

## Part IV

---

### Advanced Topics



# 11 Privacy

---

One of the motivating factors for developing distributed algorithms for optimization and learning was the desire to avoid the exchange of raw data due to privacy concerns. Indeed, all distributed algorithms we encountered so far, ranging from federated to decentralized architectures, exchange intermediate estimates in lieu of raw data. One may then conjecture that all distributed algorithms fully preserve the privacy of participating agents, since no data is exchanged. As we will see in this chapter, this is more nuanced than it may seem at first side. In order to make inference about private information (such as local data), it is sufficient to observe *correlated* information (such as models).

We must hence develop a more rigorous framework for quantifying privacy and the leakage of information. We will use *differential privacy* as a privacy notion. This concept was developed by **Cynthia Dwork** and co-authors in a sequence of works in the 2000's, and has since gained significant traction in both academia and industry. Differential privacy was originally developed in the context of probing databases, rather than a learning context, and hence we will adapt concepts and notation to our setting.

## 11.1 LEAKAGE OF PRIVATE INFORMATION

---

Let us consider a simple example to illustrate how information can leak in a distributed setting where no raw data is exchanged. Recall the centralized stochastic gradient algorithm of Chapter 4, where a fusion center sends the current model  $\mathbf{w}_{i-1}$  to all agents. Each agent then performs a local update using the gradient approximation of its local cost function:

$$\psi_{k,i} = \mathbf{w}_{i-1} - \mu \widehat{\nabla} J_k(\mathbf{w}_{i-1}) \quad (11.1)$$

The locally updated models are sent back to the fusion center, where they are aggregated according to:

$$\mathbf{w}_i = \frac{1}{K} \sum_{k=1}^K \psi_{k,i} \quad (11.2)$$

Note that from (11.1), the parameter server can reconstruct the local gradient approximation via:

$$\widehat{\nabla J}_k(\mathbf{w}_{i-1}) = \frac{\mathbf{w}_{i-1} - \psi_{k,i}}{\mu} \quad (11.3)$$

Similarly, if we implement the gradient-exchange variant (4.37)–(4.16) of the algorithm, the server will immediately have access to the gradient approximation  $\widehat{\nabla J}_k(\mathbf{w}_{i-1})$ . This means that observing the evolution of models over time allows a parameter server (or someone intercepting the message) to reconstruct the locally constructed gradient approximations. These in turn are generally correlated with the data itself. For the mean squared error, for example, where:

$$J_k(w) = \mathbb{E}Q(w; \mathbf{x}_k) = \mathbb{E} \frac{1}{2} \left( \gamma_k - \mathbf{h}_k^\top w \right)^2 \quad (11.4)$$

we have:

$$\widehat{\nabla J}_k^{\text{ord}}(w) = \nabla Q(w; \mathbf{x}_k) = -\mathbf{h}_k \left( \gamma_k - \mathbf{h}_k^\top w \right) \quad (11.5)$$

For the logistic regression problem on the other hand, where:

$$J_k(w) = \mathbb{E}Q(w; \mathbf{x}_k) = \mathbb{E} \ln \left( 1 + e^{-\gamma_k \mathbf{h}_k^\top w} \right) + \frac{\rho}{2} \|w\|^2 \quad (11.6)$$

we have:

$$\widehat{\nabla J}_k^{\text{ord}}(w) = \nabla Q(w; \mathbf{x}_k) = \rho w - \frac{\mathbf{h}_k}{1 + e^{\gamma_k \mathbf{h}_k^\top w}} \quad (11.7)$$

We observe that in both cases, an observer with access to the model  $w$  and gradient approximation  $\widehat{\nabla J}_k^{\text{ord}}(w)$  is able to reconstruct the feature vector  $\mathbf{h}_k$  up to a scaling factor  $\gamma_k - \mathbf{h}_k^\top w$  in the case of the mean squared error cost, and a factor of  $1/(1 + e^{\gamma_k \mathbf{h}_k^\top w})$  in the case of logistic regression.

These observations imply that simply sharing intermediate estimates does not guarantee privacy by default. More care needs to be taken in quantifying *exactly* how much privacy is lost by participating in a distributed learning scheme. Are certain architectures more private than others? Is there something we can do if we are unhappy with the level of privacy inherent in a particular algorithm? As we will see, the framework of differential privacy will allow us to address both of these questions.

## 11.2 DIFFERENTIAL PRIVACY

Let us consider a collection of  $K$  agents indexed by  $k$ . Each agent has access to private data, which we model through the random variable  $\mathbf{x}_k$ . Any agent will have the option of participating in collaborative effort, which we describe generically as  $\text{Alg}(\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_K)$ . Consider also an alternative scenario, where

an arbitrary agent, let us say without loss of generality, agent 1, refuses to participate in the learning protocol, and is replaced by some other agent 1' with different local data  $\mathbf{x}_{1'} \in \mathcal{X}$  from a set of permissible local data sets  $\mathcal{X}$ . We can then state the formal definition of  $\epsilon$ -differential privacy.

**DEFINITION 11.1 ( $\epsilon$ -differential privacy).** We say that an algorithm  $\mathcal{A}(\cdot)$  is  $\epsilon$ -differentially private, if it holds that:

$$\frac{f(\text{Alg}(\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_K))}{f(\text{Alg}(\mathbf{x}_{1'}, \mathbf{x}_2, \dots, \mathbf{x}_K))} \leq e^\epsilon \quad (11.8)$$

for all  $\mathbf{x}_{1'} \in \mathcal{X}$ , where  $f(\cdot)$  denotes the probability density function of the argument.  $\square$

For small  $\epsilon$  it holds that  $e^\epsilon \approx 1 + \epsilon \approx 1$ , and hence the condition of  $\epsilon$ -differential privacy essentially ensures that the distribution of the output of the algorithm is essentially unaffected by the presence of agent 1. We conclude from this definition, that if  $\epsilon$  is small, agent 1 does not relinquish much privacy by participating in the learning protocol, since the output of the algorithm is largely unaffected by its participation. In other words, agent 1 by participating does not reveal particularly unique information that would allow an observer of the output  $f(\text{Alg}(\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_K))$  to identify the presence of agent 1. A useful corollary of this fact is that if the local data sets happen to be perfectly homogeneous, we can expect little privacy loss.

---

**Example 11.1 (Homogeneous agents)** Suppose the local data distributions are i.i.d, meaning  $\mathbf{x}_1 \sim \mathbf{x}_2 \sim \dots \sim \mathbf{x}_K \sim \mathbf{x}_{1'}$ . It then follows that:

$$f(\text{Alg}(\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_K)) = f(\text{Alg}(\mathbf{x}_{1'}, \mathbf{x}_2, \dots, \mathbf{x}_K)) \quad (11.9)$$

and hence the procedure is 0-differentially private. This essentially formalized the fact that if all information provided by agents is common knowledge, there is no privacy loss incurred by any given agent's participation.

**Example 11.2 (Randomized response)** Let us consider a scenario where  $\mathbf{x}_k$  contains potentially embarrassing information, such as the answer to the question "Have you committed a crime?". For agents who have committed a crime, we will then have  $\mathbf{x}_k = 1$  with probability one, and for those who have not we will have  $\mathbf{x}_k = -1$  with probability one. Since not all agents have committed a crime, and hence the  $\mathbf{x}_k$  is not identical, agents may be hesitant to reveal this private information. We may then envision the following randomized response mechanism and algorithm for estimating the rate of criminals among the collection of agents:

- Each agent  $k$  flips a fair coin.

- If the coin falls on heads, the agent responds truthfully, i.e., returns  $\mathbf{m}_k = \mathbf{x}_k$ .
- If the coin falls on tails, the agent flips another coin and returns  $\mathbf{m}_k = 1$  on heads, and  $\mathbf{m}_k = -1$  on tails.
- The server aggregates the messages and estimates the rate of criminals as from  $\{\mathbf{m}_k\}_{k=1}^K$ .

We can show that the above mechanism is differentially private. The source of privacy protection here arises from the plausible deniability, which in turn is the result of the randomized response mechanism. If an agent answers  $\mathbf{m}_k = 1$ , this may be because it has committed a crime, but it also be merely as a result of flipping two heads in a row, which occurs with probability  $\frac{1}{4}$ . More formally, we have:

$$\text{Alg}(\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_K) = \{\mathbf{m}_k\}_{k=1}^K \quad (11.10)$$

and

$$\begin{aligned} \frac{f(\text{Alg}(\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_K))}{f(\text{Alg}(\mathbf{x}_{1'}, \mathbf{x}_2, \dots, \mathbf{x}_K))} &= \frac{f(\mathbf{m}_1, \mathbf{m}_2, \dots, \mathbf{m}_K)}{f(\mathbf{m}_1, \mathbf{m}_2, \dots, \mathbf{m}_K)} \\ &\stackrel{(a)}{=} \frac{f(\mathbf{m}_1) f(\mathbf{m}_2) \cdots f(\mathbf{m}_K)}{f(\mathbf{m}_{1'}) f(\mathbf{m}_2) \cdots f(\mathbf{m}_K)} \\ &= \frac{f(\mathbf{m}_1)}{f(\mathbf{m}_{1'})} \end{aligned} \quad (11.11)$$

where (a) follows since the messages are generated independently at each agent. Then, if  $\mathbf{x}_1 = \mathbf{x}_{1'}$ , it is immediate that  $\frac{f(\mathbf{m}_1)}{f(\mathbf{m}_{1'})} = 1$ . For the other cases, we compute:

$$\Pr\{\mathbf{m}_1 = 1 | \mathbf{x}_1 = 1\} = \Pr\{\mathbf{m}_{1'} = 1 | \mathbf{x}_{1'} = 1\} = \frac{3}{4} \quad (11.12)$$

$$\Pr\{\mathbf{m}_1 = 1 | \mathbf{x}_1 = 0\} = \Pr\{\mathbf{m}_{1'} = 1 | \mathbf{x}_{1'} = 0\} = \frac{1}{4} \quad (11.13)$$

$$\Pr\{\mathbf{m}_1 = 0 | \mathbf{x}_1 = 1\} = \Pr\{\mathbf{m}_{1'} = 0 | \mathbf{x}_{1'} = 1\} = \frac{1}{4} \quad (11.14)$$

$$\Pr\{\mathbf{m}_1 = 0 | \mathbf{x}_1 = 0\} = \Pr\{\mathbf{m}_{1'} = 0 | \mathbf{x}_{1'} = 0\} = \frac{3}{4} \quad (11.15)$$

It follows that  $\frac{1}{3} \leq \frac{f(\mathbf{m}_1)}{f(\mathbf{m}_{1'})} \leq 3$ , and hence:

$$\frac{f(\text{Alg}(\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_K))}{f(\text{Alg}(\mathbf{x}_{1'}, \mathbf{x}_2, \dots, \mathbf{x}_K))} \leq 3 = e^{\ln 3} \quad (11.16)$$

We conclude that the mechanism is  $\ln 3$ -differentially private.  $\square$

### 11.2.1 Laplace Mechanism

The randomized response mechanism we encountered in the previous section was useful, as it provides indication that adding randomness to messages can

result in quantifiable privacy gain (in the sense of differential privacy). There are two remaining limitations that need to be addressed before we can apply this concept to learning algorithms. First, we would like to be able to control the level of privacy  $\epsilon$ , rather than simply accepting  $\ln 3$ -differential privacy. Second, the randomized response mechanism is natural when dealing with binary (yes or no) questions, but while learning, we encounter continuous, vector-valued quantities. To this end, we introduce the Laplace mechanism.

**DEFINITION 11.2 ( $\ell_1$ -sensitivity).** The sensitivity of  $\text{Alg}(\cdot)$  is defined as:

$$\Delta = \max_{\mathbf{x}_{1'} \in \mathcal{X}} \|\text{Alg}(\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_K) - \text{Alg}(\mathbf{x}_{1'}, \mathbf{x}_2, \dots, \mathbf{x}_K)\|_1 \quad (11.17)$$

□

**LEMMA 11.1 (Laplace mechanism).** Suppose  $\text{Alg}(\cdot)$  has  $\ell_1$ -sensitivity  $\Delta$ , and define:

$$\text{LAlg}(\cdot) = \text{Alg}(\cdot) + \mathbf{v}_p \quad (11.18)$$

where  $\mathbf{v}_p$  is a vector of suitable dimension  $M_v$ , where each entry follows the Laplace distribution:

$$f_{\mathbf{v}_p}(v) = \frac{1}{(2b_v)^{M_v}} e^{-\frac{\|\mathbf{v}\|_1}{b_v}} \quad (11.19)$$

Then,  $\text{LAlg}(\cdot)$  is  $\left(\frac{\Delta}{b_v}\right)$ -differentially private.

**Proof:** We have:

$$\begin{aligned} \frac{f(\text{LAlg}(\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_K))}{f(\text{LAlg}(\mathbf{x}_{1'}, \mathbf{x}_2, \dots, \mathbf{x}_K))} &= \frac{e^{-\frac{\|\text{LAlg}(\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_K) - \mathbf{x}\|}{b_v}}}{e^{-\frac{\|\text{LAlg}(\mathbf{x}_{1'}, \mathbf{x}_2, \dots, \mathbf{x}_K) - \mathbf{x}\|}{b_v}}} \\ &= e^{\frac{\|\text{LAlg}(\mathbf{x}_{1'}, \mathbf{x}_2, \dots, \mathbf{x}_K) - \mathbf{x}\| - \|\text{LAlg}(\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_K) - \mathbf{x}\|}{b_v}} \\ &\stackrel{(a)}{\leq} e^{\frac{\|\text{LAlg}(\mathbf{x}_{1'}, \mathbf{x}_2, \dots, \mathbf{x}_K) - \text{LAlg}(\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_K)\|}{b_v}} \\ &\stackrel{(11.17)}{\leq} e^{\frac{\Delta}{b_v}} \end{aligned} \quad (11.20)$$

□

It follows that each base-algorithm with bounded sensitivity  $\Delta$  can be *privatized* via the addition of Laplacian noise to the output. The level of noise necessary is related to the sensitivity  $\Delta$  of the algorithm and the desired level of privacy  $\epsilon$ .

### 11.2.2 Repeated Probing

Most algorithms for distributed optimization and learning we have encountered so far have been iterative, meaning that we repeatedly sample and compute functions of the local data, and repeatedly exchange model estimates. We hence

need to quantify as well the effect of these repeated interactions on privacy. It turns out that differential privacy is quite amenable to such settings.

**LEMMA 11.2 (Multiple algorithm evaluations).** *Suppose  $\text{Alg}_1(\cdot)$  is  $\epsilon_1$ -differentially private,  $\text{Alg}_2(\cdot)$  is  $\epsilon_2$ -differentially private, and the evaluations are independent. It then follows that:*

$$\text{Alg}(\cdot) = \text{col} \{ \text{Alg}_1(\cdot), \text{Alg}_2(\cdot) \} \quad (11.21)$$

*is  $(\epsilon_1 + \epsilon_2)$ -differentially private.*

**Proof:** We have:

$$\begin{aligned} \frac{f(\text{Alg}(\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_K))}{f(\text{Alg}(\mathbf{x}_{1'}, \mathbf{x}_2, \dots, \mathbf{x}_K))} &= \frac{f(\text{Alg}_1(\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_K), \text{Alg}_2(\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_K))}{f(\text{Alg}_1(\mathbf{x}_{1'}, \mathbf{x}_2, \dots, \mathbf{x}_K), \text{Alg}_2(\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_K))} \\ &= \frac{f(\text{Alg}_1(\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_K)) \cdot f(\text{Alg}_2(\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_K))}{f(\text{Alg}_1(\mathbf{x}_{1'}, \mathbf{x}_2, \dots, \mathbf{x}_K)) \cdot f(\text{Alg}_2(\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_K))} \\ &\leq e^{\epsilon_1} \cdot e^{\epsilon_2} = e^{\epsilon_1 + \epsilon_2} \end{aligned} \quad (11.22)$$

□

### 11.3 PRIVATE DISTRIBUTED LEARNING

For simplicity, we will illustrate private distributed learning for the diffusion algorithm (8.62)–(8.63) from Chapter 8, though analogous privacy considerations can be applied to other algorithms. Let us repeat the diffusion algorithm for reference:

$$\psi_{k,i} = \mathbf{w}_{k,i-1} - \mu \widehat{\nabla J}_k(\mathbf{w}_{k,i-1}) \quad (11.23)$$

$$\mathbf{w}_{k,i} = \sum_{\ell \in \mathcal{N}_k} a_{\ell k} \psi_{\ell,i} \quad (11.24)$$

As described when motivating the privacy developments in this chapter, agents may be concerned to share their models  $\psi_{k,i}$ , as an eavesdropper or malicious agent may leverage these to reconstruct the gradient approximation  $\widehat{\nabla J}_k(\mathbf{w}_{k,i-1})$ , which will in turn reveal information about the locally available data  $\mathbf{x}_k$  at agent  $k$ . Inspired by the Laplace mechanism of Section 11.2.1, we can then consider the following variant:

$$\phi_{k,i} = \mathbf{w}_{k,i-1} - \mu \widehat{\nabla J}_k(\mathbf{w}_{k,i-1}) \quad (11.25)$$

$$\psi_{k,i} = \phi_{k,i} + \mathbf{v}_{k,i} \quad (11.26)$$

$$\mathbf{w}_{k,i} = \sum_{\ell \in \mathcal{N}_k} a_{\ell k} \psi_{\ell,i} \quad (11.27)$$

where  $\mathbf{v}_{k,i}$  follows a Laplace distribution as Lemma 11.1. We will present a theorem further below which establishes a privacy guarantee for this kind of construction. Before we proceed, however, note that if we combine (11.25) and (11.26),



we obtain after some reformulations:

$$\psi_{k,i} = \mathbf{w}_{k,i-1} - \mu \left( \widehat{\nabla J}_k(\mathbf{w}_{k,i-1}) - \frac{\mathbf{v}_{k,i}}{\mu} \right) \quad (11.28)$$

$$\mathbf{w}_{k,i} = \sum_{\ell \in \mathcal{N}_k} a_{\ell k} \psi_{\ell,i} \quad (11.29)$$

We can then define the privatized gradient approximation:

$$\widehat{\nabla J}_k^{\text{priv}}(\mathbf{w}_{k,i-1}) \triangleq \widehat{\nabla J}_k(\mathbf{w}_{k,i-1}) - \frac{\mathbf{v}_{k,i}}{\mu} \quad (11.30)$$

This adjusted gradient approximation satisfies all conditions of Chapter (3). However, in this construction, the gradient approximation  $\widehat{\nabla J}_k(\mathbf{w}_{k,i-1})$  is perturbed by the privacy noise  $\mathbf{v}_{k,i}$  divided by the step-size  $\mu$ . For small step-sizes this can result in serious amplification of the privacy noise component, yielding:

$$\sigma_{k,\text{priv}}^2 = \sigma_k^2 + \frac{\sigma_v^2}{\mu^2} \quad (11.31)$$

where  $\sigma_k^2$  denotes the absolute component of the gradient approximation  $\widehat{\nabla J}_k(\mathbf{w}_{k,i-1})$  and  $\sigma_v^2$  denotes the variance of the Laplacian privacy noise  $\mathbf{v}_{k,i}$ . We can then conclude from Theorem 9.1 that the limiting performance of the privatized diffusion algorithm (for small step-sizes) will be given by:

$$\lim_{i \rightarrow \infty} \mathbb{E} \|\tilde{\mathbf{w}}_{k,i}\|^2 = O(\mu \sigma^2) + O(\mu^{-1} \sigma_v^2) \quad (11.32)$$

We note that the privacy preserving noise added to the distributed recursion as a very detrimental effect on the limiting performance of the algorithm. This begs the question of whether we can do any better. After all, the fact that models are computed and exchanged in a decentralized manner should give us some freedom to obtain improved privacy.

### 11.3.1 Graph-Homomorphic Perturbations

The idea here will be to tune perturbations to the network topology in order to minimize their impact on learning performance while preserving privacy. The first step is to allow

$$\phi_{k,i} = \mathbf{w}_{k,i-1} - \mu \widehat{\nabla J}_k(\mathbf{w}_{k,i-1}) \quad (11.33)$$

$$\psi_{k\ell,i} = \phi_{k,i-1} + \mathbf{v}_{k\ell,i} \quad (11.34)$$

$$\mathbf{w}_{k,i} = \sum_{\ell \in \mathcal{N}_k} a_{\ell k} \psi_{\ell k,i} \quad (11.35)$$

Note that agent  $k$  in this construction sends different models  $\psi_{k\ell,i}$  to different neighbors, each perturbed by a different privacy noise  $\mathbf{v}_{k\ell,i}$ . For the network

centroid, we have:

$$\begin{aligned}
\mathbf{w}_{c,i} &\triangleq \frac{1}{K} \sum_{k=1}^K \mathbf{w}_{k,i} \\
&\stackrel{(11.35)}{=} \frac{1}{K} \sum_{k=1}^K \sum_{\ell=1}^K a_{\ell k} \phi_{\ell,i} + \frac{1}{K} \sum_{k=1}^K \sum_{\ell=1}^K a_{\ell k} \mathbf{v}_{\ell k,i} \\
&= \frac{1}{K} \sum_{\ell=1}^K \left( \sum_{k=1}^K a_{\ell k} \right) \phi_{\ell,i} + \frac{1}{K} \sum_{\ell=1}^K \sum_{k=1}^K a_{\ell k} \mathbf{v}_{\ell k,i} \\
&\stackrel{(a)}{=} \frac{1}{K} \sum_{\ell=1}^K \phi_{\ell,i} + \frac{1}{K} \sum_{\ell=1}^K \sum_{k=1}^K a_{\ell k} \mathbf{v}_{\ell k,i} \\
&\stackrel{(11.33)}{=} \mathbf{w}_{c,i-1} - \frac{\mu}{K} \sum_{\ell=1}^K \widehat{\nabla J}_k(\mathbf{w}_{k,i-1}) + \frac{1}{K} \sum_{\ell=1}^K \sum_{k=1}^K a_{\ell k} \mathbf{v}_{\ell k,i} \tag{11.36}
\end{aligned}$$

where (a) follows for symmetric adjacency matrices. We saw in Chapter 9 that the network centroid, in many cases describes the dominant dynamics of the network of learners. It is then reasonable to construct perturbations  $\mathbf{v}_{\ell k,i}$  so their impact on the centroid  $\mathbf{w}_{c,i}$  is minimized. This motivates the following definition:

**DEFINITION 11.3 (Graph-Homomorphic Perturbations).** A set of perturbations  $\mathbf{v}_{\ell k,i}$  is homomorphic for the the graph defined by the adjacency matrix  $A \triangleq [a_{\ell k}]$  if it holds with probability one that:

$$\frac{1}{K} \sum_{\ell=1}^K \sum_{k=1}^K a_{\ell k} \mathbf{v}_{\ell k,i} = 0 \tag{11.37}$$

While other constructions are possible, we present here a simple construction, which can be implemented locally and independently at every agent  $k$ .

**LEMMA 11.3 (Constructing Graph-Homomorphic Perturbations).** *Let each agent  $\ell$  sample independently from the Laplace distribution  $\mathbf{v}'_{\ell,i} \sim \text{Lap}(0, b_v)$  with variance  $\sigma_v^2 = 2b_v^2$ . Then, the construction:*

$$\mathbf{v}_{\ell k,i} = \begin{cases} \mathbf{v}'_{\ell,i}, & \text{if } k \in \mathcal{N}_\ell \text{ and } k \neq \ell, \\ -\frac{1-a_{\ell\ell}}{a_{\ell\ell}} \mathbf{v}'_{\ell,i}, & \text{if } k = \ell. \end{cases} \tag{11.38}$$

*is homomorphic for the graph described by the symmetric adjacency matrix  $A = A^\top$ .*

**Proof:** The result can be verified immediately by substitution.  $\square$

It is important to note that graph-homomorphic perturbations are not completely unimpactful to the recursion, since the network centroid recursion (11.36) is driven by gradient approximations evaluated at  $\mathbf{w}_{k,i-1}$ , rather than  $\mathbf{w}_{c,i-1}$ .

As a result, the privacy perturbation will affect the learning dynamics through the network disagreement, rather than the centroid directly. This will result in reduced, albeit non-zero effect of the privacy perturbations. Adjusting the arguments of Chapter 9 to allow for privacy perturbation of this manner is fairly straightforward, resulting in:

$$\lim_{i \rightarrow \infty} \mathbb{E} \|\tilde{\mathbf{w}}_{k,i}\|^2 = O(\mu\sigma^2) + O(\sigma_v^2) \quad (11.39)$$

Comparing the two performance expressions (11.32) and (11.39) we note that the effect of the privacy perturbations in is reduced by a factor of  $\mu$ , and now constant for small step-sizes. This results in improved performance for a given level of privacy.

### 11.3.2 Differential Privacy Guarantee for Privatized Diffusion

We motivated the privatized diffusion implementations through the Laplacian mechanism, but have not yet established an actual privacy guarantee. We will provide a sketch of the proof under the simplifying assumption that gradient approximations are bounded almost surely:

$$\|\widehat{\nabla J}_k(\mathbf{w}_{k,i-1})\| \leq G \quad (11.40)$$

This condition is in fact not necessary, and can be removed at the expense of more cumbersome privacy analysis. We refer the reader to [Rizk, Vlaski, Sayed 2023] for this.

Let us now consider the setting, where agent 1 has decided not to volunteer its private information for the diffusion of information, and its data  $\mathbf{x}_1$  is replaced by some other data  $\mathbf{x}_{1'}$ , following a different distribution. In this setting, implementing the perturbed diffusion recursions, would naturally result in a different learning trajectory  $\mathbf{w}'_{k,i}$  at every agent  $k$ , since the data  $\mathbf{x}_{1'}$  propagates through the gradient approximations and the diffusion of estimates through the entire network. We can then define and bound the sensitivity of the diffusion algorithms as.

**LEMMA 11.4 (Sensitivity of the diffusion algorithm).** *The distance between the trajectories  $\mathbf{w}_{k,i}$  and  $\mathbf{w}'_{k,i}$  is bounded with probability one by:*

$$\Delta(i) \triangleq \max_k \|\mathbf{w}_{k,i} - \mathbf{w}'_{k,i}\| \leq \mu 2Gi \quad (11.41)$$

**Proof:** For reference, we repeat:

$$\mathbf{w}_{k,i} = \sum_{\ell \in \mathcal{N}_k} a_{\ell k} \left( \mathbf{w}_{\ell,i-1} - \mu \widehat{\nabla J}_\ell(\mathbf{w}_{\ell,i-1}) \right) + \sum_{\ell \in \mathcal{N}_k} a_{\ell k} \mathbf{v}_{\ell k,i} \quad (11.42)$$

$$\mathbf{w}'_{k,i} = \sum_{\ell \in \mathcal{N}_k} a_{\ell k} \left( \mathbf{w}'_{\ell,i-1} - \mu \widehat{\nabla J}_\ell(\mathbf{w}'_{\ell,i-1}) \right) + \sum_{\ell \in \mathcal{N}_k} a_{\ell k} \mathbf{v}_{\ell k,i} \quad (11.43)$$

<sup>1</sup> E. Rizk, S. Vlaski, A. H. Sayed, “Enforcing Privacy in Distributed Learning with Performance Guarantees”, available as arXiv:2301.06412, 2023.

Where  $\mathbf{x}'_{\ell,i}$  differs from  $\mathbf{x}_{\ell,i}$  only for one agent, say agent 1, and  $\mathbf{w}'_{k,i}$  differs from  $\mathbf{w}_{k,i}$  as a result of the propagation of the alternative data  $\mathbf{x}'_{1,i}$  through  $\widehat{\nabla J}_\ell(\mathbf{w}'_{\ell,i-1})$ . The sensitivity of the algorithm is then defined as the deviation of  $\mathbf{w}'_{k,i}$  from  $\mathbf{w}_{k,i}$ , and measures the influence that agent 1 has on the evolution of the algorithm. A higher sensitivity is generally associated with a need for larger privacy perturbations, since a larger influence needs to be masked. We can bound:

$$\begin{aligned}
& \|\mathbf{w}_{k,i} - \mathbf{w}'_{k,i}\| \\
&= \left\| \sum_{\ell \in \mathcal{N}_k} a_{\ell k} \left( \mathbf{w}_{\ell,i-1} - \mu \widehat{\nabla J}_\ell(\mathbf{w}_{\ell,i-1}) \right) \right. \\
&\quad \left. - \sum_{\ell \in \mathcal{N}_k} a_{\ell k} \left( \mathbf{w}'_{\ell,i-1} - \mu \widehat{\nabla J}_\ell(\mathbf{w}'_{\ell,i-1}) \right) \right\| \\
&\stackrel{(b)}{\leq} \sum_{\ell \in \mathcal{N}_k} a_{\ell k} \|\mathbf{w}_{\ell,i-1} - \mu \widehat{\nabla J}_\ell(\mathbf{w}_{\ell,i-1}) \\
&\quad - \mathbf{w}'_{\ell,i-1} + \mu \widehat{\nabla J}_\ell(\mathbf{w}'_{\ell,i-1})\| \\
&\stackrel{(c)}{\leq} \sum_{\ell \in \mathcal{N}_k} a_{\ell k} \|\mathbf{w}_{\ell,i-1} - \mathbf{w}'_{\ell,i-1}\| \\
&\quad + \mu \sum_{\ell \in \mathcal{N}_k} a_{\ell k} \|\widehat{\nabla J}_\ell(\mathbf{w}_{\ell,i-1}) - \widehat{\nabla J}_\ell(\mathbf{w}'_{\ell,i-1})\| \\
&\stackrel{(d)}{\leq} \left( \sum_{\ell \in \mathcal{N}_k} a_{\ell k} \right) \max_{\ell} \|\mathbf{w}_{\ell,i-1} - \mathbf{w}'_{\ell,i-1}\| + \mu \left( \sum_{\ell \in \mathcal{N}_k} a_{\ell k} \right) 2G \\
&= \Delta(i-1) + \mu 2G
\end{aligned} \tag{11.44}$$

where (b) follows from Jensen's inequality, (c) makes use of the triangle inequality and (d) follows for bounded gradient approximations. We conclude that:

$$\Delta(i) \triangleq \max_k \|\mathbf{w}_{k,i} - \mathbf{w}'_{k,i}\| \leq \Delta(i-1) + \mu 2G \tag{11.45}$$

and hence (11.41) after iterating with  $\Delta(0) = 0$ .  $\square$

**DEFINITION 11.4 ( $\epsilon$ -differential privacy of diffusion).** We say that the diffusion recursion (11.33)–(11.35) is  $\epsilon(i)$ -differentially private for agent 1 at time  $i$  if:

$$\frac{f\left(\left\{\left\{\psi_{1\ell,n}\right\}_{\ell \neq 1 \in \mathcal{N}_1}\right\}_{n=0}^i\right)}{f\left(\left\{\left\{\psi'_{1\ell,n}\right\}_{\ell \neq 1 \in \mathcal{N}_1}\right\}_{n=0}^i\right)} \leq e^{\epsilon(i)} \tag{11.46}$$

where  $f(\cdot)$  denotes the probability density function and  $\left\{\left\{\psi_{1\ell,n}\right\}_{\ell \neq 1 \in \mathcal{N}_1}\right\}_{n=0}^i$  collects all quantities transmitted by agent 1 to any of its neighbors during the operation of the algorithm, while excluding its local iterates  $\psi_{11,n}$ , which are kept private.

**THEOREM 11.1 (Privacy cost of the diffusion algorithm).** Suppose (11.33)–(11.35) employs homomorphic perturbations constructed as in (11.38). Then, at

time  $i$ , algorithm (11.33)–(11.35) is  $\epsilon(i)$ -differentially private according to (11.46), with:

$$\epsilon(i) = \mu \frac{G(i^2 + i)}{b_v} \quad (11.47)$$

**Proof:** Under construction (11.38), every agent shares the same perturbed estimate  $\psi_{1\ell,n} = \phi_{1,n} + \mathbf{v}_n$  with all of its neighbors (excluding itself), and we have:

$$\frac{f\left(\left\{\left\{\psi_{1\ell,n}\right\}_{\ell \neq 1 \in \mathcal{N}_1}\right\}_{n=0}^i\right)}{f\left(\left\{\left\{\psi'_{1\ell,n}\right\}_{\ell \neq 1 \in \mathcal{N}_1}\right\}_{n=0}^i\right)} = \frac{f\left(\left\{\phi_{1,n} + \mathbf{v}_n\right\}_{n=0}^i\right)}{f\left(\left\{\phi'_{1,n} + \mathbf{v}_n\right\}_{n=0}^i\right)} \quad (11.48)$$

We construct an iterative argument by induction. Suppose we have  $\epsilon(i-1)$  differential privacy at time  $i-1$ . We can expand by Bayes' rule:

$$\begin{aligned} & \frac{f\left(\left\{\phi_{1,n} + \mathbf{v}_n\right\}_{n=0}^i\right)}{f\left(\left\{\phi'_{1,n} + \mathbf{v}_n\right\}_{n=0}^i\right)} \\ &= \frac{f\left(\left\{\phi_{1,n} + \mathbf{v}_n\right\}_{n=0}^{i-1}\right) \cdot f\left(\phi_{1,i} + \mathbf{v}_i \mid \left\{\phi_{1,n} + \mathbf{v}_n\right\}_{n=0}^{i-1}\right)}{f\left(\left\{\phi'_{1,n} + \mathbf{v}_n\right\}_{n=0}^{i-1}\right) \cdot f\left(\phi'_{1,i} + \mathbf{v}_i \mid \left\{\phi'_{1,n} + \mathbf{v}_n\right\}_{n=0}^{i-1}\right)} \\ &\stackrel{(a)}{\leq} e^{\epsilon(i-1)} \frac{f\left(\phi_{1,i} + \mathbf{v}_i \mid \left\{\phi_{1,n} + \mathbf{v}_n\right\}_{n=0}^{i-1}\right)}{f\left(\phi'_{1,i} + \mathbf{v}_i \mid \left\{\phi'_{1,n} + \mathbf{v}_n\right\}_{n=0}^{i-1}\right)} \\ &\stackrel{(b)}{\leq} e^{\epsilon(i-1)} \frac{f\left(\phi_{1,i} + \mathbf{v}_i \mid \mathcal{F}_{i-1}\right)}{f\left(\phi'_{1,i} + \mathbf{v}_i \mid \mathcal{F}'_{i-1}\right)} \end{aligned} \quad (11.49)$$

where in (a) we used the induction hypothesis and in (b) we defined  $\mathcal{F}_{i-1} \triangleq \left\{\phi_{1,n} + \mathbf{v}_n\right\}_{n=0}^{i-1}$  for brevity. Since  $\mathbf{v}_i$  is sampled in an i.i.d. fashion, it is independent of  $\mathcal{F}_{i-1}, \mathcal{F}'_{i-1}$  and  $\phi_{1,i}, \phi'_{1,i}$ , hence we have:

$$\begin{aligned} \frac{f\left(\phi_{1,i} + \mathbf{v}_i \mid \mathcal{F}_{i-1}\right)}{f\left(\phi'_{1,i} + \mathbf{v}_i \mid \mathcal{F}'_{i-1}\right)} &= \frac{\frac{1}{2b_v} e^{-\frac{\|\phi_{1,i} - \mathbf{v}_i\|}{b_v}}}{\frac{1}{2b_v} e^{-\frac{\|\phi'_{1,i} - \mathbf{v}_i\|}{b_v}}} \\ &= e^{\frac{\|\phi'_{1,i} - \mathbf{v}_i\| - \|\phi_{1,i} - \mathbf{v}_i\|}{b_v}} \\ &\stackrel{(a)}{\leq} e^{\frac{\|\phi'_{1,i} - \phi_{1,i}\|}{b_v}} \end{aligned} \quad (11.50)$$

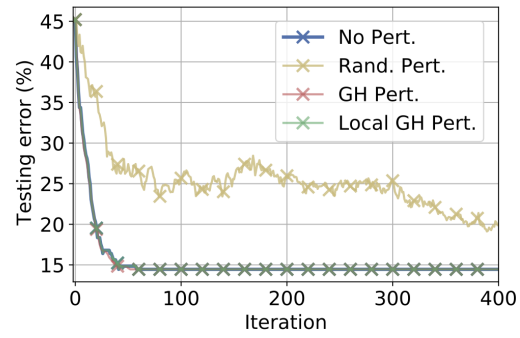
where (a) follows from  $\|x\| = \|x - y + y\| \leq \|x - y\| + \|y\|$ . We can then bound with probability one:

$$\begin{aligned} & \|\phi'_{1,i} - \phi_{1,i}\| \\ &= \|\mathbf{w}'_{k,i-1} - \mu \nabla Q_k(\mathbf{w}'_{k,i-1}; \mathbf{x}'_{k,i}) - \mathbf{w}_{k,i-1} + \mu \nabla Q_k(\mathbf{w}_{k,i-1}; \mathbf{x}_{k,i})\| \\ &\leq \|\mathbf{w}'_{k,i-1} - \mathbf{w}_{k,i-1}\| + \mu \|\nabla Q_k(\mathbf{w}'_{k,i-1}; \mathbf{x}'_{k,i}) - \nabla Q_k(\mathbf{w}_{k,i-1}; \mathbf{x}_{k,i})\| \\ &\leq \Delta(i-1) + \mu 2G = \mu 2Gi \end{aligned} \quad (11.51)$$

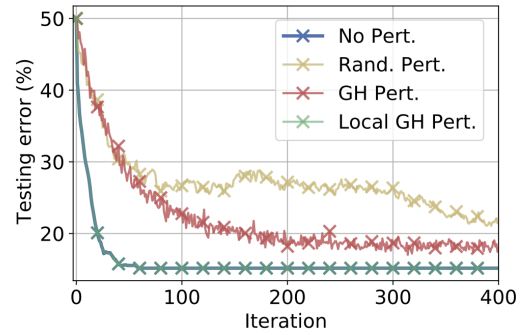
We conclude:

$$\epsilon(i) \leq \epsilon(i-1) + \frac{\mu 2Gi}{b_v} \implies \epsilon(i) \leq \frac{\mu 2G}{b_v} \sum_{n=1}^i i = \mu \frac{G(i^2 + i)}{b_v} \quad (11.52)$$

□



(a) Centroid testing error



(b) Average individual testing error

**Figure 11.1** Performance of privatized variants of the diffusion algorithm. Taken from [Rizk, Vlaski, Sayed 2023].