1,mp} \mathbf{w}_{m,i-1}, \quad (21)$$

$$\psi_{p,i} = \sum_{m \in \mathcal{N}_p} a_{0,mp} \phi_{m,i-1} - \mu \widehat{\nabla}_{w^\top} J_p(\phi_{p,i-1}), \quad (22)$$

$$\mathbf{w}_{p,i} = \sum_{m \in \mathcal{N}_p} a_{2,mp} \psi_{m,i}. \quad (23)$$

where we are introducing three combination matrices,  $\{A_0, A_1, A_2\}$ . By setting  $A_0 = A$  and  $A_1 = A_2 = I$ , we obtain consensus, while  $A_1 = A$  and  $A_0 = A_2 = I$  leads to CTA diffusion, and  $A_2 = A$  and  $A_0 = A_1 = I$  leads to ATC diffusion. Other choices are possible.

#### B. Privacy Learning

We now examine differentially private algorithms to safeguard the privacy of the information that is shared among the agents. For illustration purposes, assume the data  $\{x_{1,n}\}$  at agent 1 is replaced by a different set  $\{x'_{1,n}\}$ . The algorithm will thus take a new trajectory, which we denote by  $\{\phi'_{p,i-1}, \psi'_{p,i}, \mathbf{w}'_{p,i}\}$ . In a private implementation, an external observer should be oblivious to this change at agent 1. Concretely, all we need to do is add noise to the messages that need privatization. Most commonly, noise with exponential distributions, such as Laplacian or Gaussian, is added [36]. Thus, we are motivated initially to consider a privatized distributed implementation of the following form:

$$\phi_{p,i-1} = \sum_{m \in \mathcal{N}_p} a_{1,mp} (\mathbf{w}_{m,i-1} + \mathbf{g}_{1,mp,i}), \quad (24)$$

$$\psi_{p,i} = \sum_{m \in \mathcal{N}_p} a_{0,mp} (\phi_{m,i-1} + \mathbf{g}_{0,mp,i}) - \mu \widehat{\nabla}_{w^\top} J_p(\phi_{p,i-1}), \quad (25)$$

$$\mathbf{w}_{p,i} = \sum_{m \in \mathcal{N}_p} a_{2,mp} (\psi_{m,i} + \mathbf{g}_{2,mp,i}), \quad (26)$$

where the  $\mathbf{g}_{j,mp,i}$  denote zero-mean Laplacian random noises for  $j = 0, 1, 2$  for every  $m, p = 1, 2, \dots, P$ . For example, in (24), agent  $m$  shares  $\mathbf{w}_{m,i-1}$  with agent  $p$  over the edge that links them. During this transmission, an amount of Laplacian noise  $\mathbf{g}_{1,mp,i}$  is added. The subscript  $mp$  is used to denote that this noise is for the directed communication from  $m$  to  $p$ . Similarly, for the other noises.

We next define differential privacy formally [36], and show that the above algorithm is indeed differentially private.

**Definition 1 ( $\epsilon(i)$ -Differential privacy).** *We say that the algorithm given by (24)–(26) is  $\epsilon(i)$ -differentially private for agent  $p$  at time  $i$  if the following condition on the probabilities for observing the respective events holds on the*

joint distribution  $f(\cdot)$  where the notation  $\mathbf{y}_{p,j-1}$  represents any of the shared messages  $\{\mathbf{w}_{p,j-1}, \phi_{p,j}, \psi_{p,j}\}$ ,  $\mathbf{g}_{\cdot,pm,j}$  the corresponding added noise  $\{\mathbf{g}_{1,pm,j}, \mathbf{g}_{0,pm,j}, \mathbf{g}_{2,pm,j}\}$ :

$$\frac{f\left(\left\{\left\{\mathbf{y}_{p,j-1} + \mathbf{g}_{\cdot,pm,j}\right\}_{m \in \mathcal{N}_p \setminus \{p}\right\}_{j=1}^i\right)}{f\left(\left\{\left\{\mathbf{y}'_{p,j-1} + \mathbf{g}_{\cdot,pm,j}\right\}_{m \in \mathcal{N}_p \setminus \{p}\right\}_{j=1}^i\right)} \leq e^{\epsilon(i)}. \quad (27)$$

□

The above bounds ensure that for small  $\epsilon(i)$ , the distributions of the original and modified trajectories are close to each other. This makes it difficult to infer information about the data at the agents since we cannot distinguish the trajectories of the algorithm for different combinations of participating agents. In other words, if agent 1 chooses to replace its original dataset by  $\{x'_{1,n}\}$ , then the resulting models  $\{\mathbf{w}'_{p,j-1}, \phi'_{p,j}, \psi'_{p,j}\}$  are close enough in distribution to the original models  $\{\mathbf{w}_{p,j-1}, \phi_{p,j}, \psi_{p,j}\}$ , and an outside observer will not be able to conclude what dataset was used. The two model trajectories resulting from the use of the original and the alternative dataset are indistinguishable.

To show that algorithm (24)–(26) satisfies condition (27), we first calculate the sensitivity of the algorithm. The sensitivity at time  $i$  is defined in Appendix A as the change in the trajectory of the algorithm if instead of using the original dataset, agent 1 uses the alternative dataset  $\{x'_{1,n}\}$ . In Appendix A the sensitivity is shown to satisfy:

$$\Delta(i) \triangleq \|\mathbf{w}_i - \mathbf{w}'_i\| \leq B + B' + \sqrt{P}\|w^o - w'^o\|, \quad (28)$$

for some constants  $B$  and  $B'$  and with high probability. That is, it holds that:

$$\begin{aligned} & \mathbb{P}(\Delta(i) \leq B + B' + \sqrt{P}\|w^o - w'^o\|) \\ & \geq \left( 1 - \frac{\kappa_2^2 \mathbf{1}^\top \Gamma^i \left[ \frac{\mathbb{E}\|\tilde{\mathbf{w}}_0\|^2}{\mathbb{E}\|\tilde{\mathbf{v}}_0\|^2} \right] + O(\mu) + O(\mu^{-1})}{B^2} \right) \\ & \times \left( 1 - \frac{\kappa_2'^2 \mathbf{1}^\top \Gamma'^i \left[ \frac{\mathbb{E}\|\tilde{\mathbf{w}}'_0\|^2}{\mathbb{E}\|\tilde{\mathbf{v}}'_0\|^2} \right] + O(\mu) + O(\mu^{-1})}{B'^2} \right), \quad (29) \end{aligned}$$

where  $\mathbf{w}_i \triangleq \text{col}\{\mathbf{w}_{p,i}\}_{p=1}^P$ , the model error at time zero is denoted by  $\tilde{\mathbf{w}}_0 \triangleq \text{col}\{w^o - \mathbf{w}_{p,0}\}_{p=1}^P$ , and the variables  $\{\tilde{\mathbf{w}}_0, \tilde{\mathbf{v}}_0\}$  arise from the partitioning  $\mathcal{V}_\theta^\top \tilde{\mathbf{w}}_0 = \text{col}\{\tilde{\mathbf{w}}_0, \tilde{\mathbf{v}}_0\}$  with the matrix  $\Gamma$  and the constant  $\kappa_2$  defined in Appendix B. Result (29) means that the sensitivity  $\Delta(i)$  is bounded with high probability. The bound constants  $B$  and  $B'$  are chosen by the user: larger values for  $B$  and  $B'$  result in higher probability of bounded sensitivity but, as shown in (31), they result in a larger privacy bound. In other words, the values of  $B$  and  $B'$  can be controlled to balance the trade-off between the privacy level and the likelihood of bounded sensitivity. Next, if we denote the variance of the Laplacian noise  $\mathbf{g}_{j,mp,i}$  by  $\sigma_g^2$ ,

with  $\mathbf{y}_{p,j-1} = \mathbf{w}_{p,j-1}$  the fraction in (27) can be bounded as follows with high probability:

$$\begin{aligned} & \frac{f\left(\left\{\left\{\mathbf{w}_{p,j-1} + \mathbf{g}_{0,pm,j}\right\}_{m \in \mathcal{N}_p \setminus \{p}\right\}_{j=1}^i\right)}{f\left(\left\{\left\{\mathbf{w}'_{p,j-1} + \mathbf{g}_{0,pm,j}\right\}_{m \in \mathcal{N}_p \setminus \{p}\right\}_{j=1}^i\right)} \\ & \stackrel{(a)}{=} \prod_{j=1}^i \frac{f\left(\left\{\left\{\mathbf{w}_{p,j-1} + \mathbf{g}_{0,pm,j}\right\}_{m \in \mathcal{N}_p \setminus \{p}\right\} | \mathcal{X}_{j-1}\right)}{f\left(\left\{\left\{\mathbf{w}'_{p,j-1} + \mathbf{g}_{0,pm,j}\right\}_{m \in \mathcal{N}_p \setminus \{p}\right\} | \mathcal{X}'_{j-1}\right)} \\ & \stackrel{(b)}{=} \prod_{j=1, m \in \mathcal{N}_p \setminus \{p\}}^i \frac{\exp(-\sqrt{2}\|\mathbf{w}_{p,j-1} + \mathbf{g}_{0,pm,j}\|/\sigma_g)}{\exp(-\sqrt{2}\|\mathbf{w}'_{p,j-1} + \mathbf{g}_{0,pm,j}\|/\sigma_g)} \\ & \leq \exp\left(\frac{\sqrt{2}}{\sigma_g} \sum_{j=1, m \in \mathcal{N}_p \setminus \{p\}}^i \|\mathbf{w}_{p,j-1} - \mathbf{w}'_{p,j-1}\|\right) \\ & \leq \exp\left(\frac{\sqrt{2}P}{\sigma_g} \sum_{j=1}^i \|\mathbf{w}_{j-1} - \mathbf{w}'_{j-1}\|\right), \quad (30) \end{aligned}$$

where the first equality (a) follows from applying Bayes' rule with  $\mathcal{X}_{j-1} \triangleq \{\mathbf{w}_{p,j-1}\} \cup \left\{ \left\{ \mathbf{w}_{p,o-1} + \mathbf{g}_{0,pm,o} \right\}_{m \in \mathcal{N}_p \setminus \{p\}} \right\}_{o=1}^{j-1}$ , and the second equality (b) follows from the independence of  $\mathbf{w}_{p,j-1} + \mathbf{g}_{0,pm,j}$  for  $m \in \mathcal{N}_p \setminus \{p\}$  conditioned on  $\mathbf{w}_{p,j-1}$ . A similar bound can be found for  $\mathbf{y}_{p,j-1} \in \{\phi_{p,j}, \psi_{p,j}\}$ .

Thus, the level of privacy is defined by the following choice for  $\epsilon(i)$  in terms of the running  $\Delta(j)$  values:

$$\epsilon(i) = \frac{\sqrt{2}P}{\sigma_g} \sum_{j=0}^{i-1} \Delta(j) \leq \frac{\sqrt{2}P}{\sigma_g} (B + B' + \sqrt{P}\|w^o - w'^o\|)i. \quad (31)$$

These results show that in order to arrive at an  $\epsilon(i)$ -differentially private algorithm, it is sufficient to select the variance of the Laplacian noise to satisfy (31). Expression (31) shows that  $\epsilon(i)$  is a linear function of the iterations. This means that the process becomes less private at a rate no greater than a linear rate. It is important to note here that most earlier studies on differentially private schemes for multi-agent systems [38], [41], [42] assume bounded gradients for the risk function. However, this condition is rarely satisfied in practice. For instance, even quadratic risks have unbounded gradients. For this reason, in our approach, we have avoided relying on this assumption. Instead, we are able to establish that differential privacy continues to hold with high probability.

We still need to examine the effect of the added noises on performance. To do so, we introduce the extended model  $\mathbf{w}_i \triangleq \text{col}\{\mathbf{w}_{p,i}\}_{p=1}^P$  and write the three-step algorithm (24)–(26) using one single recursion as follows:

$$\begin{aligned} \mathbf{w}_i & = \mathcal{A}_2^\top \mathcal{A}_0^\top \mathcal{A}_1^\top \mathbf{w}_{i-1} - \mu \mathcal{A}_2^\top \text{col}\left\{\widehat{\nabla_{\mathbf{w}^\top} J_p}(\phi_{p,i-1})\right\}_{p=1}^P \\ & \quad + \mathcal{A}_2^\top \mathcal{A}_0^\top \text{diag}(\mathcal{A}_1^\top \mathcal{G}_{1,i}) + \mathcal{A}_2^\top \text{diag}(\mathcal{A}_0^\top \mathcal{G}_{0,i}) \\ & \quad + \text{diag}(\mathcal{A}_2^\top \mathcal{G}_{2,i}), \quad (32) \end{aligned}$$

where for  $j = 0, 1, 2$ ,  $\mathcal{A}_j \triangleq A_j \otimes I_M$ , and  $\mathcal{G}_{j,i}$  is a matrix whose entries are the added noises  $\mathbf{g}_{j,mp,i}$ . We denote the

model error by  $\tilde{\mathbf{w}}_i \triangleq \text{col}\{w^o - \mathbf{w}_{p,i}\}_{p=1}^P$ , and introduce the local gradient noise:

$$\mathbf{s}_{p,i} \triangleq \widehat{\nabla_{w^\tau} J_p}(\phi_{p,i-1}) - \nabla_{w^\tau} J_p(\phi_{p,i-1}). \quad (33)$$

It is customary to assume that this gradient noise process has zero mean and bounded second-order moment (see, e.g., [26], [49], where this property is actually shown to hold in many important cases of interest and similar arguments can be applied to the current case), namely:

$$\mathbb{E}\{\|\mathbf{s}_{p,i}\|^2 | \mathcal{F}_{i-1}\} \leq \beta_{s,p}^2 \|\tilde{\phi}_{p,i-1}\|^2 + \sigma_{s,p}^2, \quad (34)$$

for some nonnegative constants  $\beta_{s,p}^2$  and  $\sigma_{s,p}^2$ , and where the conditioning is taken over all past models  $\mathcal{F}_{i-1} \triangleq \text{filtration}\{\mathbf{w}_{p,j}\}_{p=1, j=0}^{P, i-1}$ . Then, using the extended gradient noise  $\mathbf{s}_i \triangleq \text{col}\{\mathbf{s}_{p,i}\}_{p=1}^P$ , the error recursion corresponding to (32) is given by:

$$\begin{aligned} \tilde{\mathbf{w}}_i &= \mathcal{A}_2^\top \mathcal{A}_0^\top \mathcal{A}_1^\top \tilde{\mathbf{w}}_{i-1} + \mu \mathcal{A}_2^\top \text{col}\{\nabla_{w^\tau} J_p(\phi_{p,i-1})\}_{p=1}^P \\ &\quad + \mu \mathcal{A}_2^\top \mathbf{s}_i - \mathcal{A}_2^\top \mathcal{A}_0^\top \text{diag}(\mathcal{A}_1^\top \mathcal{G}_{1,i}) - \mathcal{A}_2^\top \text{diag}(\mathcal{A}_0^\top \mathcal{G}_{0,i}) \\ &\quad - \text{diag}(\mathcal{A}_2^\top \mathcal{G}_{2,i}). \end{aligned} \quad (35)$$

Since  $J_p(\cdot)$  are twice differentiable, we appeal to the mean-value theorem to express the gradient in the form [26]:

$$\nabla_{w^\tau} J_p(\phi_{p,i-1}) = -\mathbf{H}_{p,i-1} \tilde{\phi}_{p,i-1} - \nabla_{w^\tau} J_p(w^o), \quad (36)$$

where:

$$\mathbf{H}_{p,i-1} \triangleq \int_0^1 \nabla_{w^\tau}^2 J_p(w^o - t\phi_{p,i-1}) dt. \quad (37)$$

Then, introducing the quantities:

$$\mathcal{B}_{i-1} \triangleq \mathcal{A}_2^\top (\mathcal{A}_0^\top - \mu \mathcal{H}_{i-1}) \mathcal{A}_1^\top, \quad (38)$$

$$\mathcal{H}_{i-1} \triangleq \text{diag}\{\mathbf{H}_{p,i-1}\}_{p=1}^P, \quad (39)$$

$$\mathbf{b} \triangleq \text{col}\{\nabla_{w^\tau} J_p(w^o)\}_{p=1}^P, \quad (40)$$

we rewrite (35) as:

$$\begin{aligned} \tilde{\mathbf{w}}_i &= \mathcal{B}_{i-1} \tilde{\mathbf{w}}_{i-1} + \mu \mathcal{A}_2^\top \mathbf{s}_i - \mu \mathcal{A}_2^\top \mathbf{b} + \mu \mathcal{A}_2^\top \mathcal{H}_{i-1} \text{diag}(\mathcal{A}_1^\top \mathcal{G}_{1,i}) \\ &\quad - \text{diag}(\mathcal{A}_2^\top \mathcal{G}_{2,i}) - \mathcal{A}_2^\top \text{diag}(\mathcal{A}_0^\top \mathcal{G}_{0,i}) - \mathcal{A}_2^\top \mathcal{A}_0^\top \text{diag}(\mathcal{A}_1^\top \mathcal{G}_{1,i}). \end{aligned} \quad (41)$$

We show in the next theorem that the weight-error size converges to the neighbourhood of zero, with the size of the neighbourhood determined by the step-size and the added noise variance.

**Theorem 1 (MSE convergence of privatized distributed learning).** *Under assumptions 1 and 2, the distributed recursions (24)–(26) converge exponentially fast for a small enough step-size to a neighbourhood of the optimal model, i.e.:*

$$\limsup_{i \rightarrow \infty} \mathbb{E}\|\tilde{\mathbf{w}}_i\|^2 \leq O(\mu)\sigma_s^2 + O(\mu) + (O(\mu^{-1}) + O(\mu))\sigma_g^2. \quad (42)$$

*Proof.* See Appendix B.  $\square$

By examining the bound in (42) on the mean-square error (MSE), we observe that the noise variance  $\sigma_g^2$  appears multiplied by a term on the order of  $\mu^{-1}$ , which is detrimental

to performance when  $\mu$  is small. Therefore, the traditional approach of adding Laplacian noise over the edges to ensure privacy is *calamitous* to performance and needs to be improved. We describe next an alternative approach.

### C. Graph-Homomorphic Noise

The noises added to the communication links in the previous section did not take into account the graph topology. As a result, their effect gets magnified by  $O(\mu^{-1})$  as shown in (42). We now examine another strategy for adding noise, which relies on a graph-homomorphic construction from [46]. Specifically, the noises are now constructed to satisfy the following condition:

$$\sum_{p,m=1}^P q_p a_{mp} \mathbf{g}_{j,mp,i} = 0, \quad (43)$$

for  $j = 0, 1, 2$ , and where  $q = \text{col}\{q_p\}_{p=1}^P$  is the Perron eigenvector of  $\mathcal{A}_2^\top \mathcal{A}_0^\top \mathcal{A}_1^\top$ . This can be satisfied if we continue to choose zero-mean Laplacian noises  $\mathbf{g}_{j,p,i}$  with variance  $\sigma_g^2$  and then set:

$$\mathbf{g}_{j,pm,i} = \begin{cases} \frac{a_{pm}}{a_{mp}} \mathbf{g}_{j,p,i}, & m \neq p \\ -\frac{1-a_{j,pp}}{a_{j,pp}} \mathbf{g}_{j,p,i}, & m = p. \end{cases} \quad (44)$$

Condition (43), along with construction (44), ensure that the net effect of the additional noises cancel out over the entire graph during the local aggregation steps. We show in the next theorem that the MSE bound improves in this case. To see this, we first introduce the network centroid  $\mathbf{w}_{c,i}$  and study its convergence under graph-homomorphic perturbations. Let:

$$\begin{aligned} \mathbf{w}_{c,i} &\triangleq \sum_{p=1}^P q_p \mathbf{w}_{p,i} \\ &= \mathbf{w}_{c,i-1} - \mu \sum_{p=1}^P q_p \mathbf{s}_{p,i} \\ &\quad - \mu \sum_{p=1}^P q_p \nabla_{w^\tau} J_p \left( \sum_{m \in \mathcal{N}_p} a_{1,mp} (\mathbf{w}_{m,i-1} + \mathbf{g}_{1,mp,i}) \right) \\ &\quad + \sum_{p,m} q_p (a_{1,mp} \mathbf{g}_{1,mp,i} + a_{0,mp} \mathbf{g}_{0,mp,i} + a_{2,mp} \mathbf{g}_{2,mp,i}). \end{aligned} \quad (45)$$

Since we are using graph-homomorphic perturbations, the sum of the noise terms in the last line cancels out. We can therefore write the following error recursion:

$$\begin{aligned} \tilde{\mathbf{w}}_{c,i} &= \tilde{\mathbf{w}}_{c,i-1} + \mu (q^\top \otimes I) \mathbf{s}_i + \mu (q^\top \otimes I) \mathbf{b} \\ &\quad - \mu \sum_{p=1}^P q_p \mathbf{H}_{p,i-1} \sum_{m \in \mathcal{N}_p} a_{1,mp} (\tilde{\mathbf{w}}_{m,i-1} - \mathbf{g}_{1,mp,i}). \end{aligned} \quad (46)$$

Before stating the theorem on the MSE convergence, we bound the network disagreement defined as the average second-order moment of the difference between the local models and the centroid model.

**Lemma 1 (Network disagreement).** *The average deviation from the centroid is uniformly bounded during each iteration  $i$ , and, moreover, it holds asymptotically that:*

$$\limsup_{i \rightarrow \infty} \frac{1}{P} \sum_{p=1}^P \mathbb{E} \|\mathbf{w}_{p,i} - \mathbf{w}_{c,i}\|^2 \leq O(1)\sigma_g^2 + O(\mu) \quad (47)$$

*Proof.* See Appendix C.  $\square$

Expression (47) shows that the local models will be at most  $O(1)\sigma_g^2$  away from the centroid model. Thus, if the centroid model manages to converge to the optimal model  $w^o$  with only a slight variation, then the local models will always be a constant, proportional to the noise variance  $\sigma_g^2$ , away from the true model. In the next theorem, we show that the added noise only alters the centroid model by an  $O(1)$  factor.

**Theorem 2 (MSE convergence of the network centroid).** *Under assumptions 1 and 2, the network centroid defined in (45) converges exponentially fast for a small enough step-size to a neighbourhood of the optimal model:*

$$\limsup_{i \rightarrow \infty} \mathbb{E} \|\tilde{\mathbf{w}}_{c,i}\|^2 \leq O(\mu)\sigma_s^2 + O(1)\sigma_g^2 + O(\mu^2). \quad (48)$$

*Proof.* See Appendix D.  $\square$

Thus, the network centroid is at most  $O(1)\sigma_g^2$  away from the true minimizer  $w^o$ , even with added noise, as opposed to  $O(\mu^{-1})\sigma_g^2$ . In Lemma 1, we showed that the individual models  $\mathbf{w}_{p,i}$  are  $O(1)\sigma_g^2$  away from the centroid model. Thus, by using the graph-homomorphic perturbations (43)–(44), the MSE is not inversely proportional to  $\mu$  anymore, which is an improvement relative to (42).

#### D. Local Graph-Homomorphic Noise

We explain how to improve on the  $O(1)\sigma_g^2$  deviation and replace it by  $O(\mu)\sigma_g^2$ , by relying on the use of *local* graph-homomorphic noise. To do so, we construct the noises to satisfy the following alternative condition as opposed to (43):

$$\sum_{m \in \mathcal{N}_p} a_{mp} \mathbf{g}_{j,mp,i} = 0. \quad (49)$$

Observe that we are requiring the sum of the noises to cancel out *locally*, rather than globally as required in the previous section. The neighbours of every agent  $p$  must collaborate together to generate dependent random noises that will cancel out locally at  $p$ . The collaboration will occur through agent  $p$ , since a direct link might not exist amongst these neighbours. A similar problem exists in blockchain applications where the generation of a random number is required to occur in a distributed manner [53].

We now devise a distributed scheme that leads to noises that satisfy condition (49). For the sake of demonstration, we describe the protocol through an example. Thus, assume agent 1 is connected to 5 agents labeled 2, 3, 4, 5, 6 (Fig. 2, left). Since the neighbours of agent 1 need not be connected to each other through direct links, all communications will take place through agent 1. In this subnetwork, we allow agent 1 to be the orchestrator of the scheme. The first step is for agent 1 to split its neighbours into two disjoint sets  $\mathcal{N}_1 = \mathcal{N}_+ \cup \mathcal{N}_-$ .

For example, we may collect the even numbered agents into  $\mathcal{N}_+$ , and the odd numbered agents into  $\mathcal{N}_-$ . Then, we allow every pair of agents from the two disjoint sets to agree on a noise value they will add to their message such that they will cancel out at agent 1. We force agents from  $\mathcal{N}_-$  to multiply the noise they will add to their messages by a negative sign. Therefore, agent 2 will add to its message two noise terms, one generated with agent 3 and another with agent 5. We denote the noise term generated by agents 2 and 3 that will be sent to agent 1 by  $\mathbf{g}_{\{23\}1,i}$ . Since the messages are scaled by the weights attributed by agent 1 to its neighbours, the added noise must then be divided by the weights, i.e., the message sent by agent 2 to agent 1 is the original message  $\mathbf{w}_{2,i}$  and the two generated noises by agent 2 with agents 3 and 5 scaled by the corresponding weight  $a_{12}$ :

$$\mathbf{w}_{2,i} + \frac{\mathbf{g}_{\{23\}1,i}}{a_{12}} + \frac{\mathbf{g}_{\{25\}1,i}}{a_{12}}. \quad (50)$$

However, this requires that agent 2 know the weight attributed to its messages by agent 1. Thus, agent 1 will have to make the weights public in case of a non-doubly stochastic combination matrix. The same process occurs between agents 4 and 6 with both 3 and 5. Then, the aggregate messages sent to agent 1 will end up being the sum of the unmasked weights:

$$\begin{aligned} & \sum_{k \in \mathcal{N}_+} a_{1k} \left( \mathbf{w}_{k,i} + \sum_{\ell \in \mathcal{N}_-} \frac{\mathbf{g}_{\{k\ell\}1,i}}{a_{1k}} \right) \\ & + \sum_{k \in \mathcal{N}_-} a_{1k} \left( \mathbf{w}_{k,i} - \sum_{\ell \in \mathcal{N}_+} \frac{\mathbf{g}_{\{k\ell\}1,i}}{a_{1k}} \right) = \sum_{k \in \mathcal{N}_1} a_{1k} \mathbf{w}_{k,i}. \end{aligned} \quad (51)$$

We move to the method used to generate the pairwise noise terms  $\mathbf{g}_{\{k\ell\}1,i}$ . We rely on the Diffie-Helmann key exchange protocol where each pair of agents shares a secret key that is used to generate the added noise. Given two agents, say 2 and 3, we assume they have individual secret keys  $v_2$  and  $v_3$ , respectively. A known modulus  $\pi$  and base  $b$  is agreed upon amongst the agents. Then, agent 2 broadcasts its public key  $V_2 = (b^{v_2} \bmod \pi)$  and agent 3 does the same  $V_3 = (b^{v_3} \bmod \pi)$ . Agent 2 then calculates,  $v_{23} = (V_3^{v_2} \bmod \pi) = (b^{v_2 v_3} \bmod \pi)$  which is the same as what agent 3 calculates  $v_{23} = (V_2^{v_3} \bmod \pi) = (b^{v_3 v_2} \bmod \pi)$ . Thus, the two agents now share a secret key  $v_{23}$  only known to them. This secret key can then be used as the added noise to mask the messages, i.e.,  $\mathbf{g}_{\{23\}1,i} = v_{23}$ . However, to make the process differentially private we need the resulting added noise to be Laplacian,  $\text{Lap}(0, \sigma_g/\sqrt{2})$ . In what follows, we describe a scheme, which we call the local graph-homomorphic processing scheme that ensures the added noise is Laplacian. An illustration of this process is found in the right subfigure of Fig. 2.

**Definition 2 (Local graph-homomorphic process).** *We are given a subnetwork of agent  $k$ , and neighbours  $\ell \in \mathcal{N}_+$  and  $m \in \mathcal{N}_-$ . Let agent  $\ell$  sample two secret keys  $v_\ell$  and  $v'_\ell$  from a uniform distribution on  $[0, 1]$ , and let agent  $m$  sample its keys  $v_m$  and  $v'_m$  from a gamma distribution  $\Gamma(2, 1)$ . Let  $\pi$  be some large prime number and let  $a$  be a multiple of  $\pi$ . Then, for:*

$$\mathbf{v}_{\ell m} = a e^{-v_\ell v_m} \bmod \pi, \quad (52)$$

$$\mathbf{v}'_{\ell m} = a e^{-v'_\ell v'_m} \bmod \pi, \quad (53)$$

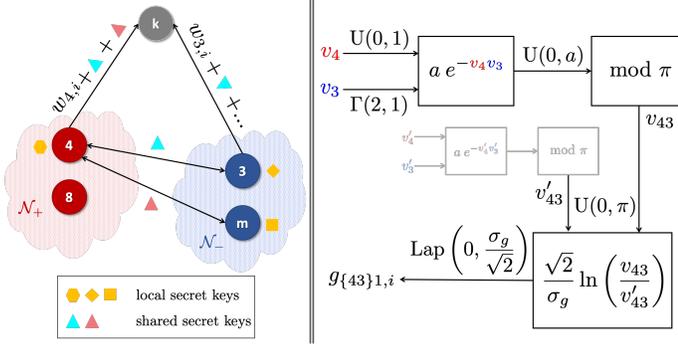


Fig. 2: Illustration of the local graph-homomorphic process. The figure on the left describes the Diffie-Helman key exchange procedure. The figure on the right shows the transformation the random variable goes through.

the desired Laplacian noise can be constructed as:

$$\mathbf{g}_{\{\ell m\}k,i} = \frac{\sqrt{2}}{\sigma_g} \ln \left( \frac{\mathbf{v}_{\ell m}}{\mathbf{v}'_{\ell m}} \right). \quad (54)$$

□

The local graph-homomorphic process proposed here is related to methods that fall under secure aggregation (see [31]). However, the main difference between our method and earlier investigations is that we devise a scheme for the more general *distributed* setting, while other works focus largely on the particular case of federated learning with its specialized structure with a central processor. Furthermore, while we generate random numbers making our scheme more secure, the work [31] adds pseudo-random numbers to the shared messages. Since pseudo-random numbers are generated by deterministic algorithms, it makes the noise predictable and susceptible to attacks, contrary to random numbers. Furthermore, we quantify the privacy of our scheme as opposed to [31]. In the next theorem, we show that using construction (54) results in a differentially private algorithm.

**Theorem 3 (Privacy of distributed learning with local graph-homomorphic perturbations).** *Under the local graph-homomorphic process defined by (54) the resulting privatized algorithm is  $\epsilon(i)$ -differentially private with high probability and with  $\epsilon(i)$  defined in (31).*

*Proof.* See Appendix E. □

Note that since the noises cancel out locally during each iteration, the algorithm follows the same trajectory as the non-privatized algorithm. This implies that the MSD performance of the privatized distributed learning algorithm with the local graph-homomorphic perturbation (54) is the same as the non-privatized distributed learning algorithm. In particular, the results on the convergence of the non-privatized algorithm from Theorem 9.1 in [26] will continue to hold. This implies that the MSE will now be on the order of  $O(\mu)\sigma_g^2$ .

The above result highlights one difficulty with differential privacy. Note that, in principle, as the variance  $\sigma_g^2$  of the added noise is increased, the level of privacy is also increased. However, this process introduces an additional communication

cost. For example, agent 1 needs to communicate to its neighbourhood the splitting into the positive agents and negative agents. It will also need to communicate with almost half the neighbourhood of its neighbours to agree on an added noise  $\mathbf{g}_{\{1m\}k,i}$  for  $k \in \mathcal{N}_1$  and  $m \in \mathcal{N}_k$ . Thus, in total, the communication cost for agent 1 will increase by at most  $|\mathcal{N}_1| + \sum_{k=2}^6 |\mathcal{N}_k|/2$ . The additional communication cost is not captured by the privacy measure even though it clearly affects the level of privacy<sup>1</sup>. This reduces its functionality and motivates the search for other privacy metrics, such as those based on information-theoretic measures [55]. For example, one metric could be the mutual information between the original message and the perturbed shared message [54]. Thus, if we assume the individual messages  $\{\mathbf{w}_{k,i}\}$  are Gaussian random variables with variance  $\sigma_w^2$ , and if we perturb them with a total Gaussian noise  $\mathbf{g}_{k,i}$  of variance  $\sigma_g^2$ , then the mutual information is given by:

$$\begin{aligned} I(\mathbf{w}_{k,i}; \mathbf{w}_{k,i} + \mathbf{g}_{k,i}) &= H(\mathbf{w}_{k,i} + \mathbf{g}_{k,i}) - H(\mathbf{w}_{k,i} + \mathbf{g}_{k,i} | \mathbf{w}_{k,i}) \\ &= H(\mathbf{w}_{k,i} + \mathbf{g}_{k,i}) - H(\mathbf{g}_{k,i}) \\ &= \frac{1}{2} \log \left( 1 + \frac{\sigma_w^2}{\sigma_g^2} \right). \end{aligned} \quad (55)$$

Observe again that as we increase the noise variance  $\sigma_g^2$ , mutual information decreases while privacy increases. Mutual information again fails to capture the communication cost incurred by the process. It appears that no metric capturing the communication-privacy trade-off exists as of yet in the literature. This calls for the search of a more appropriate privacy metric for secure aggregation methods.

#### IV. EXPERIMENTAL ANALYSIS

We run two experiments. In the first experiment we focus on a linear regression problem with simulated data. We then study a classification problem on real data.

##### A. Generalized Distributed Privacy Learning

For each of consensus, CTA, and ATC diffusion, we compare four algorithms: the standard distributed algorithm, the privatized algorithm with random perturbations, the privatized algorithm with graph-homomorphic perturbations, and the privatized algorithm with local graph-homomorphic perturbations. We consider a network of 30 agents (Fig. 3) and a regularized quadratic loss function:

$$\min_{\mathbf{w} \in \mathbb{R}^2} \frac{1}{30} \sum_{p=1}^{30} \frac{1}{100} \sum_{n=1}^{100} (\mathbf{d}_p(n) - \mathbf{u}_{p,n}^\top \mathbf{w})^2 + 0.01 \|\mathbf{w}\|^2. \quad (56)$$

We generate a random dataset  $\{\mathbf{u}_{p,n}, \mathbf{d}_p(n)\}_{n=1}^{100}$  as follows: we let the two-dimensional feature vector  $\mathbf{u}_{p,n} \sim \mathcal{N}(0; R_u)$ , and add noise  $\mathbf{o}_p(n) \sim \mathcal{N}(0; \sigma_{o,p}^2)$  such that  $\mathbf{d}_p(n) = \mathbf{u}_{p,n}^\top \mathbf{w}^* + \mathbf{o}_p(n)$ , for some generative model  $\mathbf{w}^* \in \mathbb{R}^2$  and randomly set variance  $R_u$  and added noise variance  $\sigma_{o,p}^2$ . To make the data distributions non-iid, we use different noise

<sup>1</sup>If we were to decrease the number of times a random noise is generated by the local graph-homomorphic process and instead re-use the noise, then we would be decreasing the communication cost but increasing the chance of an attacker learning the noise and unmasking the messages.

variances  $\sigma_{o,p}^2$  across the agents. The optimal global model has a closed form solution with  $\widehat{R}_u$  and  $\widehat{r}_{uo}$  as defined previously:

$$w^o = (\widehat{R}_u + 0.01I)^{-1}(\widehat{R}_u w^* + \widehat{r}_{uo}). \quad (57)$$

We set the step-size  $\mu = 0.4$ , the noise variance  $\sigma_g^2 = 0.01$ , and run the algorithms for 1000 iterations. We repeat the experiment 20 times and plot the MSD in the log domain of the centroid model and the average of the individual MSDs:

$$\text{MSD}_i \triangleq \|\mathbf{w}_{c,i} - w^o\|^2, \quad (58)$$

$$\text{MSD}_{\text{avg},i} \triangleq \frac{1}{P} \sum_{p=1}^P \|\mathbf{w}_{p,i} - w^o\|^2. \quad (59)$$

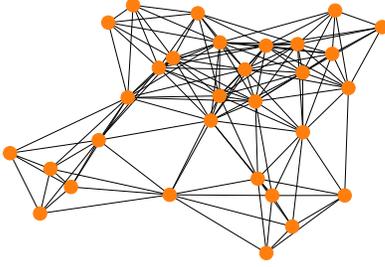


Fig. 3: The generated network of agents.

As we observe in Fig. 4, graph-homomorphic perturbations do not hinder the performance of the algorithm in approximating the true model as do random perturbations. The random perturbations MSD curve (yellow) is significantly higher than the graph homomorphic perturbations MSD curve (red) which is close to the non-perturbed MSD curve (blue). If we examine the average MSD of the individual models, we observe that the decay in performance is not as much as that for random perturbations. Moreover, since local graph-homomorphic perturbations do not affect the performance of the algorithms, we observe that the MSD curve follows that of the non-privatized algorithm.

### B. Classification in Distributed Learning

We next run an experiment on a classification dataset. We use the Avazu click through dataset [55], which contains a set of online add clicks. We distribute the data among  $P = 50$  agents. To get non-iid data, we add non-iid Gaussian noise to each agent's dataset. We let  $\mu = 0.5$ ,  $\rho = 0.001$ , and  $\sigma_g^2 = 0.8$ . We plot the testing error in Fig. 5 of the standard algorithm, the privatized algorithm with random perturbations, the privatized algorithm with graph-homomorphic perturbations, and the privatized algorithm with local graph-homomorphic perturbations. It comes as no surprise that using graph-homomorphic perturbations does not hinder the testing error as random perturbations do, and local graph-homomorphic perturbations do not change the testing error.

## V. CONCLUSION

The goal of this work has been to study the effect of privacy in distributed learning. We established the superiority of graph-homomorphic perturbations in the model performance, as opposed to random perturbations. We then designed local graph-homomorphic perturbations that ensure the added noise does

not affect the model performance. Thus, the main takeaway from this work is that graph-homomorphic perturbations are better than random perturbations in distributed learning.

## APPENDIX A

### SENSITIVITY OF THE DISTRIBUTED ALGORITHM

We study the sensitivity of the distributed learning algorithm (21)–(23), which is defined at each time instant by the expression:

$$\Delta(i) = \|\mathbf{w}_i - \mathbf{w}'_i\|. \quad (60)$$

This definition captures the change when the data samples of a single agent are changed. The prime symbol represents the new trajectory. We can bound the sensitivity using the triangle inequality by the individual errors and the difference in the optimal models:

$$\Delta(i) \leq \|\widetilde{\mathbf{w}}'_i\| + \|\widetilde{\mathbf{w}}_i\| + \sqrt{P}\|w^o - w'^o\|. \quad (61)$$

Then, for any constants  $B$  and  $B'$  chosen by the designer, we can use Markov's inequality to get the bounds:

$$\mathbb{P}(\|\widetilde{\mathbf{w}}_i\| \geq B) \leq \frac{\mathbb{E}\|\widetilde{\mathbf{w}}_i\|^2}{B^2}, \quad (62)$$

$$\mathbb{P}(\|\widetilde{\mathbf{w}}'_i\| \geq B') \leq \frac{\mathbb{E}\|\widetilde{\mathbf{w}}'_i\|^2}{B'^2}. \quad (63)$$

Now we recall from Theorem 1 that:

$$\mathbb{E}\|\widetilde{\mathbf{w}}_i\|^2 \leq \kappa_2^2 \mathbf{1}^\top \Gamma^i \begin{bmatrix} \mathbb{E}\|\widetilde{\mathbf{w}}_0\|^2 \\ \mathbb{E}\|\widetilde{\mathbf{w}}_0\|^2 \end{bmatrix} + O(\mu) + O(\mu^{-1}), \quad (64)$$

$$\mathbb{E}\|\widetilde{\mathbf{w}}'_i\|^2 \leq \kappa_2'^2 \mathbf{1}^\top \Gamma'^i \begin{bmatrix} \mathbb{E}\|\widetilde{\mathbf{w}}'_0\|^2 \\ \mathbb{E}\|\widetilde{\mathbf{w}}'_0\|^2 \end{bmatrix} + O(\mu) + O(\mu^{-1}). \quad (65)$$

It follows that the sensitivity is bounded by:

$$\Delta(i) \leq B + B' + \sqrt{P}\|w^o - w'^o\| \quad (66)$$

with high probability given by:

$$\begin{aligned} & \mathbb{P}(\Delta(i) \leq B + B' + \sqrt{P}\|w^o - w'^o\|) \\ & \geq \left( 1 - \frac{\kappa_2^2 \mathbf{1}^\top \Gamma^i \begin{bmatrix} \mathbb{E}\|\widetilde{\mathbf{w}}_0\|^2 \\ \mathbb{E}\|\widetilde{\mathbf{w}}_0\|^2 \end{bmatrix} + O(\mu) + O(\mu^{-1})}{B^2} \right) \\ & \quad \times \left( 1 - \frac{\kappa_2'^2 \mathbf{1}^\top \Gamma'^i \begin{bmatrix} \mathbb{E}\|\widetilde{\mathbf{w}}'_0\|^2 \\ \mathbb{E}\|\widetilde{\mathbf{w}}'_0\|^2 \end{bmatrix} + O(\mu) + O(\mu^{-1})}{B'^2} \right). \end{aligned} \quad (67)$$

## APPENDIX B

### PROOF OF THEOREM 1

The following proof follows similar steps to those used in [26] for non-private algorithms. Using the Jordan decomposition of  $A_2^\top A_0^\top A_1^\top$ :

$$A_2^\top A_0^\top A_1^\top = V_\theta J V_\theta^{-1}, \quad (68)$$

$$V_\theta \triangleq [q \quad V_R], \quad (69)$$

$$V_\theta^{-1} \triangleq \begin{bmatrix} \mathbf{1}^\top \\ V_L^\top \end{bmatrix}, \quad (70)$$

$$J \triangleq \begin{bmatrix} 1 & 0 \\ 0 & J_\theta \end{bmatrix}, \quad (71)$$

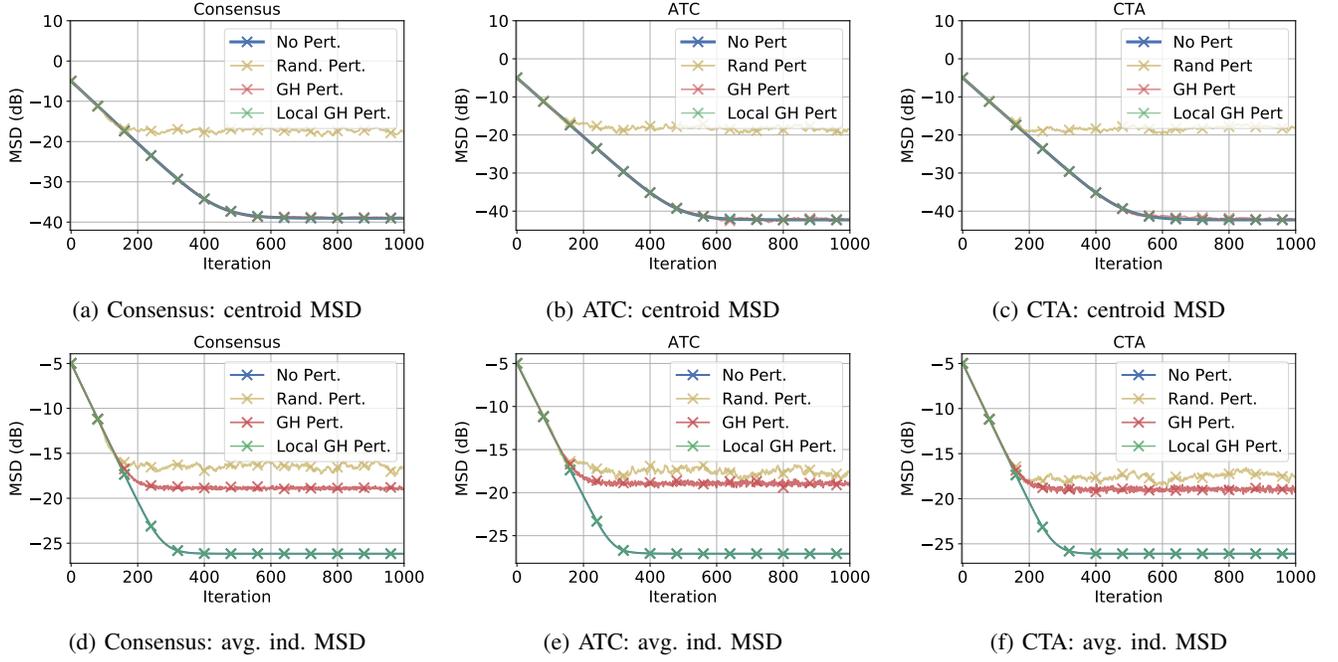


Fig. 4: MSD plots for the three distributed learning algorithms.

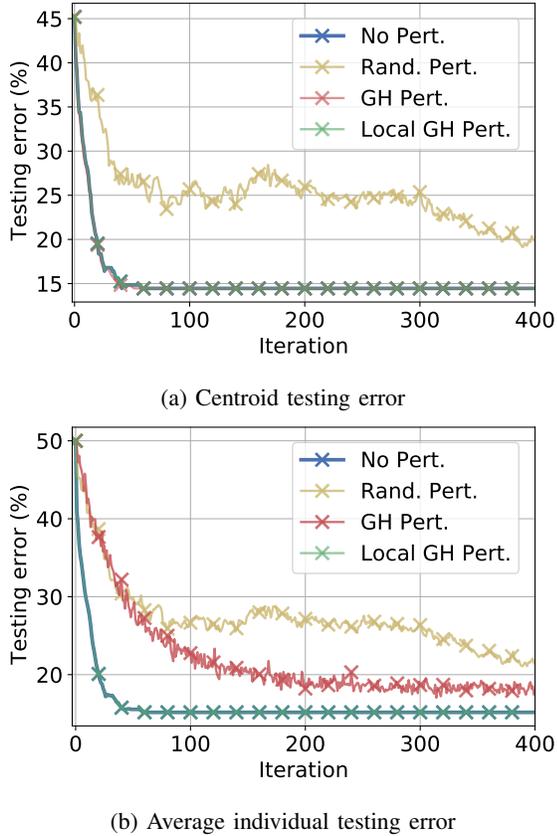


Fig. 5: Testing error of standard ATC, privatized ATC with random perturbations, and privatized ATC with graph-homomorphic perturbations.

where  $q$  is the Perron eigenvector of  $A_2^T A_0^T A_1^T$  and  $J_\theta$  contains Jordan blocks of the corresponding eigenvalues  $\lambda$  of the form

(example of a  $3 \times 3$  matrix):

$$\begin{bmatrix} \lambda & 0 & 0 \\ \theta & \lambda & 0 \\ 0 & \theta & \lambda \end{bmatrix}, \quad (72)$$

with a constant  $\theta$  in the subdiagonal. We first write:

$$\mathbf{B}_{i-1} = (\mathcal{V}_\theta^{-1})^T (\mathcal{J} - \mathcal{D}_{i-1}^T) \mathcal{V}_\theta^T, \quad (73)$$

where:

$$\mathcal{V}_\theta^{-1} \triangleq \mathcal{V}_\theta^{-1} \otimes I_M, \quad (74)$$

$$\mathcal{V}_\theta \triangleq \mathcal{V}_\theta \otimes I_M, \quad (75)$$

$$\mathcal{J} \triangleq J \otimes I_M = \begin{bmatrix} I_M & 0 \\ 0 & \mathcal{J}_\theta \end{bmatrix}, \quad (76)$$

$$\mathcal{D}_{i-1}^T \triangleq \mu \mathcal{V}_\theta^T \mathcal{A}_2^T \mathcal{H}_{i-1} \mathcal{A}_1^T (\mathcal{V}_\theta^{-1})^T = \begin{bmatrix} \mathcal{D}_{11,i-1}^T & \mathcal{D}_{21,i-1}^T \\ \mathcal{D}_{12,i-1}^T & \mathcal{D}_{22,i-1}^T \end{bmatrix}. \quad (77)$$

It is shown in the proof of Theorem 9.1 in [26] that:

$$\|I_M - \mathcal{D}_{11,i-1}\| \leq 1 - \sigma_{11}\mu, \quad (78)$$

$$\|\mathcal{D}_{ij}\| \leq \sigma_{ij}\mu, \quad (79)$$

for some positive constants  $\sigma_{ij}$  for  $i, j = 1, 2$ . Multiplying both sides of the error recursion (41) from the left by  $\mathcal{V}_\theta^T$ :

$$\begin{aligned} \mathcal{V}_\theta^T \tilde{\mathbf{w}}_i &= \mathcal{V}_\theta^T \mathbf{B}_{i-1} (\mathcal{V}_\theta^{-1})^T \mathcal{V}_\theta^T \tilde{\mathbf{w}}_{i-1} + \mu \mathcal{V}_\theta^T \mathcal{A}_2^T \mathbf{s}_i - \mu \mathcal{V}_\theta^T \mathcal{A}_2^T \mathbf{b} \\ &\quad + \mu \mathcal{V}_\theta^T \mathcal{A}_2^T \mathcal{H}_{i-1} \text{diag}(\mathcal{A}_1^T \mathcal{G}_{1,i}) - \mathcal{V}_\theta^T \text{diag}(\mathcal{A}_2^T \mathcal{G}_{2,i}) \\ &\quad - \mathcal{V}_\theta^T \mathcal{A}_2^T \text{diag}(\mathcal{A}_0^T \mathcal{G}_{0,i}) - \mathcal{V}_\theta^T \mathcal{A}_2^T \mathcal{A}_0^T \text{diag}(\mathcal{A}_1^T \mathcal{G}_{1,i}), \end{aligned} \quad (80)$$

and introducing the new notation:

$$\mathcal{V}_\theta^\top \tilde{\mathbf{w}}_i = \begin{bmatrix} (q^\top \otimes I_M) \tilde{\mathbf{w}}_i \\ (V_R^\top \otimes I) \tilde{\mathbf{w}}_i \end{bmatrix} \triangleq \begin{bmatrix} \tilde{\mathbf{w}}_i \\ \check{\mathbf{w}}_i \end{bmatrix}, \quad (81)$$

$$\mu \mathcal{V}_\theta^\top \mathcal{A}_2^\top \mathbf{s}_i = \mu \begin{bmatrix} (q^\top \otimes I_M) \mathcal{A}_2^\top \mathbf{s}_i \\ (V_R^\top \otimes I) \mathcal{A}_2^\top \mathbf{s}_i \end{bmatrix} \triangleq \begin{bmatrix} \tilde{\mathbf{s}}_i \\ \check{\mathbf{s}}_i \end{bmatrix}, \quad (82)$$

$$\mu \mathcal{V}_\theta^\top \mathcal{A}_2^\top \mathbf{b} = \mu \begin{bmatrix} (q^\top \otimes I_M) \mathcal{A}_2^\top \mathbf{b} \\ (V_R^\top \otimes I) \mathcal{A}_2^\top \mathbf{b} \end{bmatrix} \triangleq \begin{bmatrix} 0 \\ \check{\mathbf{b}} \end{bmatrix}, \quad (83)$$

we get:

$$\begin{bmatrix} \tilde{\mathbf{w}}_i \\ \check{\mathbf{w}}_i \end{bmatrix} = \begin{bmatrix} I_M - \mathbf{D}_{11,i-1}^\top & -\mathbf{D}_{21,i-1}^\top \\ -\mathbf{D}_{12,i-1}^\top & \mathcal{J}_\epsilon \end{bmatrix} \begin{bmatrix} \tilde{\mathbf{w}}_{i-1} \\ \check{\mathbf{w}}_{i-1} \end{bmatrix} + \begin{bmatrix} \tilde{\mathbf{s}}_i \\ \check{\mathbf{s}}_i \end{bmatrix} + \begin{bmatrix} 0 \\ \check{\mathbf{b}} \end{bmatrix} \\ + \mu \mathcal{V}_\theta^\top \mathcal{A}_2^\top \mathcal{H}_{i-1} \text{diag}(\mathcal{A}_1^\top \mathcal{G}_{1,i}) - \mathcal{V}_\theta^\top \text{diag}(\mathcal{A}_2^\top \mathcal{G}_{2,i}) \\ - \mathcal{V}_\theta^\top \mathcal{A}_2^\top \text{diag}(\mathcal{A}_0^\top \mathcal{G}_{0,i}) - \mathcal{V}_\theta^\top \mathcal{A}_2^\top \mathcal{A}_0^\top \text{diag}(\mathcal{A}_1^\top \mathcal{G}_{1,i}). \quad (84)$$

Then, taking the expectation of the  $\ell_2$ -norm, and using Jensen's inequality, we have:

$$\begin{aligned} \mathbb{E} \|\tilde{\mathbf{w}}_i\|^2 &\leq (1 - \sigma_{11}\mu) \mathbb{E} \|\tilde{\mathbf{w}}_{i-1}\|^2 + \frac{\sigma_{21}^2 \mu}{\sigma_{11}} \mathbb{E} \|\check{\mathbf{w}}_{i-1}\|^2 + \mathbb{E} \|\tilde{\mathbf{s}}_i\|^2 \\ &\quad + 2\mu^2 \mathbb{E} \|(q^\top \otimes I_M) \mathcal{A}_2^\top \mathcal{H}_{i-1} \text{diag}(\mathcal{A}_1^\top \mathcal{G}_{1,i})\|^2 + \\ &\quad + 2\mathbb{E} \|(q^\top \otimes I_M) \mathcal{A}_2^\top \mathcal{A}_0^\top \text{diag}(\mathcal{A}_1^\top \mathcal{G}_{1,i})\|^2 \\ &\quad + \mathbb{E} \|(q^\top \otimes I_M) \text{diag}(\mathcal{A}_2^\top \mathcal{G}_{2,i})\|^2 \\ &\quad + \mathbb{E} \|(q^\top \otimes I_M) \mathcal{A}_2^\top \text{diag}(\mathcal{A}_0^\top \mathcal{G}_{0,i})\|^2, \quad (85) \end{aligned}$$

and:

$$\begin{aligned} \mathbb{E} \|\check{\mathbf{w}}_i\|^2 &\leq \left( \rho(J_\theta) + \theta + \frac{2\sigma_{22}^2 \mu^2}{1 - \rho(J_\theta) - \theta} \right) \mathbb{E} \|\check{\mathbf{w}}_{i-1}\|^2 \\ &\quad + \frac{3\sigma_{21}^2 \mu^2}{1 - \rho(J_\theta) - \theta} \mathbb{E} \|\tilde{\mathbf{w}}_{i-1}\|^2 + \frac{3\|\check{\mathbf{b}}\|^2}{1 - \rho(J_\theta) - \theta} \\ &\quad + \mathbb{E} \|\check{\mathbf{s}}_i\|^2 \\ &\quad + 2\mu^2 \mathbb{E} \|(V_R^\top \otimes I_M) \mathcal{A}_2^\top \mathcal{H}_{i-1} \text{diag}(\mathcal{A}_1^\top \mathcal{G}_{1,i})\|^2 \\ &\quad + \mathbb{E} \|J_\epsilon^\top (V_R^\top \otimes I_M) \mathcal{A}_2^\top \mathcal{A}_0^\top \text{diag}(\mathcal{A}_1^\top \mathcal{G}_{1,i})\|^2 \\ &\quad + \mathbb{E} \|(V_R^\top \otimes I_M) \text{diag}(\mathcal{A}_2^\top \mathcal{G}_{2,i})\|^2 \\ &\quad + \mathbb{E} \|(V_R^\top \otimes I_M) \mathcal{A}_2^\top \text{diag}(\mathcal{A}_0^\top \mathcal{G}_{0,i})\|^2, \quad (86) \end{aligned}$$

with the cross terms equal to zero due to the independence of the zero-mean random variables. Then, we bound the sum of the gradient noise:

$$\begin{aligned} &\mathbb{E} \|\tilde{\mathbf{s}}_i\|^2 + \mathbb{E} \|\check{\mathbf{s}}_i\|^2 \\ &\leq \|\mathcal{V}_\theta\|^2 \mu^2 \sum_{p=1}^P \mathbb{E} \|\mathbf{s}_{p,i}\|^2 \\ &\leq \kappa_1^2 \mu^2 \sum_{p=1}^P \beta_{s,p}^2 \mathbb{E} \|\tilde{\phi}_{p,i-1}\|^2 + \sigma_{s,p}^2 \\ &\leq \kappa_1^2 \mu^2 \sum_{p=1}^P \beta_{s,p}^2 \sum_{m=1}^P (\mathbb{E} \|\tilde{\mathbf{w}}_{m,i-1}\|^2 + \mathbb{E} \|\mathbf{g}_{1,mp,i}\|^2) + \sigma_{s,p}^2 \\ &\leq \kappa_1^2 \mu^2 \beta_s^2 \mathbb{E} \|\tilde{\mathbf{w}}_{i-1}\|^2 + \kappa_1^2 \mu^2 \sigma_s^2 + \kappa_1^2 \mu^2 \sum_{p,m=1}^P \beta_{s,p}^2 \mathbb{E} \|\mathbf{g}_{1,mp,i}\|^2 \\ &\leq \kappa_1^2 \kappa_2^2 \mu^2 \beta_s^2 (\mathbb{E} \|\tilde{\mathbf{s}}_i\|^2 + \mathbb{E} \|\check{\mathbf{s}}_i\|^2) + \kappa_1^2 \mu^2 \sigma_s^2 \\ &\quad + \kappa_1^2 \mu^2 \sum_{p,m=1}^P \beta_{s,p}^2 \mathbb{E} \|\mathbf{g}_{1,mp,i}\|^2, \quad (87) \end{aligned}$$

where we introduced the constants  $\beta_s^2$  and  $\sigma_s^2$ , which are the sums of  $\beta_{s,p}^2$  and  $\sigma_{s,p}^2$ , respectively. Then, going back:

$$\begin{aligned} \mathbb{E} \|\tilde{\mathbf{w}}_i\|^2 &\leq (1 - \sigma_{11}\mu + \kappa_1^2 \kappa_2^2 \beta_s^2 \mu^2) \mathbb{E} \|\tilde{\mathbf{w}}_{i-1}\|^2 \\ &\quad + \left( \frac{\sigma_{21}^2 \mu}{\sigma_{11}} + \kappa_1^2 \kappa_2^2 \mu^2 \right) \mathbb{E} \|\check{\mathbf{w}}_{i-1}\|^2 + \kappa_1^2 \mu^2 \sigma_s^2 \\ &\quad + \kappa_1^2 \sum_{p,m=1}^P \beta_{s,p}^2 \mu^2 \mathbb{E} \|\mathbf{g}_{1,mp,i}\|^2 \\ &\quad + 2\mu^2 \mathbb{E} \|(q^\top \otimes I_M) \mathcal{A}_2^\top \mathcal{H}_{i-1} \text{diag}(\mathcal{A}_1^\top \mathcal{G}_{1,i})\|^2 \\ &\quad + 2\mathbb{E} \|(q^\top \otimes I_M) \mathcal{A}_2^\top \mathcal{A}_0^\top \text{diag}(\mathcal{A}_1^\top \mathcal{G}_{1,i})\|^2 \\ &\quad + \mathbb{E} \|(q^\top \otimes I_M) \text{diag}(\mathcal{A}_2^\top \mathcal{G}_{2,i})\|^2 \\ &\quad + \mathbb{E} \|(q^\top \otimes I_M) \mathcal{A}_2^\top \text{diag}(\mathcal{A}_0^\top \mathcal{G}_{0,i})\|^2, \quad (88) \end{aligned}$$

and:

$$\begin{aligned} \mathbb{E} \|\check{\mathbf{w}}_i\|^2 &\leq \left( \rho(J_\theta) + \theta + \frac{3\sigma_{22}^2 \mu^2}{1 - \rho(J_\theta) - \theta} + \kappa_1^2 \kappa_2^2 \beta_s^2 \mu^2 \right) \mathbb{E} \|\check{\mathbf{w}}_{i-1}\|^2 \\ &\quad + \left( \frac{3\sigma_{21}^2 \mu^2}{1 - \rho(J_\theta) - \theta} + \kappa_1^2 \kappa_2^2 \beta_s^2 \mu^2 \right) \mathbb{E} \|\tilde{\mathbf{w}}_{i-1}\|^2 \\ &\quad + \frac{3\|\check{\mathbf{b}}\|^2}{1 - \rho(J_\theta) - \theta} + \kappa_1^2 \mu^2 \sigma_s^2 \\ &\quad + \kappa_1^2 \mu^2 \sum_{p,m=1}^P \beta_{s,p}^2 \mathbb{E} \|\mathbf{g}_{1,mp,i}\|^2 \\ &\quad + 2\mu^2 \mathbb{E} \|(V_R^\top \otimes I_M) \mathcal{A}_2^\top \mathcal{H}_{i-1} \text{diag}(\mathcal{A}_1^\top \mathcal{G}_{1,i})\|^2 \\ &\quad + \mathbb{E} \|J_\theta^\top (V_R^\top \otimes I_M) \mathcal{A}_2^\top \mathcal{A}_0^\top \text{diag}(\mathcal{A}_1^\top \mathcal{G}_{1,i})\|^2 \\ &\quad + \mathbb{E} \|(V_R^\top \otimes I_M) \text{diag}(\mathcal{A}_2^\top \mathcal{G}_{2,i})\|^2 \\ &\quad + \mathbb{E} \|(V_R^\top \otimes I_M) \mathcal{A}_2^\top \text{diag}(\mathcal{A}_0^\top \mathcal{G}_{0,i})\|^2. \quad (89) \end{aligned}$$

Adding the two bounds:

$$\begin{aligned} \mathbb{E} \|\tilde{\mathbf{w}}_i\|^2 &\leq \kappa_2^2 \left( \bar{\gamma} \mathbb{E} \|\tilde{\mathbf{w}}_{i-1}\|^2 + \check{\gamma} \mathbb{E} \|\check{\mathbf{w}}_{i-1}\|^2 + \frac{3\|\check{\mathbf{b}}\|^2}{1 - \rho(J_\theta) - \theta} \right. \\ &\quad + 2\kappa_1^2 \mu^2 \sigma_s^2 + 2\kappa_1^2 \mu^2 \sum_{p,m=1}^P \beta_{s,p}^2 \mathbb{E} \|\mathbf{g}_{1,mp,i}\|^2 \\ &\quad + 2\mu^2 \mathbb{E} \|\mathcal{V}_\theta^\top \mathcal{A}_2^\top \mathcal{H}_{i-1} \text{diag}(\mathcal{A}_1^\top \mathcal{G}_{1,i})\|^2 \\ &\quad + \mathbb{E} \|\mathcal{V}_\theta^\top \mathcal{A}_2^\top \mathcal{A}_0^\top \text{diag}(\mathcal{A}_1^\top \mathcal{G}_{1,i})\|^2 \\ &\quad + \mathbb{E} \|\mathcal{V}_\theta^\top \text{diag}(\mathcal{A}_2^\top \mathcal{G}_{2,i})\|^2 \\ &\quad \left. + \mathbb{E} \|\mathcal{V}_\theta^\top \mathcal{A}_2^\top \text{diag}(\mathcal{A}_0^\top \mathcal{G}_{0,i})\|^2 \right) \\ &\leq \kappa_2^2 \left( \bar{\gamma} \mathbb{E} \|\tilde{\mathbf{w}}_{i-1}\|^2 + \check{\gamma} \mathbb{E} \|\check{\mathbf{w}}_{i-1}\|^2 + \frac{3\|\check{\mathbf{b}}\|^2}{1 - \rho(J_\theta) - \theta} \right. \\ &\quad + 2\kappa_1^2 \mu^2 \sigma_s^2 + 2\kappa_1^2 \mu^2 \sum_{p,m=1}^P \beta_{s,p}^2 \mathbb{E} \|\mathbf{g}_{1,mp,i}\|^2 \\ &\quad + \kappa_1^2 (1 + 2\delta^2 \mu^2) \mathbb{E} \|\text{diag}(\mathcal{A}_1^\top \mathcal{G}_{1,i})\|^2 \\ &\quad + \kappa_1^2 \mathbb{E} \|\text{diag}(\mathcal{A}_2^\top \mathcal{G}_{2,i})\|^2 \\ &\quad \left. + \kappa_1^2 \mathbb{E} \|\text{diag}(\mathcal{A}_0^\top \mathcal{G}_{0,i})\|^2 \right). \quad (90) \end{aligned}$$

Then, recursively bounding the MSE and taking the limit as  $i$  tends to infinity, we get:

$$\begin{aligned} & \limsup_{i \rightarrow \infty} \mathbb{E} \|\tilde{\mathbf{w}}_i\|^2 \\ & \leq \kappa_2^2 \mathbf{1}^\top (I - \Gamma)^{-1} \\ & \times \left[ \frac{\kappa_1^2 \mu^2 \sigma_s^2 + (4 + (2\delta^2 + P\beta_s^2)\mu^2) \sigma_g^2}{1 - \rho(J_\theta) - \theta} + \kappa_1^2 \mu^2 \sigma_s^2 + \kappa_1^2 (3 + (2\delta^2 + P\beta_s^2)\mu^2) \sigma_g^2 \right] \\ & = O(\mu) \sigma_s^2 + O(\mu) + (O(\mu^{-1}) + O(\mu)) \sigma_g^2, \end{aligned} \quad (91)$$

where:

$$\Gamma \triangleq \begin{bmatrix} \bar{\gamma} & \frac{\sigma_{21}^2 \mu}{\sigma_{11}} + \kappa_1^2 \kappa_2^2 \mu^2 \\ \frac{3\sigma_{12}^2 \mu^2}{1 - \rho(J_\theta) - \theta} + \kappa_1^2 \kappa_2^2 \beta_s^2 \mu^2 & \bar{\gamma} \end{bmatrix}. \quad (92)$$

### APPENDIX C PROOF OF LEMMA 1

We define  $\mathbf{w}_{c,i} \triangleq (q\mathbf{1}^\top \otimes I)\mathbf{w}_i$  and write:

$$\begin{aligned} \mathbf{w}_i - \mathbf{w}_{c,i} &= (I - q\mathbf{1}^\top \otimes I) \mathbf{w}_i \\ &= (V_L^\top \otimes I)(V_R \otimes I)\mathbf{w}_i \\ &= (V_L^\top \otimes I)J_\theta(V_R \otimes I)\mathbf{w}_{i-1} \\ &\quad - \mu(V_L^\top \otimes I)(V_R \otimes I)\text{col}\{\nabla_{w^\top} J_p(\phi_{p,i-1})\} \\ &\quad - \mu(V_L^\top \otimes I)(V_R \otimes I)\mathbf{s}_i \\ &\quad + (V_L^\top \otimes I)(V_R \otimes I)(\text{diag}(\mathcal{A}_2^\top \mathcal{G}_{2,i})) \\ &\quad + \mathcal{A}_2^\top \text{diag}(\mathcal{A}_0^\top \mathcal{G}_{0,i}) + \mathcal{A}_2^\top \mathcal{A}_0^\top \text{diag}(\mathcal{A}_1^\top \mathcal{G}_{1,i}). \end{aligned} \quad (93)$$

We bound  $\mathbb{E}\|(V_R \otimes I)\mathbf{w}_i\|^2$  by using Jensen's inequality with a constant  $\rho(J_\theta) < 1$  and define  $\kappa_1^2 = \|\mathcal{V}_\theta\|^2$  and  $\kappa_2^2 = \|\mathcal{V}_\theta^{-1}\|^2$ :

$$\begin{aligned} & \mathbb{E}\|(V_R \otimes I)\mathbf{w}_i\|^2 \\ & \leq \rho(J_\theta) \mathbb{E}\|(V_R \otimes I)\mathbf{w}_{i-1}\|^2 \\ & \quad + \frac{\kappa_1^2 \kappa_2^2 \mu^2}{1 - \rho(J_\theta)} \sum_{p=1}^P \mathbb{E}\|\mathbf{H}_{p,i-1} \tilde{\phi}_{p,i-1} + \nabla_{w^\top} J_p(w^o)\|^2 \\ & \quad + \kappa_1^2 \kappa_2^2 \mu^2 \sum_{p=1}^P \mathbb{E}\|\mathbf{s}_{p,i}\|^2 + v_1^2 v_2^2 (\mathbb{E}\|\text{diag}(\mathcal{A}_2^\top \mathcal{G}_{2,i})\|^2) \\ & \quad + \mathbb{E}\|\mathcal{A}_2^\top \text{diag}(\mathcal{A}_0^\top \mathcal{G}_{0,i})\|^2 + \mathbb{E}\|\mathcal{A}_2^\top \mathcal{A}_0^\top \text{diag}(\mathcal{A}_1^\top \mathcal{G}_{1,i})\|^2 \\ & \leq \rho(J_\theta) \mathbb{E}\|(V_R \otimes I)\mathbf{w}_{i-1}\|^2 + \frac{2\kappa_1^2 \kappa_2^2 \mu^2 \|b\|^2}{1 - \rho(J_\theta)} \\ & \quad + \kappa_1^2 \kappa_2^2 \mu^2 \sum_{p=1}^P \beta_{s,p}^2 \sum_{m \in \mathcal{N}_p} a_{1,mp} \mathbb{E}\|\tilde{\mathbf{w}}_{m,i-1}\|^2 \\ & \quad + \kappa_1^2 \kappa_2^2 \mu^2 \sum_{p,m=1}^P a_{1,mp}^2 \sigma_g^2 + \kappa_1^2 \kappa_2^2 \mu^2 \sigma_s^2 \\ & \quad + \kappa_1^2 \kappa_2^2 \sum_{p,m=1}^P (a_{2,mp}^2 + a_{0,mp}^2 + a_{1,mp}^2) \sigma_g^2. \end{aligned} \quad (94)$$

Then, the individual errors  $\mathbb{E}\|\tilde{\mathbf{w}}_{m,i-1}\|^2$  can be bounded as shown in Theorem 1:

$$\begin{aligned} \mathbb{E}\|\tilde{\mathbf{w}}_{m,i-1}\|^2 & \leq \mathbb{E}\|\tilde{\mathbf{w}}_{i-1}\|^2 \\ & \leq \kappa_2^2 \mathbf{1}^\top \left( \Gamma^{i-1} \begin{bmatrix} \mathbb{E}\|\tilde{\mathbf{w}}_0\|^2 \\ \mathbb{E}\|\tilde{\mathbf{w}}_0\|^2 \end{bmatrix} \right. \\ & \quad \left. + (I - \Gamma)^{-1} (I - \Gamma^{i-1}) \begin{bmatrix} O(\mu^2) + O(1) \\ O(\mu^2) + O(1) \end{bmatrix} \right), \end{aligned} \quad (95)$$

where the  $O(\mu^2)$  and  $O(1)$  terms are constants depending on the gradient noise variance, the bias term  $b$ , and the noise variance. Also, the matrix  $\Gamma$  captures the rate of the recursion and was previously defined in Appendix B.

Then, we plug back this bound into the main inequality (94) and recursively bound over  $i$ . The network disagreement is then bounded as:

$$\frac{1}{P} \sum_{p=1}^P \mathbb{E}\|\mathbf{w}_{p,i} - \mathbf{w}_{c,i}\|^2 \leq \frac{\kappa_2^2}{P} \mathbb{E}\|(V_R \otimes I)\mathbf{w}_i\|^2, \quad (96)$$

and in the limit:

$$\begin{aligned} & \limsup_{i \rightarrow \infty} \frac{1}{P} \sum_{p=1}^P \mathbb{E}\|\mathbf{w}_{p,i} - \mathbf{w}_{c,i}\|^2 \\ & \leq \frac{2\kappa_1^2 \kappa_2^4 \mu^2 \|b\|^2}{P(1 - \rho(J_\theta))^2} + \frac{\kappa_1^2 \kappa_2^4}{P(1 - \rho(J_\theta))} \sigma_g^2 + O(\mu) \\ & \quad + \frac{\kappa_1^2 \kappa_2^4}{P(1 - \rho(J_\theta))} (O(\mu^2) \sigma_s^2 + O(\mu^2) \sigma_g^2). \end{aligned} \quad (97)$$

### APPENDIX D PROOF OF THEOREM 2

Starting from (46) and taking the conditional mean of the squared Euclidean norm over the past models, we can split the norm into three independent terms: the model error, the gradient noise, and the added noise. Taking again expectations and using Jensen's with  $\alpha = \sqrt{1 - 2\nu\mu + \delta^2\mu^2}$ , we have:

$$\begin{aligned} & \mathbb{E}\|\tilde{\mathbf{w}}_{c,i}\|^2 \\ & \leq \alpha \mathbb{E}\|\tilde{\mathbf{w}}_{c,i-1}\|^2 + \mu^2 \mathbb{E}\|(q^\top \otimes I)\mathbf{s}_i\|^2 \\ & \quad + \frac{\mu^2}{1 - \alpha} \mathbb{E} \left\| \sum_{p=1}^P q_p \mathbf{H}_{p,i-1} \sum_{m \in \mathcal{N}_p} a_{1,mp} (\mathbf{w}_{m,i-1} - \mathbf{w}_{c,i-1}) \right\|^2 \\ & \quad + \mu^2 \mathbb{E} \left\| \sum_{p=1}^P q_p \mathbf{H}_{p,i-1} \sum_{m \in \mathcal{N}_p} a_{1,mp} \mathbf{g}_{1,mp,i} \right\|^2 \\ & \leq \alpha \mathbb{E}\|\tilde{\mathbf{w}}_{c,i-1}\|^2 + \mu^2 \mathbb{E}\|(q^\top \otimes I)\mathbf{s}_i\|^2 \\ & \quad + \frac{\delta^2 \mu^2}{1 - \alpha} \sum_{p=1}^P \mathbb{E}\|\mathbf{w}_{p,i-1} - \mathbf{w}_{c,i-1}\|^2 \\ & \quad + \mu^2 \mathbb{E} \left\| \sum_{p=1}^P q_p \mathbf{H}_{p,i-1} \sum_{m \in \mathcal{N}_p} a_{1,mp} \mathbf{g}_{1,mp,i} \right\|^2. \end{aligned} \quad (98)$$

We bound the gradient noise by starting from (34) and using Jensen's inequality to introduce  $\tilde{\mathbf{w}}_{c,i-1}$ :

$$\begin{aligned}
& \mathbb{E} \left\| (q^\top \otimes I) \mathbf{s}_i^2 \right\|^2 \\
&= \sum_{p=1}^P q_p^2 \mathbb{E} \left\| \mathbf{s}_{p,i} \right\|^2 \\
&\leq \sum_{p=1}^P q_p^2 \beta_{s,p}^2 \mathbb{E} \left\| \tilde{\boldsymbol{\phi}}_{p,i-1} \right\|^2 + q_p^2 \sigma_{s,p}^2 \\
&\leq \sum_{p=1}^P q_p^2 \beta_{s,p}^2 \sum_{m \in \mathcal{N}_p} a_{1,mp} \mathbb{E} \left\| \tilde{\mathbf{w}}_{m,i-1} \right\|^2 + a_{1,mp} \mathbb{E} \left\| \mathbf{g}_{1,mp,i} \right\|^2 \\
&\quad + \sigma_s^2 \\
&\leq 2\beta_s^2 \mathbb{E} \left\| \tilde{\mathbf{w}}_{c,i-1} \right\|^2 + \sigma_s^2 + \beta_s^2 \sigma_g^2 \\
&\quad + 2 \sum_{p=1}^P q_p^2 \beta_{s,p}^2 \sum_{m \in \mathcal{N}_p} a_{1,mp} \mathbb{E} \left\| \mathbf{w}_{m,i-1} - \mathbf{w}_{c,i-1} \right\|^2 \\
&\leq 2\beta_s^2 \mathbb{E} \left\| \tilde{\mathbf{w}}_{c,i-1} \right\|^2 + \sigma_s^2 + \beta_s^2 \sigma_g^2 \\
&\quad + 2\beta_s^2 \sum_{p=1}^P \mathbb{E} \left\| \mathbf{w}_{p,i-1} - \mathbf{w}_{c,i-1} \right\|^2. \tag{99}
\end{aligned}$$

The noise term can be bounded as follows by using twice Jensen's inequality:

$$\mathbb{E} \left\| \sum_{p=1}^P q_p \mathbf{H}_{p,i-1} \sum_{m \in \mathcal{N}_p} a_{1,mp} \mathbf{g}_{1,mp,i} \right\|^2 \leq \delta^2 \sigma_g^2. \tag{100}$$

We plug the bounds on the gradient noise and the added privacy noise in (98):

$$\begin{aligned}
\mathbb{E} \left\| \tilde{\mathbf{w}}_{c,i} \right\|^2 &\leq (\alpha + 2\beta_s^2 \mu^2) \mathbb{E} \left\| \tilde{\mathbf{w}}_{c,i-1} \right\|^2 + \mu^2 \sigma_s^2 + (\beta_s^2 + \delta^2) \mu^2 \sigma_g^2 \\
&\quad + \left( 2\beta_s^2 + \frac{\delta^2}{1-\alpha} \right) \mu^2 \sum_{p=1}^P \mathbb{E} \left\| \mathbf{w}_{p,i-1} - \mathbf{w}_{c,i-1} \right\|^2. \tag{101}
\end{aligned}$$

We use the bound from Lemma 1. Recursively bounding the second-order moment of the error and taking the limit:

$$\begin{aligned}
\limsup_{i \rightarrow \infty} \mathbb{E} \left\| \tilde{\mathbf{w}}_{c,i} \right\|^2 &\leq \frac{\mu^2 (\sigma_s^2 + (\beta_s^2 + \delta^2) \sigma_g^2)}{1 - \gamma_c} + O(\mu^2) \\
&\quad + \frac{\mu}{1 - \gamma_c} \left( 2\beta_s^2 + \frac{\delta^2}{1-\alpha} \right) \frac{\kappa_1^2 \kappa_2^4}{1 - \rho(J_\theta)} \sigma_g^2 \\
&= O(\mu) \sigma_s^2 + O(1) \sigma_g^2 + O(\mu^2). \tag{102}
\end{aligned}$$

#### APPENDIX E PROOF OF THEOREM 3

It suffices to show the noise generated from the local graph-homomorphic process is Laplacian since we already know that adding Laplacian noise makes the algorithm differentially private (see [36], [46]) with high probability. Thus, it is well known that the product of a uniform random variable  $U(0, 1)$  with a gamma random variable  $\Gamma(2, 1)$  results in an exponential random variable  $\text{Exp}(1)$  [56]. Then  $e^{-\mathbf{v}_\ell \mathbf{v}_m}$  is uniformly distributed on  $[0, 1]$ :

$$\mathbb{P}(e^{-\mathbf{v}_\ell \mathbf{v}_m} \leq c) = \mathbb{P}(\mathbf{v}_\ell \mathbf{v}_m \geq -\ln c) = e^{\ln c} = c. \tag{103}$$

But multiplying it by  $a$  makes the resulting variable uniformly distributed on  $[0, a]$ . The modulo  $p$  of a uniform random variable is uniform on  $[0, \pi]$  so long as  $a$  is a multiple of  $\pi$ . Let  $a = t\pi$  for some integer  $t$  and  $\mathbf{x} \sim U(0, a)$ . We divide the interval  $[0, a]$  into  $t$  disjoint sub-intervals of length  $\pi$ ,  $[0, a] = [0, \pi) \cup [\pi, 2\pi) \cdots \cup [(t-1)\pi, a]$ . On each of these sub-intervals  $[i\pi, (i+1)\pi)$ ,  $\mathbf{x}$  is uniformly distributed  $\mathbb{P}(\mathbf{x} \leq x | \mathbf{x} \in [i\pi, (i+1)\pi)) = x - i\pi$ , and so will  $\mathbf{x} \bmod \pi = \mathbf{x} - \lfloor \mathbf{x}/\pi \rfloor \pi = \mathbf{x} - i\pi$  on  $[0, \pi]$ . Thus since  $a = t\pi$  we get:

$$\begin{aligned}
\mathbb{P}(\mathbf{x} \leq x) &= \sum_{i=0}^{t-1} \mathbb{P}(\mathbf{x} \leq x | \mathbf{x} \in [i\pi, (i+1)\pi)) \\
&\quad \times \mathbb{P}(\mathbf{x} \in [i\pi, (i+1)\pi)) \\
&= \sum_{i=0}^{t-1} x \frac{\pi}{a} = x. \tag{104}
\end{aligned}$$

This now means that  $\mathbf{v}_{\ell m} \sim U(0, \pi)$ . Then, taking the difference of two exponential random variables results in a Laplacian. Thus, we require to transform two uniform random variables to two exponential random variables with parameter  $\frac{\sigma_g}{\sqrt{2}}$ . Taking  $-\frac{\sqrt{2}}{\sigma_g} \ln \mathbf{v}_{\ell m}$  results in an exponential random variable:

$$\mathbb{P} \left( -\frac{\sqrt{2}}{\sigma_g} \ln(\mathbf{v}_{\ell m}) \leq c \right) = \mathbb{P} \left( \mathbf{v}_{\ell m} \geq e^{-\frac{c\sigma_g}{\sqrt{2}}} \right) = 1 - e^{-\frac{\sigma_g c}{\sqrt{2}}}. \tag{105}$$

#### REFERENCES

- [1] D. Bertsekas, "A new class of incremental gradient methods for least squares problems," *SIAM J. Optim.*, vol. 7, pp. 913–926, Nov. 1997.
- [2] D. Blatt, A. O. Hero, and H. Gauchman, "A convergent incremental gradient method with a constant step size," *SIAM J. Optim.*, vol. 18, pp. 29–51, 2007.
- [3] F. S. Cattivelli and A. H. Sayed, "Analysis of spatial and incremental lms processing for distributed estimation," *IEEE Transactions on Signal Processing*, vol. 59, no. 4, pp. 1465–1480, 2011.
- [4] C. G. Lopes and A. H. Sayed, "Incremental adaptive strategies over distributed networks," *IEEE Transactions on Signal Processing*, vol. 55, no. 8, pp. 4064–4077, 2007.
- [5] B. Johansson, M. Rabi, and M. Johansson, "A randomized incremental subgradient method for distributed optimization in networked systems," *SIAM J. Optim.*, vol. 20, pp. 1157–1170, Jan 2009.
- [6] E. S. Helou and A. R. De Pierro, "Incremental subgradients for constrained convex optimization: A unified framework and new methods," *SIAM J. Optim.*, vol. 20, no. 3, p. 1547–1572, Dec 2009.
- [7] M. G. Rabbat and R. D. Nowak, "Quantized incremental algorithms for distributed optimization," *IEEE Journal on Selected Areas in Communications*, vol. 23, no. 4, pp. 798–808, 2005.
- [8] A. Nedic and D. P. Bertsekas, "Incremental subgradient methods for nondifferentiable optimization," *SIAM Journal on Optimization*, vol. 12, no. 1, pp. 109–138, 2001.
- [9] M. H. DeGroot, "Reaching a consensus," *Journal of the American Statistical Association*, vol. 69, no. 345, pp. 118–121, 1974.
- [10] 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.
- [11] L. Xiao and S. Boyd, "Fast linear iterations for distributed averaging," *Systems & Control Letters*, vol. 53, no. 1, pp. 65–78, Sep 2004.
- [12] S. Barbarossa and G. Scutari, "Bio-inspired sensor network design," *IEEE Signal Processing Magazine*, vol. 24, no. 3, pp. 26–35, May 2007.
- [13] B. Johansson, T. Keviczky, M. Johansson, and K. H. Johansson, "Subgradient methods and consensus algorithms for solving convex optimization problems," in *Proc. IEEE Conf. Dec. Control (CDC)*, Cancun, Mexico, December 2008, pp. 4185–4190.
- [14] W. Ren and R. Beard, "Consensus seeking in multiagent systems under dynamically changing interaction topologies," *IEEE Transactions on Automatic Control*, vol. 50, no. 5, pp. 655–661, 2005.

- [15] S. Boyd, A. Ghosh, B. Prabhakar, and D. Shah, "Randomized gossip algorithms," *IEEE Transactions on Information Theory*, vol. 52, no. 6, pp. 2508–2530, Jun 2006.
- [16] S. Kar and J. M. F. Moura, "Distributed consensus algorithms in sensor networks with imperfect communication: Link failures and channel noise," *IEEE Transactions on Signal Processing*, vol. 57, no. 1, pp. 355–369, Jan 2009.
- [17] K. Srivastava and A. Nedic, "Distributed asynchronous constrained stochastic optimization," *IEEE Journal of Selected Topics in Signal Processing*, vol. 5, no. 4, pp. 772–790, Aug 2011.
- [18] O. Hlinka, O. Slučiak, F. Hlawatsch, P. M. Djurić, and M. Rupp, "Likelihood consensus and its application to distributed particle filtering," *IEEE Transactions on Signal Processing*, vol. 60, no. 8, pp. 4334–4349, Aug 2012.
- [19] J. Chen and A. H. Sayed, "Diffusion adaptation strategies for distributed optimization and learning over networks," *IEEE Transactions on Signal Processing*, vol. 60, no. 8, pp. 4289–4305, Aug 2012.
- [20] S.-Y. Tu and A. H. Sayed, "Diffusion strategies outperform consensus strategies for distributed estimation over adaptive networks," *IEEE Transactions on Signal Processing*, vol. 60, no. 12, pp. 6217–6234, Dec 2012.
- [21] J. Chen and A. H. Sayed, "On the learning behavior of adaptive networks—part i: Transient analysis," *IEEE Transactions on Information Theory*, vol. 61, no. 6, pp. 3487–3517, Dec 2015.
- [22] S. Chouvardas, K. Slavakis, Y. Kopsinis, and S. Theodoridis, "A sparsity promoting adaptive algorithm for distributed learning," *IEEE Transactions on Signal Processing*, vol. 60, no. 10, pp. 5412–5425, 2012.
- [23] S. Vlaski and A. H. Sayed, "Distributed learning in non-convex environments—part I: Agreement at a linear rate," *IEEE Transactions on Signal Processing*, vol. 69, pp. 1242–1256, 2021.
- [24] 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, Jul 2008.
- [25] A. H. Sayed, "Adaptive networks," *Proceedings of the IEEE*, vol. 102, no. 4, pp. 460–497, 2014.
- [26] —, "Adaptation, learning, and optimization over networks," *Foundations and Trends in Machine Learning*, vol. 7, no. 4-5, pp. 311–801, 2014.
- [27] B. Hitaj, G. Ateniese, and F. Perez-Cruz, "Deep models under the gan: Information leakage from collaborative deep learning," in *Proceedings of ACM SIGSAC Conference on Computer and Communications Security*, New York, NY, USA, 2017, p. 603–618.
- [28] L. Melis, C. Song, E. De Cristofaro, and V. Shmatikov, "Exploiting unintended feature leakage in collaborative learning," in *IEEE Symposium on Security and Privacy*, San Francisco, CA, USA, May 2019, pp. 691–706.
- [29] M. Nasr, R. Shokri, and A. Houmansadr, "Comprehensive privacy analysis of deep learning: Passive and active white-box inference attacks against centralized and federated learning," in *IEEE Symposium on Security and Privacy*, San Jose, CA, May 2019, pp. 739–753.
- [30] L. Zhu and S. Han, "Deep leakage from gradients," in *Advances in Neural Information Processing Systems*, Vancouver, Canada, Dec 2019, pp. 17–31.
- [31] K. Bonawitz, V. Ivanov, B. Kreuter, A. Marcedone, H. B. McMahan, S. Patel, D. Ramage, A. Segal, and K. Seth, "Practical secure aggregation for privacy-preserving machine learning," in *Proc. of ACM SIGSAC Conference on Computer and Communications Security*, New York, USA, 2017, p. 1175–1191.
- [32] A. Gascón, P. Schoppmann, B. Balle, M. Raykova, J. Doerner, S. Zahur, and D. Evans, "Privacy-preserving distributed linear regression on high-dimensional data," *Proceedings on Privacy Enhancing Technologies*, vol. 2017, no. 4, pp. 345–364, 2017.
- [33] P. Mohassel and Y. Zhang, "Secureml: A system for scalable privacy-preserving machine learning," in *IEEE Symposium on Security and Privacy*, San Jose, CA, May 2017, pp. 19–38.
- [34] D. Froelicher, J. R. Troncoso-Pastoriza, A. Pyrgelis, S. Sav, J. S. Sousa, J.-P. Bossuat, and J.-P. Hubaux, "Scalable privacy-preserving distributed learning," *Proceedings on Privacy Enhancing Technologies*, vol. 2021, no. 2, pp. 323–347, 2021.
- [35] V. Nikolaenko, U. Weinsberg, S. Ioannidis, M. Joye, D. Boneh, and N. Taft, "Privacy-preserving ridge regression on hundreds of millions of records," in *IEEE Symposium on Security and Privacy*, Berkeley, CAA, May 2013, pp. 334–348.
- [36] C. Dwork and A. Roth, "The algorithmic foundations of differential privacy," *Found. Trends Theor. Comput. Sci.*, vol. 9, no. 3-4, pp. 211–407, 2014.
- [37] R. C. Geyer, T. Klein, and M. Nabi, "Differentially private federated learning: A client level perspective," *arXiv:1712.07557*, 2017.
- [38] R. Hu, Y. Guo, H. Li, Q. Pei, and Y. Gong, "Personalized federated learning with differential privacy," *IEEE Internet of Things Journal*, vol. 7, no. 10, pp. 9530–9539, 2020.
- [39] A. Triastcyn and B. Faltings, "Federated learning with bayesian differential privacy," in *IEEE International Conference on Big Data*, Los Angeles, CA, Dec 2019, pp. 2587–2596.
- [40] S. Truex, L. Liu, K.-H. Chow, M. E. Gursay, and W. Wei, "Ldp-fed: Federated learning with local differential privacy," in *Proceedings of the Third ACM International Workshop on Edge Systems, Analytics and Networking*, Heraklion, Greece, April 2020, pp. 61–66.
- [41] K. Wei, J. Li, M. Ding, C. Ma, H. H. Yang, F. Farokhi, S. Jin, T. Q. Quek, and H. V. Poor, "Federated learning with differential privacy: Algorithms and performance analysis," *IEEE Transactions on Information Forensics and Security*, vol. 15, pp. 3454–3469, 2020.
- [42] B. Jayaraman, L. Wang, D. Evans, and Q. Gu, "Distributed learning without distrust: Privacy-preserving empirical risk minimization," in *Advances in Neural Information Processing Systems*, vol. 31, Montreal, Canada, Dec 2018.
- [43] C. Li, P. Zhou, L. Xiong, Q. Wang, and T. Wang, "Differentially private distributed online learning," *IEEE Transactions on Knowledge and Data Engineering*, vol. 30, no. 8, pp. 1440–1453, 2018.
- [44] J. Zhu, C. Xu, J. Guan, and D. O. Wu, "Differentially private distributed online algorithms over time-varying directed networks," *IEEE Transactions on Signal and Information Processing over Networks*, vol. 4, no. 1, pp. 4–17, 2018.
- [45] M. A. Pathak, S. Rane, and B. Raj, "Multiparty differential privacy via aggregation of locally trained classifiers," in *Advances in Neural Information Processing Systems*, Vancouver, Canada, Dec 2010, pp. 1876–1884.
- [46] S. Vlaski and A. H. Sayed, "Graph-homomorphic perturbations for private decentralized learning," in *Proc. ICASSP*, Toronto, Canada, June 2021, pp. 5240–5244.
- [47] B. Ying and A. H. Sayed, "Performance limits of stochastic sub-gradient learning, part ii: Multi-agent case," *Signal Processing*, vol. 144, pp. 253–264, 2018.
- [48] S. Vlaski and A. H. Sayed, "Distributed learning in non-convex environments—part i: Agreement at a linear rate," *IEEE Transactions on Signal Processing*, vol. 69, pp. 1242–1256, 2021.
- [49] A. H. Sayed, *Inference and Learning from Data*. Cambridge University Press, 2023.
- [50] R. Nassif, S. Vlaski, C. Richard, J. Chen, and A. H. Sayed, "Multitask learning over graphs: An approach for distributed, streaming machine learning," *IEEE Signal Processing Magazine*, vol. 37, no. 3, pp. 14–25, 2020.
- [51] S. Wang, T. Tuor, T. Salonidis, K. K. Leung, C. Makaya, T. He, and K. Chan, "Adaptive federated learning in resource constrained edge computing systems," *IEEE Journal on Selected Areas in Communications*, vol. 37, no. 6, pp. 1205–1221, 2019.
- [52] X. Li, K. Huang, W. Yang, S. Wang, and Z. Zhang, "On the convergence of fedavg on non-iid data," in *International Conference on Learning Representations*, Addis Ababa, Ethiopia, 2020, pp. 1–26.
- [53] Ginar, "A review of random number generator (rng) on blockchain," Dec 2019. [Online]. Available: <https://medium.com/ginar-io/a-review-of-random-number-generator-rng-on-blockchain-fe342d76261b>
- [54] Y. Wang and H. V. Poor, "Decentralized stochastic optimization with inherent privacy protection," *IEEE Transactions on Automatic Control*, pp. 1–16, 2022.
- [55] Avazu and Kaggle, "Avazu's click-through rate prediction," 2014. [Online]. Available: <http://www.csie.ntu.edu.tw/~cj1in/libsvmtools/>
- [56] N. L. Johnson, S. Kotz, and N. Balakrishnan, *Continuous Univariate Distributions*. Wiley, 1995.