

# Group Privacy for Personalized Federated Learning

Filippo Galli<sup>1</sup><sup>a</sup>, Sayan Biswas<sup>2,4</sup><sup>b</sup>, Kangsoo Jung<sup>2</sup><sup>c</sup>, Tommaso Cucinotta<sup>3</sup><sup>d</sup>  
and Catuscia Palamidessi<sup>2,4</sup><sup>e</sup>

<sup>1</sup>*Scuola Normale Superiore, Pisa, Italy*

<sup>2</sup>*INRIA, Palaiseau, France*

<sup>3</sup>*Scuola Superiore Sant'Anna, Pisa, Italy*

<sup>4</sup>*LIX, École Polytechnique, Palaiseau, France*

*filippo.galli@sns.it, {sayan.biswas, gangsoo.zeong}@inria.fr, tommaso.cucinotta@santannapisa.it, catuscia@lix.polytechnique.fr*


Keywords: Federated Learning, Differential Privacy,  $d$ -Privacy, Personalized Models.


Abstract: Federated learning (FL) is a particular type of distributed, collaborative machine learning, where participating clients process their data locally, sharing only updates of the training process. Generally, the goal is the privacy-aware optimization of a statistical model's parameters by minimizing a cost function of a collection of datasets which are stored locally by a set of clients. This process exposes the clients to two issues: leakage of private information and lack of personalization of the model. To mitigate the former, differential privacy and its variants serve as a standard for providing formal privacy guarantees. But often the clients represent very heterogeneous communities and hold data which are very diverse. Therefore, aligned with the recent focus of the FL community to build a framework of personalized models for the users representing their diversity, it is of utmost importance to protect the clients' sensitive and personal information against potential threats. To address this goal we consider  $d$ -privacy, also known as metric privacy, which is a variant of local differential privacy, using a metric-based obfuscation technique that preserves the topological distribution of the original data. To cope with the issues of protecting the privacy of the clients and allowing for personalized model training, we propose a method to provide group privacy guarantees exploiting some key properties of  $d$ -privacy which enables personalized models under the framework of FL. We provide theoretical justifications to the applicability and experimental validation on real-world datasets to illustrate the working of the proposed method.


## 1 INTRODUCTION


With the modern developments in machine learning, user data collection has become ubiquitous, often disclosing sensitive personal information with increasing risks of users' privacy violations (Le Métayer and De, 2016; NIST, 2021). To try and curb such threats, Federated Learning (McMahan et al., 2017a) was introduced as a collaborative machine learning paradigm where the users' devices, on top of harvesting user data, directly train a global predictive model, without ever sending the raw data to a central server. On the one hand, this paradigm has received much at-


tention with the appealing promises of guaranteeing user privacy and model performance. On the other hand, given the heterogeneity of the data distributions among clients, training convergence is not guaranteed and model utility may be reduced by local updates. Many works have thus focused on the topic of personalized federated learning, to tailor a set of models to clusters of users with similar data distributions (Ghosh et al., 2020; Mansour et al., 2020; Sattler et al., 2020). On a similar note, other lines of work have also showed that relying on avoiding the release of user's raw data only provides a lax protection to potential attacks violating the users' privacy (Hitaj et al., 2017), (Nasr et al., 2019), (Zhu et al., 2019). To tackle this problem, researchers have been exploring the application of Differential Privacy (DP) (Dwork et al., 2006b; Dwork et al., 2006a) to federated learning, in order to quantify and provide

<sup>a</sup> <https://orcid.org/0000-0002-2279-3545>

<sup>b</sup> <https://orcid.org/0000-0002-2115-1495>

<sup>c</sup> <https://orcid.org/0000-0003-2070-1050>

<sup>d</sup> <https://orcid.org/0000-0002-0362-0657>

<sup>e</sup> <https://orcid.org/0000-0003-4597-7002>

privacy to users participating in the optimization. The goal of differential privacy mechanisms is to introduce randomness in the information released by the clients, such that each user’s contribution to the final model can be made probabilistically indistinguishable up to a certain likelihood factor. To bound this factor, the domain of secrets (i.e. the parameter space in FL) is artificially bounded, be it to provide central (Andrew et al., 2021; McMahan et al., 2017b) or local DP guarantees (Truex et al., 2020; Zhao et al., 2020). When users share with the central server their locally updated models for averaging, constraining the optimization to a subset of  $\mathbb{R}^n$  can have destructive effects, e.g. when the optimal model parameters for a certain cluster of users may be found outside such bounded domain. Therefore, in this work we aim to address the problems of personalization and local privacy protection by adopting a generalization of DP, i.e.  $d$ -privacy or metric-based privacy (Chatzikokolakis et al., 2013). This notion of privacy does not require a bounded domain and provides guarantees dependent on the distance between any two points in the parameter space. Thus, under the minor assumption that clients with similar data distributions will have similar optimal fitting parameters,  $d$ -privacy will provide them with stronger indistinguishability guarantees. Conversely, privacy guarantees degrade gracefully for clients whose data distributions are vastly different. More precisely, our key contributions in this paper are outlined as follows:

1. We provide an algorithm for the collaborative training of machine learning models, which builds on top of state-of-the-art strategies for model personalization.
2. We formalize the privacy guarantees in terms of  $d$ -privacy. To the best of our knowledge, this is the first time that  $d$ -privacy is used in the context of machine learning for protecting user information.
3. We study the Laplace mechanism on high dimensions, under Euclidean distance, based on a generalization of the Laplace distribution in  $\mathbb{R}$ , and we give a closed form expression.
4. We provide an efficient procedure for sampling from such

The rest of the paper is organized as follows. Section 2 introduces fundamental notions for federated learning and differential privacy and discusses related work. Section 3 explains the proposed algorithm for personalized federated learning with group privacy. Section 4 validates the proposed procedure through experimental results. Section 5 concludes and discusses future work.

## 2 BACKGROUND

### 2.1 Related Works

Federated optimization has shown to be underperforming when the local datasets are samples of non-congruent distributions, failing to minimize both the local and global objectives at the same time. In (Ghosh et al., 2020; Mansour et al., 2020; Sattler et al., 2020), the authors investigate different meta-algorithms for personalization. Claims of user privacy preservation are based solely on the clients releasing updated models (or model updates) instead of transferring the raw data to the server, with potentially dramatic effects. To confront this issue, a number of works have focused on the privatization of the (federated) optimization algorithm under the framework of DP (Abadi et al., 2016; Geyer et al., 2017; McMahan et al., 2017b; Andrew et al., 2021) who adopt DP to provide defenses against an *honest-but-curious* adversary. Even in this setting though, no protection is guaranteed against sample reconstruction from the local datasets (Zhu et al., 2019), using the client updates. Different strategies have been tried to provide local privacy guarantees, either from the perspective of cryptography (Bonawitz et al., 2016), or under the framework of local DP (Truex et al., 2020; Agarwal et al., 2018; Hu et al., 2020). In particular in (Hu et al., 2020) the authors address the problem of personalized and locally differentially private federated learning, but for the simple case of convex, 1-Lipschitz cost functions of the inputs. Note that this assumption is unrealistic in most machine learning models, and it excludes many statistical modeling techniques, notably neural networks.

### 2.2 Personalized Federated Learning

The problem can be cast under the framework of stochastic optimization and we adopt the notation of (Ghosh et al., 2020) to find the set of minimizers  $\theta_j^* \in \mathbb{R}^n$  with  $j \in \{1, \dots, k\}$  of the cost functions

$$F(\theta_j) = \mathbb{E}_{z \sim \mathcal{D}_j} [f(\theta_j; z)], \quad (1)$$

where  $\{\mathcal{D}_1, \dots, \mathcal{D}_k\}$  are the data distributions which can only be accessed through a collection of client datasets  $Z_c = \{z | z \sim \mathcal{D}_j, z \in \mathbb{D}\}$  for some  $j \in \{1, \dots, k\}$  with  $c \in C = \{1, \dots, N\}$  the set of clients, and  $\mathbb{D}$  a generic domain of data points.  $C$  is partitioned in  $k$  disjoint sets

$$S_j^* = \{c \in C | \forall z \in Z_c, z \sim \mathcal{D}_j\} \quad \forall j \in \{1, \dots, k\} \quad (2)$$

The mapping  $c \rightarrow j$  is unknown and we rely on estimates  $S_j$  of the membership of  $Z_c$  to compute the

empirical cost functions

$$\begin{aligned}\tilde{F}(\theta_j) &= \frac{1}{|S_j|} \sum_{c \in S_j} \tilde{F}_c(\theta_j; Z_c); \\ \tilde{F}_c(\theta_j; Z_c) &= \frac{1}{|Z_c|} \sum_{z_i \in Z_c} f(\theta; z_i)\end{aligned}\quad (3)$$

The cost function  $f: \mathbb{R}^n \times \mathbb{D} \mapsto \mathbb{R}_{\geq 0}$  is applied on  $z \in \mathbb{D}$ , parametrized by the vector  $\theta_j \in \mathbb{R}^n$ . Thus, the optimization aims to find,  $\forall j \in \{1, \dots, k\}$ ,

$$\tilde{\theta}_j^* = \arg \min_{\theta_j} \tilde{F}(\theta_j) \quad (4)$$

### 2.3 Privacy

$d$ -privacy (Chatzikokolakis et al., 2013) is a generalization of DP for any domain  $\mathcal{X}$ , representing the space of original data, endowed with a distance measure  $d: \mathcal{X}^2 \mapsto \mathbb{R}_{\geq 0}$ , and any space of secrets  $\mathcal{Y}$ . A random mechanism  $\mathcal{R}: \mathcal{X} \mapsto \mathcal{Y}$  is called  $\varepsilon$ - $d$ -private if for all  $x_1, x_2 \in \mathcal{X}$  and measurable  $S \subseteq \mathcal{Y}$ :

$$\mathbb{P}[\mathcal{R}(x_1) \in S] \leq e^{\varepsilon d(x_1, x_2)} \mathbb{P}[\mathcal{R}(x_2) \in S] \quad (5)$$

Note that when  $\mathcal{X}$  is the domain of databases, and  $d$  is the distance on the Hamming graph of their adjacency relation, then Equation (5) results in the standard definition of DP in (Dwork et al., 2006b; Dwork et al., 2006a). In this work we will have though that  $\theta \in \mathbb{R}^n = \mathcal{X} = \mathcal{Y}$ . The main motivation behind the use of  $d$ -privacy is to preserve the topology of the parameter distributions among clients, i.e. to have that, in expectation, clients with close model parameters in the non-privatized space  $\mathcal{X}$  will communicate close model parameters in the privatized space  $\mathcal{Y}$ .

## 3 AN ALGORITHM FOR PRIVATE AND PERSONALIZED FEDERATED LEARNING

We propose an algorithm for personalized federated learning with local guarantees to provide group privacy (Algorithm 1). Locality refers to the sanitization of the information released by the client to the server, whereas group privacy refers to indistinguishability with respect to a neighborhood of clients, defined with respect to a certain distance metric. Thus we proceed to define *neighborhood* and *group*.

**Definition 3.1.** For any model parametrized by  $\theta_0 \in \mathbb{R}^n$ , we define its  $r$ -neighborhood as the set of points in the parameter space which are at a  $L_2$  distance of at most  $r$  from  $\theta_0$ , i.e.,  $\{\theta \in \mathbb{R}^n: \|\theta_0 - \theta\|_2 \leq r\}$ . Clients whose models are parametrized by  $\theta \in \mathbb{R}^n$  in the same

$r$ -neighborhood are said to be in the same *group*, or *cluster*.

Algorithm 1 is motivated by the Iterative Federated Clustering Algorithm (IFCA) (Ghosh et al., 2020) and builds on top of it to provide formal privacy guarantees. The main differences lie in the introduction of the `SanitizeUpdate` function described in Algorithm 2 and  $k$ -means for server-side clustering of the updated models.

### 3.1 The Laplace Mechanism Under Euclidean Distance in $\mathbb{R}^n$

Algorithm 2's `SanitizeUpdate` is based on a generalization of the Laplace mechanism under Euclidean distance to  $\mathbb{R}^n$ , introduced in (Andrés et al., 2013) for geo-indistinguishability in  $\mathbb{R}^2$ . The motivation to adopt the  $L_2$  norm as distance measure is twofold. First, clustering is performed on  $\theta$  with the  $k$ -means algorithm under Euclidean distance. Since we define clusters or groups of users based on how close their model parameters are under  $L_2$  norm, we are looking for a  $d$ -privacy mechanism that obfuscates the reported values within a certain group and allows the server to discern among users belonging to different clusters. Second, parameters that are sanitized by equidistant noise vectors in  $L_2$  norm are also equiprobable by construction and lead to the same bound in the increase of the cost function in first order approximation, as shown in Proposition 3.2. The Laplace mechanism under Euclidean distance in a generic space  $\mathbb{R}^n$  is defined in Proposition 3.1. Proofs of all Propositions and Theorems are included in Appendix 5.

**Proposition 3.1.** Let  $\mathcal{L}_\varepsilon: \mathbb{R}^n \mapsto \mathbb{R}^n$  be the Laplace mechanism with distribution  $\mathcal{L}_{x_0, \varepsilon}(x) = \mathbb{P}[\mathcal{L}_\varepsilon(x_0) = x] = K e^{-\varepsilon d(x, x_0)}$  with  $d(\cdot)$  being the Euclidean distance. If  $\rho \sim \mathcal{L}_{x_0, \varepsilon}(x)$ , then:

1.  $\mathcal{L}_{x_0, \varepsilon}$  is  $\varepsilon$ - $d$ -private and  $K = \frac{\varepsilon^n \Gamma(\frac{n}{2})}{2\pi^{\frac{n}{2}} \Gamma(n)}$
2.  $\|\rho\|_2 \sim \gamma_{\varepsilon, n}(r) = \frac{\varepsilon^n e^{-\varepsilon r} r^{n-1}}{\Gamma(n)}$
3. The  $i^{\text{th}}$  component of  $\rho$  has variance  $\sigma_{\rho_i}^2 = \frac{n+1}{\varepsilon^2}$

where  $\Gamma(n)$  is the Gamma function defined for positive reals as  $\int_0^\infty t^{n-1} e^{-t} dt$  which reduces to the factorial function whenever  $n \in \mathbb{N}$ .

**Proposition 3.2.** Let  $y = f(x, \theta)$  be the fitting function of a machine learning model parameterized by  $\theta$ , and  $(X, Y) = Z$  the dataset over which the RMSE loss function  $F(Z, \theta)$  is to be minimized, with  $x \in X$  and  $y \in Y$ . If  $\rho \sim \mathcal{L}_{0, \varepsilon}$ , the bound on the increase of the cost function does not depend on the direction of  $\rho$ ,

---

Algorithm 1: An algorithm for personalized federated learning with formal privacy guarantees in local neighborhoods.

---

Input: number of clusters  $k$ ; initial hypotheses  $\theta_j^{(0)}, j \in \{1, \dots, k\}$ ; number of rounds  $T$ ; number of users per round  $U$ ; number of local epochs  $E$ ; local step size  $s$ ; user batch size  $B_s$ ; noise multiplier  $v$ ; local dataset  $Z_c$  held by user  $c$ .

**for**  $t = \{0, 1, \dots, T - 1\}$  **do** ▷ Server-side loop

$C^{(t)} \leftarrow \text{SampleUserSubset}(U)$

BroadcastParameterVectors( $C^{(t)}; \theta_j^{(t)}, j \in \{1, \dots, k\}$ )

**for**  $c \in C^{(t)}$  **do in parallel** ▷ Client-side loop

$\bar{j} = \arg \min_{j \in \{1, \dots, k\}} F_c(\theta_j^{(t)}; Z_c)$

$\theta_{\bar{j}, c}^{(t)} \leftarrow \text{LocalUpdate}(\theta_{\bar{j}}^{(t)}; s; E; Z_c)$

$\hat{\theta}_{\bar{j}, c}^{(t)} \leftarrow \text{SanitizeUpdate}(\theta_{\bar{j}, c}^{(t)}; v)$

**end for**

$\{S_1, \dots, S_k\} = \text{k-means}(\hat{\theta}_{\bar{j}, c}^{(t)}, c \in C^{(t)}; \theta_j^{(t)}, j \in \{1, \dots, k\})$

$\theta_j^{(t+1)} \leftarrow \frac{1}{|S_j|} \sum_{c \in S_j} \hat{\theta}_{\bar{j}, c}^{(t)}, \quad \forall j \in \{1, \dots, k\}$

**end for**

---

Algorithm 2: SanitizeUpdate obfuscates a vector  $\theta \in \mathbb{R}^n$ , with a Laplacian noise tuned on the radius of a certain neighborhood and centered in 0.

---

**function** SANITIZEUPDATE( $\theta_{\bar{j}}^{(t)}; \theta_{\bar{j}, c}^{(t)}; v$ )

$\delta_c^{(t)} = \theta_{\bar{j}, c}^{(t)} - \theta_{\bar{j}}^{(t)}$

$\epsilon = \frac{n}{v \|\delta_c^{(t)}\|}$

Sample  $\rho \sim \mathcal{L}_{0, \epsilon}(x)$

$\hat{\theta}_{\bar{j}, c}^{(t)} = \theta_{\bar{j}, c}^{(t)} + \rho$

**return**  $\hat{\theta}_{\bar{j}, c}^{(t)}$

**end function**

---

in first order approximation, and:

$$\begin{aligned} \|F(Z, \theta + \rho)\|_2 - \|F(Z, \theta)\|_2 &\leq \\ \|J_f(X, \theta)\|_2 \|\rho\|_2 + o(\|J_f(X, \theta) \cdot \rho\|_2) \end{aligned} \quad (6)$$

The results in Proposition 3.1 allow to reduce the problem of sampling a point from Laplace to i) sampling the norm of such point according to the result in Item 2 of Proposition 3.1 and then ii) sample uniformly a unit (directional) vector from the hypersphere in  $\mathbb{R}^n$ . Much like DP,  $d$ -privacy provides a means to compute the total privacy parameters in case of repeated queries, a result known as Compositionality Theorem for  $d$ -privacy 3.1. Although it was known as a folk result, we provide a formal proof in Appendix 5.

**Theorem 3.1.** *Let  $\mathcal{K}_i$  be  $(\epsilon_i)$ - $d$ -private mechanism for  $i \in \{1, 2\}$ . Then their independent composition is  $(\epsilon_1 + \epsilon_2)$ - $d$ -private.*

### 3.2 A Heuristic for Defining the Neighborhood of a Client

At the  $t^{\text{th}}$  iteration, when a user  $c$  calls the SanitizeUpdate routine in Algorithm 2, it has already received a set of hypotheses, optimized  $\theta_{\bar{j}}^{(t)}$  (the one that fits best its data distribution), and got  $\theta_{\bar{j}, c}^{(t)}$ . It is reasonable to assume that clients whose datasets are sampled from the same underlying data distribution  $\mathcal{D}_j$  will perform an update similar to  $\delta_c^{(t)}$ . Therefore, we enforce points which are within the  $\delta_c^{(t)}$ -neighborhood of  $\hat{\theta}_{\bar{j}, c}^{(t)}$  to be indistinguishable. To provide this guarantee, we tune the Laplace mechanism such that the points within the neighborhood are  $\epsilon \|\delta_c^{(t)}\|_2$  differentially private. With the choice of  $\epsilon = n / (v \|\delta_c^{(t)}\|_2)$ , one finds that  $\epsilon \|\delta_c^{(t)}\|_2 = n/v$ , and we call  $v$  the *noise multiplier*. It is straightforward to observe that the larger the value of  $v$  gets, the stronger is the privacy guarantee. This results from the norm of the noise vector sampled from the Laplace distribution being distributed according to Equation (12) whose expected value is  $\mathbb{E}[\gamma_{\epsilon, n}(r)] = n/\epsilon$ .

## 4 EXPERIMENTS

### 4.1 Synthetic Data

We generate data according to  $k = 2$  different distributions:  $y = x^T \theta_i^* + u$  and  $u \sim \text{Uniform}[0, 1], \forall i \in \{1, 2\}$  and  $\theta_1^* = [+5, +6]^T, \theta_2^* = [+4, -4.5]^T$ . We then as-

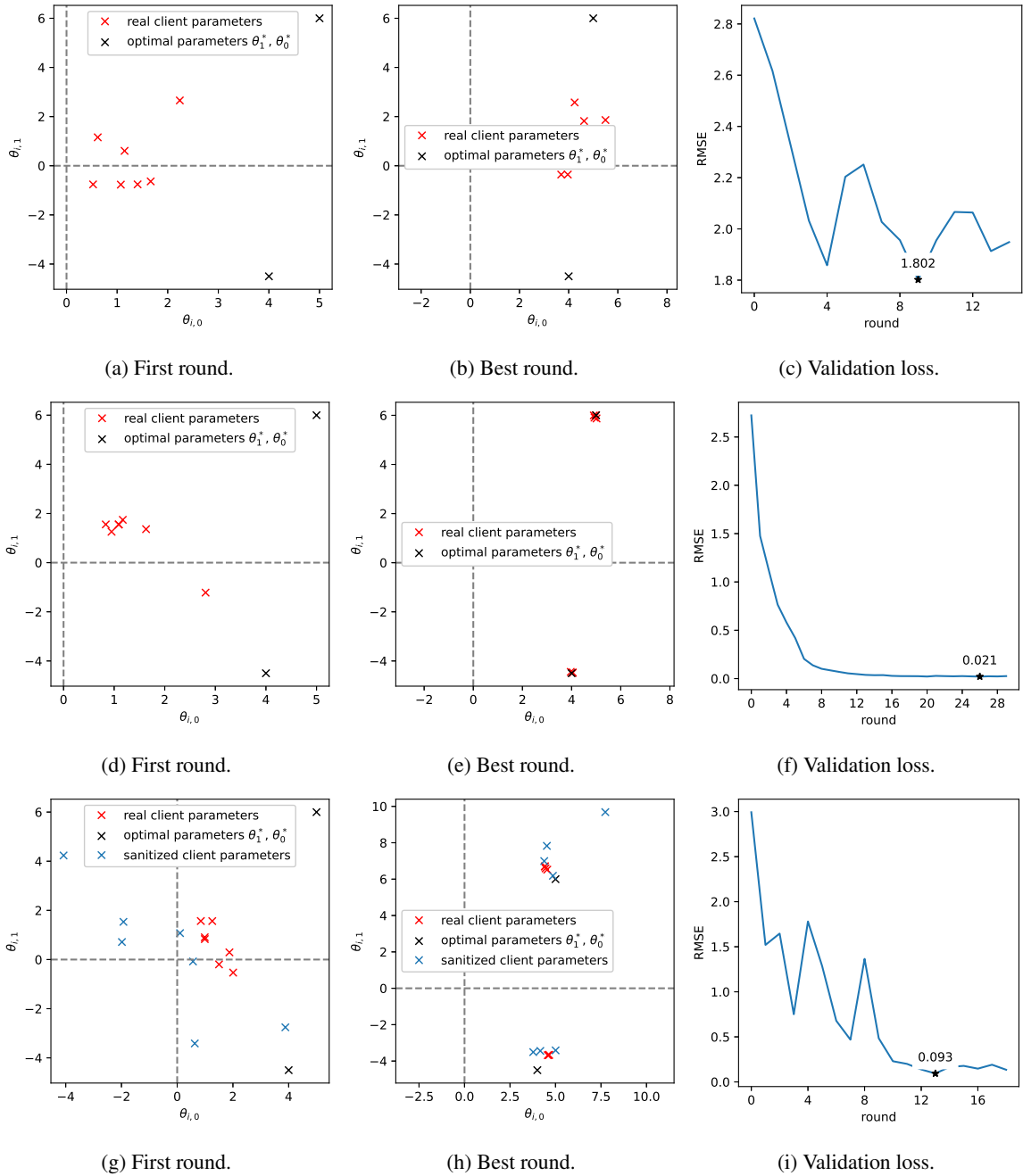


Figure 1: Learning federated linear models with: (a, b, c) one initial hypothesis and non-sanitized communication, (d, e, f) two initial hypotheses and non-sanitized communication, (g, h, i) two initial hypotheses and sanitized communication. The first two figures of each row show the parameter vectors released by the clients to the server.

sess how training progresses as we move from the Federated Averaging (Konečný et al., 2016) (Figure 1a, 1b, 1c), to IFCA (Figure 1d, 1e, 1f), and finally Algorithm 1 (Figure 1g, 1h, 1i). When using Federated Averaging, there seems to be an obvious problem, that is using one single hypothesis is not enough to capture the diversity in the data distributions, re-

sulting in the final parameters settling somewhere in between the optimal parameters (Figure 1b). On the contrary using IFCA, shows that having multiple initial hypotheses helps in improving the performance when the clients have heterogeneous data, as the optimized clients parameters almost overlap the optimal parameters (Figure 1e). Adopting our algorithm

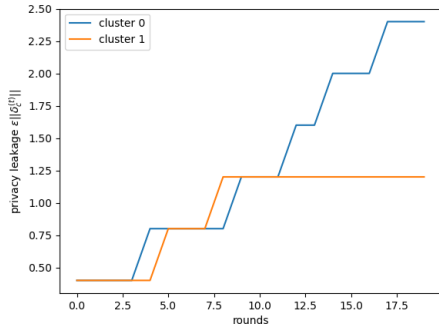


Figure 2: Synthetic data: max privacy leakage among clients. Privacy leakage is constant when clients with the largest privacy leakage are not sampled (by chance) to participate in those rounds.

shows that on top of providing formal guarantees, we can still achieve great results in terms of proximity to the optimal parameters (Figure 1h) and reduction of the loss function (Figure 1i). Figure 2 provides the maximum value of privacy leakage clients incur into, per cluster. Further details about the experimental settings are provided in Appendix 5.

## 4.2 Hospital Charge Data

This experiment is performed on the Hospital Charge Dataset by the Centers for Medicare and Medicaid Services of the US Government (CMMS, 2021). The healthcare providers are considered the set of clients willing to train a machine learning model with federated learning. The goal is to predict the cost of a service given where it is performed in the country, and what kind of procedure it is. More details on the preprocessing and training settings are included in Appendix 5. To assess the trade-off between privacy, personalization and accuracy, a different number of initial hypotheses has been checked, as it is not known a-priori how many distributions generated the data. Accuracy has been evaluated at different levels of the noise multiplier  $\nu$ . Note that, using Algorithm 1 with 1 hypothesis results in the Federated Averaging algorithm. Figure 3 shows that adopting multiple hypotheses drastically reduces the RMSE loss function. This is especially true when moving from 1 to 3 hypotheses. Additionally, we highlight how increasing the number of hypotheses also helps in curbing the effects of the noise multiplier even when it reaches high levels, on the right hand side of the picture, making a compelling case for adopting formal privacy guarantees when a slight increase in the cost function is admissible. Figure 4 provides the empirical privacy leakage distribution of the clients involved in a particular training configuration. Table 1 shows privacy leakage statistics over multiple rounds and for all con-

figurations.

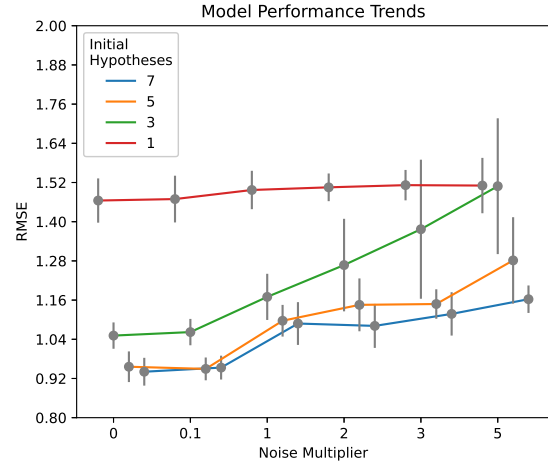


Figure 3: RMSE for models trained with Algorithm 1 on the Hospital Charge Dataset. Error bars show  $\pm\sigma$ , with  $\sigma$  the empirical standard deviation. Lower RMSE values are better for accuracy.

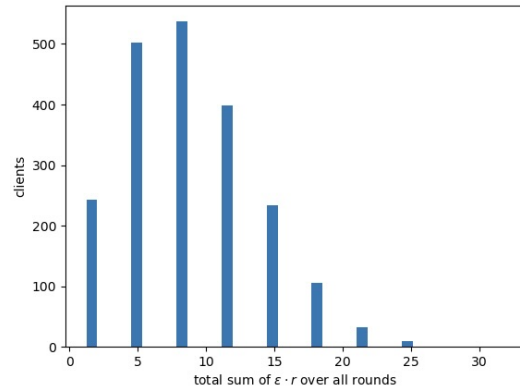


Figure 4: Hospital charge data: the empirical distribution of the privacy budget over the clients for:  $\nu = 3$ , 5 initial hypotheses,  $\text{seed} = 3$ ,  $r$  is the radius of the neighborhood, the total number of clients is 2062.

## 4.3 FEMNIST Image Classification

This task consists of image-based character recognition on the FEMNIST dataset (Caldas et al., 2018). Details on the experimental settings are in Appendix 5. With the choice of the range of noise multipliers  $\nu$  the corresponding value for the privacy leakage  $\epsilon \|\delta_c^{(t)}\|_2 = n/\nu$  would be enormous, considering a CNN with  $n = 206590$  parameters, providing no meaningful theoretical privacy guarantees. This is a common issue for local privacy mechanisms (Bassily et al., 2017), and it comes from the linear dependence of the expected value of the norm of the noise vector on  $n$ :  $\mathbb{E}[\gamma_{\epsilon,n}(r)] = n/\epsilon$ . Still, it is possible to validate,

Table 1: Hospital charge data: median and maximum local privacy budgets over the whole set of clients, averaged over 10 runs with different seeds.  $v = 0$  means no privacy guarantee.

v	Hypotheses			
	7	5	3	1
0	-, -	-, -	-, -	-, -
0.1	517.0, 1551.0	418.0, 1342.0	473.0, 1386.0	528.0, 1540.0
1	36.3, 126.5	40.7, 127.6	44.0, 138.6	49.5, 147.4
2	15.4, 57.8	14.3, 54.5	22.0, 69.3	21.5, 66.6
3	7.7, 32.3	8.4, 36.7	12.5, 40.0	12.1, 40.0
5	5.7, 21.3	5.9, 22.0	5.5, 21.6	5.3, 20.9

Table 2: Effects of increasing the noise multiplier on the validation accuracy and standard deviation.

v	Cross Entropy loss		RMSE loss	
	Average Accuracy	Standard Deviation	Average Accuracy	Standard Deviation
0	0.832	$\pm 0.012$	0.801	$\pm 0.001$
0.001	0.843	$\pm 0.006$	0.813	$\pm 0.014$
0.01	0.832	$\pm 0.017$	0.805	$\pm 0.008$
0.1	0.834	$\pm 0.026$	0.808	$\pm 0.019$
1	0.834	$\pm 0.014$	0.814	$\pm 0.012$
3	0.835	$\pm 0.017$	0.825	$\pm 0.010$
5	0.812	$\pm 0.016$	0.787	$\pm 0.003$
10	0.692	$\pm 0.002$	0.687	$\pm 0.014$
15	0.561	$\pm 0.005$	0.622	$\pm 0.003$

in practice, whether this particular generalization of the Laplace mechanism can protect against a *specific* attack: DLG (Zhu et al., 2019). Figure 5 and Table 2 report the results of varying the noise multiplier values. When  $v = 10^{-3}$  the ground truth image is fully reconstructed. Up to  $v = 10^{-1}$  we see that at least partial reconstruction is possible. For  $v \geq 1$  we see that, experimentally, the DLG attack fails to reconstruct input samples when we protect the client-server communication with the mechanism in Proposition 3.1.

## 5 CONCLUSIONS

We use the framework of  $d$ -privacy to sanitize points in the parameter space of machine learning models, which are then communicated to a central server for aggregation in order to converge to the optimal parameters and, thus, obtain the personalized models for the diverse datasets. Given that the distribution of the data among individuals is unknown, it is reasonable to assume a mixture of multiple distributions. Clustering the sanitized parameter vectors released by the clients with the  $k$ -means algorithm shows to be a good proxy for aggregating clients with similar data distri-

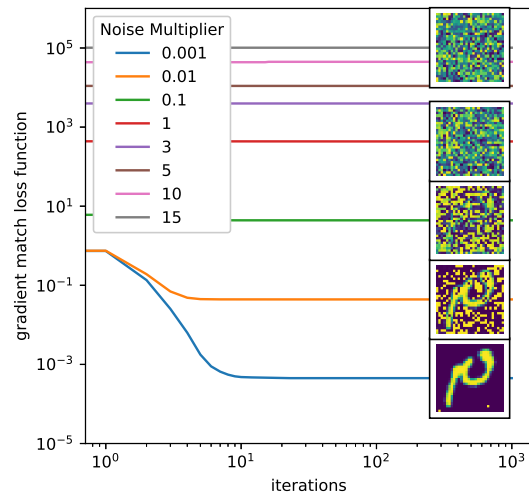


Figure 5: Effects of the Laplace mechanism in Proposition 3.1 with different noise multipliers as a defense strategy against the DLG attack.

butions. This is possible because  $d$ -private mechanisms preserve the topology of the domain of true values. Our mechanism shows to be promising when machine learning models have a *small* number of parameters. Although formal privacy guarantees degrade sharply with large machine learning models, we show

experimentally that the Laplace mechanism is effective against the DLG attack. As future work, we want to explore other privacy mechanisms, which may be more effective in providing a good trade-off between privacy and accuracy in the context of machine learning. Furthermore, we are interested in studying more complex federated learning scenarios where participants and datasets may change over time.

## ACKNOWLEDGEMENTS

The work of Sayan Biswas and Catuscia Palamidessi was supported by the European Research Council (ERC) project HYPATIA under the European Union’s Horizon research and innovation programme, grant agreement n. 835294. The work of Kangsoo Jung was supported by ELSA - The European Lighthouse on Secure and Safe AI, a Network of Excellence funded by the European Union under the Horizon research and innovation programme.

## REFERENCES

- Abadi, M., Chu, A., Goodfellow, I., McMahan, H. B., Mironov, I., Talwar, K., and Zhang, L. (2016). Deep learning with differential privacy. In *Proceedings of the 2016 ACM SIGSAC conference on computer and communications security*, pages 308–318.
- Agarwal, N., Suresh, A. T., Yu, F. X. X., Kumar, S., and McMahan, B. (2018). cpsgd: Communication-efficient and differentially-private distributed SGD. *Advances in Neural Information Processing Systems*, 31.
- Andrés, M. E., Bordenabe, N. E., Chatzikokolakis, K., and Palamidessi, C. (2013). Geo-indistinguishability: Differential privacy for location-based systems. In *Proceedings of the 2013 ACM SIGSAC conference on Computer & communications security*, pages 901–914.
- Andrew, G., Thakkar, O., McMahan, B., and Ramaswamy, S. (2021). Differentially private learning with adaptive clipping. *Advances in Neural Information Processing Systems*, 34.
- Bassily, R., Nissim, K., Stemmer, U., and Guha Thakurta, A. (2017). Practical locally private heavy hitters. In Guyon, I., Luxburg, U. V., Bengio, S., Wallach, H., Fergus, R., Vishwanathan, S., and Garnett, R., editors, *Advances in Neural Information Processing Systems*, volume 30. Curran Associates, Inc.
- Bonawitz, K., Ivanov, V., Kreuter, B., Marcedone, A., McMahan, H. B., Patel, S., Ramage, D., Segal, A., and Seth, K. (2016). Practical secure aggregation for federated learning on user-held data. *arXiv preprint arXiv:1611.04482*.
- Caldas, S., Duddu, S. M. K., Wu, P., Li, T., Konečný, J., McMahan, H. B., Smith, V., and Talwalkar, A. (2018). Leaf: A benchmark for federated settings. *arXiv preprint arXiv:1812.01097*.
- Chatzikokolakis, K., Andrés, M. E., Bordenabe, N. E., and Palamidessi, C. (2013). Broadening the scope of differential privacy using metrics. In *International Symposium on Privacy Enhancing Technologies Symposium*, pages 82–102. Springer.
- CMMS (2021). Centers for medicare and medicaid services. Accessed: 2022-09-21.
- Cohen, G., Afshar, S., Tapson, J., and Van Schaik, A. (2017). Emnist: Extending mnist to handwritten letters. In *2017 international joint conference on neural networks (IJCNN)*, pages 2921–2926. IEEE.
- Dwork, C., Kenthapadi, K., McSherry, F., Mironov, I., and Naor, M. (2006a). Our data, ourselves: Privacy via distributed noise generation. In Vaudenay, S., editor, *Advances in Cryptology - EUROCRYPT 2006*, pages 486–503, Berlin, Heidelberg. Springer Berlin Heidelberg.
- Dwork, C., McSherry, F., Nissim, K., and Smith, A. (2006b). Calibrating noise to sensitivity in private data analysis. In Halevi, S. and Rabin, T., editors, *Theory of Cryptography*, pages 265–284, Berlin, Heidelberg. Springer Berlin Heidelberg.
- Geyer, R. C., Klein, T., and Nabi, M. (2017). Differentially private federated learning: A client level perspective. *arXiv preprint arXiv:1712.07557*.
- Ghosh, A., Chung, J., Yin, D., and Ramchandran, K. (2020). An efficient framework for clustered federated learning. *Advances in Neural Information Processing Systems*, 33:19586–19597.
- Goodfellow, I. J., Mirza, M., Xiao, D., Courville, A., and Bengio, Y. (2013). An empirical investigation of catastrophic forgetting in gradient-based neural networks. *arXiv preprint arXiv:1312.6211*.
- Hitaj, B., Ateniese, G., and Perez-Cruz, F. (2017). Deep models under the gan: information leakage from collaborative deep learning. In *Proceedings of the 2017 ACM SIGSAC conference on computer and communications security*, pages 603–618.
- Hu, R., Guo, Y., Li, H., Pei, Q., and Gong, Y. (2020). Personalized federated learning with differential privacy. *IEEE Internet of Things Journal*, 7(10):9530–9539.
- Kirkpatrick, J., Pascanu, R., Rabinowitz, N., Veness, J., Desjardins, G., Rusu, A. A., Milan, K., Quan, J., Ramalho, T., Grabska-Barwinska, A., et al. (2017). Overcoming catastrophic forgetting in neural networks. *Proceedings of the national academy of sciences*, 114(13):3521–3526.
- Konečný, J., McMahan, H. B., Yu, F. X., Richtárik, P., Suresh, A. T., and Bacon, D. (2016). Federated learning: Strategies for improving communication efficiency. *arXiv preprint arXiv:1610.05492*.
- Le Métayer, D. and De, S. J. (2016). PRIAM: a Privacy Risk Analysis Methodology. In Livraga, G., Torra, V., Aldini, A., Martinelli, F., and Suri, N., editors, *Data Privacy Management and Security Assurance*, Heraklion, Greece. Springer.
- Lopez-Paz, D. and Ranzato, M. (2017). Gradient episodic



memory for continual learning. *Advances in neural information processing systems*, 30.

- Mansour, Y., Mohri, M., Ro, J., and Suresh, A. T. (2020). Three approaches for personalization with applications to federated learning. *arXiv preprint arXiv:2002.10619*.
- McMahan, B., Moore, E., Ramage, D., Hampson, S., and y Arcas, B. A. (2017a). Communication-efficient learning of deep networks from decentralized data. In *Artificial intelligence and statistics*, pages 1273–1282. PMLR.
- McMahan, H. B., Ramage, D., Talwar, K., and Zhang, L. (2017b). Learning differentially private recurrent language models. *arXiv preprint arXiv:1710.06963*.
- Nasr, M., Shokri, R., and Houmansadr, A. (2019). Comprehensive privacy analysis of deep learning: Passive and active white-box inference attacks against centralized and federated learning. In *2019 IEEE symposium on security and privacy (SP)*, pages 739–753. IEEE.
- NIST (2021). Nist privacy framework core.
- Sattler, F., Müller, K.-R., and Samek, W. (2020). Clustered federated learning: Model-agnostic distributed multitask optimization under privacy constraints. *IEEE transactions on neural networks and learning systems*, 32(8):3710–3722.
- Truex, S., Liu, L., Chow, K.-H., Gursoy, M. E., and Wei, W. (2020). Ldp-fed: Federated learning with local differential privacy. In *Proceedings of the Third ACM International Workshop on Edge Systems, Analytics and Networking*, pages 61–66.
- Zhang, Z. and Sabuncu, M. (2018). Generalized cross entropy loss for training deep neural networks with noisy labels. *Advances in neural information processing systems*, 31.
- Zhao, Y., Zhao, J., Yang, M., Wang, T., Wang, N., Lyu, L., Niyato, D., and Lam, K.-Y. (2020). Local differential privacy-based federated learning for internet of things. *IEEE Internet of Things Journal*, 8(11):8836–8853.
- Zhu, L., Liu, Z., and Han, S. (2019). Deep leakage from gradients. *Advances in Neural Information Processing Systems*, 32.

## APPENDIX

### Proofs

**Proposition 3.1.** Let  $\mathcal{L}_\epsilon: \mathbb{R}^n \mapsto \mathbb{R}^n$  be the Laplace mechanism with distribution  $\mathcal{L}_{x_0, \epsilon}(x) = \mathbb{P}[\mathcal{L}_\epsilon(x_0) = x] = Ke^{-\epsilon d(x, x_0)}$  with  $d(\cdot)$  being the Euclidean distance. If  $\rho \sim \mathcal{L}_{x_0, \epsilon}(x)$ , then:

1.  $\mathcal{L}_{x_0, \epsilon}$  is  $\epsilon$ - $d$ -private and  $K = \frac{\epsilon^n \Gamma(\frac{n}{2})}{2\pi^{\frac{n}{2}} \Gamma(n)}$
2.  $\|\rho\|_2 \sim \gamma_{\epsilon, n}(r) = \frac{\epsilon^n e^{-\epsilon r} r^{n-1}}{\Gamma(n)}$
3. The  $i^{\text{th}}$  component of  $\rho$  has variance  $\sigma_{\rho_i}^2 = \frac{n+1}{\epsilon^2}$

where  $\Gamma(n)$  is the Gamma function defined for positive reals as  $\int_0^\infty t^{n-1} e^{-t} dt$  which reduces to the factorial function whenever  $n \in \mathbb{N}$ .

*Proof.* We provide proof of the three statements separately:

1. If  $\mathcal{L}_{x_0, \epsilon}(x) = Ke^{-\epsilon d(x, x_0)}$  is a probability density function of a point  $x \in \mathbb{R}^n$  then  $K$  should be such that  $\int_{\mathbb{R}^n} \mathcal{L}_{x_0}(x) dx = 1$ . We note that it depends only on the distance  $x$  and  $x_0$  and we can write  $Ke^{-\epsilon d(x, x_0)}$  as  $Ke^{-\epsilon r}$  where  $r$  is the radius of the ball in  $\mathbb{R}^n$  centered in  $x_0$ . Without loss of generality, let us now take  $x_0 = 0$ . The probability density of the event  $x \in \mathbb{S}_n(r) = \{x : \|x\|_2 = r\}$  is then  $p(x \in \mathbb{S}_n(r)) = Ke^{-\epsilon r} S_n(1) r^{n-1}$  where  $S_n(1)$  is the surface of the unitary ball in  $\mathbb{R}^n$  and  $S_n(r) = S_n(1) r^{n-1}$  is the surface of a generic ball of radius  $r$ . Given that

$$S_n(1) = \frac{2\pi^{n/2}}{\Gamma(\frac{n}{2})} \quad (7)$$

solving

$$\int_0^{+\infty} \mathbb{P}[x \in \mathbb{S}_n(r)] dr = \int_0^{+\infty} Ke^{-\epsilon r} S_n(1) r^{n-1} dr = K \frac{2\pi^{n/2} \Gamma(n)}{\epsilon^n \Gamma(\frac{n}{2})} = 1 \quad (8)$$

results in

$$K = \frac{\epsilon^n \Gamma(\frac{n}{2})}{2\pi^{\frac{n}{2}} \Gamma(n)} \quad (9)$$

where  $\Gamma(\cdot)$  denotes the gamma function. By plugging  $\mathcal{L}_{x_0, \epsilon}(x) = Ke^{-\epsilon d(x, x_0)}$  in Equation 5:

$$Ke^{-\epsilon d(x, x_1)} \leq e^{\epsilon d(x_1, x_2)} Ke^{-\epsilon d(x, x_2)} \quad (10)$$

$$e^{\epsilon(\|x-x_2\|_2 - \|x-x_1\|_2)} \leq e^{\epsilon\|x_1-x_2\|} = e^{\epsilon d(x_1, x_2)} \quad (11)$$

which completes the poof of the first statement.

2. Without loss of generality, let us take  $x_0 = 0$ . Exploiting the radial symmetry of the Laplace distribution, we note that, in order to sample a point  $\rho \sim \mathcal{L}_{x_0, \epsilon}(x)$  in  $\mathbb{R}^n$ , it is possible to first sample the set of points distant  $d(x, 0) = r$  from  $x_0$  and then sample uniformly from the resulting hypersphere. Accordingly, the p.d.f. of the  $L_2$ -norm of  $\rho$  is the p.d.f. of the event  $\rho \in \mathbb{S}_n(r) = \{\rho : \|\rho\|_2 = r\}$  which is then  $\mathbb{P}[\rho \in \mathbb{S}_n(r)] = Ke^{-\epsilon r} S_n(1) r^{n-1}$ , where  $\mathbb{S}_n(r)$  is the surface of the sphere with radius  $r$  in  $\mathbb{R}^n$ . Hence, we can write

$$\|\rho\|_2 \sim \gamma_{\epsilon, n}(r) = \frac{\epsilon^n e^{-\epsilon r} r^{n-1}}{\Gamma(n)} \quad (12)$$

which completes the proof of the second statement.

3. With  $\rho \sim \gamma_{\epsilon, n}$  we have that, by construction,

$$\mathbb{E} [\rho^2] = \mathbb{E} \left[ \sum_{i=1}^n \rho_i^2 \right] = n \mathbb{E} [\rho_i^2] = n \sigma_{\rho_i}^2 \quad (13)$$

With the last equality holding since  $\mathcal{L}_{0, \epsilon}$  is isotropic and centered in zero. Recalling that

$$\mathbb{E} [\rho^2] = \frac{d^2}{dt^2} M_{\rho}(t) \Big|_{t=0} \quad (14)$$

with  $M_{\rho}(t)$  the moment generating function of the gamma distribution  $\gamma_{\epsilon, n}$ ,

$$\begin{aligned} & \frac{d^2}{dt^2} \left( \left(1 - \frac{t}{\epsilon}\right)^{-n} \right) \Big|_{t=0} = \\ & = \frac{n(n+1)}{\epsilon^2} \left(1 - \frac{t}{\epsilon}\right)^{-(n+2)} \Big|_{t=0} = \\ & = \frac{n(n+1)}{\epsilon^2} \end{aligned}$$

which leads to  $\sigma_{\rho_i}^2 = \frac{n+1}{\epsilon^2}$ , completing the proof of the third statement and of the Proposition.  $\square$

*Proof.* The Root Mean Square Error loss function is defined as:

$$F = \sqrt{\frac{\sum_{i=1}^{|Z|} (f(x_i, \theta) - y_i)^2}{|Z|}} = \frac{\|f(X, \theta) - Y\|_2}{\sqrt{|Z|}} \quad (15)$$

If the model parameters  $\theta$  are sanitized by the addition of a random vector  $\rho \sim \mathcal{L}_{0, \epsilon}$ , we can evaluate how the cost function would change with respect to the non-sanitized parameters. Dropping the multiplicative constant we find:

$$\begin{aligned} & \|f(X, \theta + \rho) - Y\|_2 - \|f(X, \theta) - Y\|_2 \leq \\ & \|f(X, \theta + \rho) - Y - f(X, \theta) + Y\|_2 = \\ & \|f(X, \theta + \rho) - f(X, \theta)\|_2 = \\ & \|f(X, \theta) + J_f(X, \theta) \cdot \rho - f(X, \theta) + o(J_f(X, \theta) \cdot \rho)\|_2 = \\ & \|J_f(X, \theta) \cdot \rho + o(J_f(X, \theta) \cdot \rho)\|_2 \leq \\ & \|J_f(X, \theta)\|_2 \|\rho\|_2 + o(\|J_f(X, \theta) \cdot \rho\|_2) \end{aligned}$$

$\square$

with  $J_f(X, \theta)$  being the Jacobian of  $f$  with respect to  $X$  and  $o(\cdot)$  being higher terms coming from the Taylor expansion. Thus we proved that the bound on the increase of the cost function does not depend on the direction of the additive noise, but on its norm, in first order approximation.

**Theorem 3.1.** *Let  $\mathcal{K}_i$  be  $(\epsilon_i)$ - $d$ -private mechanism for  $i \in \{1, 2\}$ . Then their independent composition is  $(\epsilon_1 + \epsilon_2)$ - $d$ -private.*

*Proof.* Let us simplify the notation and denote:

$$P_i = \mathbb{P}_{\mathcal{K}_i} [y_i \in S_i | x_i]$$

$$P'_i = \mathbb{P}_{\mathcal{K}'_i} [y_i \in S_i | x'_i]$$

for  $i \in \{1, 2\}$ . As mechanisms  $\mathcal{K}_1$  and  $\mathcal{K}_2$  are applied independently, we have:

$$\mathbb{P}_{\mathcal{K}_1, \mathcal{K}_2} [(y_1, y_2) \in S_1 \times S_2 | (x_1, x_2)] = P_1 \cdot P_2$$

$$\mathbb{P}_{\mathcal{K}'_1, \mathcal{K}'_2} [(y_1, y_2) \in S_1 \times S_2 | (x'_1, x'_2)] = P'_1 \cdot P'_2$$

Therefore, we obtain:

$$\begin{aligned} & \mathbb{P}_{\mathcal{K}_1, \mathcal{K}_2} [(y_1, y_2) \in S_1 \times S_2 | (x_1, x_2)] = P_1 \cdot P_2 \\ & \leq \left( e^{\epsilon_1 d(x_1, x'_1)} P'_1 \right) \left( e^{\epsilon_2 d(x_2, x'_2)} P'_2 \right) \\ & \leq e^{\epsilon_1 d(x_1, x'_1) + \epsilon_2 d(x_2, x'_2)} \mathbb{P}_{\mathcal{K}'_1, \mathcal{K}'_2} [(y_1, y_2) \in S_1 \times S_2 | (x'_1, x'_2)] \end{aligned} \quad \square$$

## Experimental Settings

All the following experiments are run on a local server running Ubuntu 20.04.3 LTS with an AMD EPYC 7282 16-Core processor, 1.5TB of RAM and NVIDIA A100 GPUs. Python and PyTorch are the main software tools adopted for simulating the federation of clients and their corresponding collaborative training.

## Synthetic Data

A total of 100 users holding 10 samples each, drawn from either one of the distributions, participate in a training of two initial hypotheses which are sampled from a Gaussian distribution centered in 0 and unit variance at iteration  $t = 0$ . A total of  $U = 7$  users are asked to participate in the optimization at each round and train locally the hypothesis that fits better their dataset for  $E = 1$  epochs each time. The noise multiplier is set to  $\nu = 5$ . Local step size  $s = 0.1$  and a batch size  $B_s = 10$  complete the required inputs to the algorithm. To verify the training process, another set of users with the same characteristics is held out from training to perform validation and stop the federated optimization once there is no improvement in the loss function in Equation (15) for 6 consecutive rounds. Although at first the updates seem to be distributed all over the domain, in just a few rounds of training the process converges to values very close to the two optimal parameters. With the heuristic presented in Section 3.2 it is easy to find that whenever a user participates in an optimization round it incurs in a privacy leakage of at most  $n/\nu = 2/5 = 0.4$ , in a differential private sense, with respect to points in its neighborhood. Using the result in Theorem 3.1 clients can

compute the overall privacy leakage of the optimization process, should they be required to participate multiple times. For any user, whether to participate or not in a training round can be decided right before releasing the updated parameters, in case that would increase the privacy leakage above a threshold value decided beforehand.

### Hospital Charge Data

The dataset contains details about charges for the 100 most common inpatient services and the 30 most common outpatient services. It shows a great variety of charges applied by healthcare providers with details mostly related to the type of service and the location of the provider. Preprocessing of the dataset includes a number of procedures, the most important of which are described here:

- i) Selection of the 4 most widely treated conditions, which amount to simple pneumonia; kidney and urinary tract infections; hart failure and shock; esophagitis and digestive system disorders.
- ii) Transformation of ZIP codes into numerical coordinates in terms of longitude and latitude.
- iii) Setting as target the Average Total Payments, i.e. the cost of the service averaged among the times it was given by a certain provider.
- iv) As it is a standard procedure in the context of gradient-based optimization, dependent and independent variables are brought to be in the range of the *units* before being fed to the machine learning model. Note that this point takes the spot of the common feature normalization and standardization procedures, which we decided not to perform here to keep the setting as realistic as possible. In fact, both would require the knowledge of the empirical distribution of all the data. Although it is available in simulation, it would not be available in a real scenario, as each user would only have access to their dataset.

Given the preprocessing described above, the dataset results in 2947 clients, randomly split in train and validation subsets with 70 and 30 per cent of the total clients each. The goal is being able to predict the cost that a service would require given where it is performed in the country, and what kind of procedure it is. The model that was adopted in this context is a fully connected neural network (NN) of two layers, with a total of 11 parameters and Rectified Linear Unit (ReLU) activation function. Inputs to the model are an increasing index which uniquely defines the healthcare service, the longitude and latitude of the provider. Output of the model is the expected

cost. Tests have been performed to minimize the RMSE loss on the clients selected for training (100 per round) and at each round the performance of the model is checked against a held-out set of validation clients, from where 200 are sampled every time. If 30 validation rounds are passed without improvement in the cost function, the optimization process is terminated. In order to decrease the variability of the results, a total of 10 runs have been performed with different seeds for every combination of number of hypotheses and noise multiplier.

Table 3: NN architecture adopted in the experiments of Section 4.3.

Layer	Properties
2D Convolution	kernel size: (2,2) stride: (1,1) nonlinearity: ReLU output features: 32
2D Convolution	kernel size: (2,2) stride: (1,1) nonlinearity: ReLU output features: 64
2D Max Pool	kernel size: (2,2) stride: (2,2) nonlinearity: ReLU
Fully Connected	nonlinearity: ReLU units: 128
Fully Connected	nonlinearity: ReLU units: 62

### FEMNIST Image Classification

The task consists in performing image classification on the FEMNIST (Caldas et al., 2018) dataset, which is a standard benchmark dataset for federated learning, based on EMNIST (Cohen et al., 2017) and with the data points grouped by user. It consists of a large number of images of handwritten digits, lower and upper case letters of the Latin alphabet. As a preprocessing step, images of client  $c$  are rotated 90 degrees counter-clockwise depending on the realization of the random variable  $\text{rot}_c \sim \text{Bernoulli}(0.5)$ . This is a common practice in machine learning to simulate local datasets held by different clients being generated by different distributions (Ghosh et al., 2020; Goodfellow et al., 2013; Kirkpatrick et al., 2017; Lopez-Paz and Ranzato, 2017). The chosen architecture is described in Table 3 and yields a parameter vector  $\theta \in \mathbb{R}^{n_0}$ ,  $n_0 = 1206590$ . Runs are performed with a maximum of 500 rounds of federated optimization, unless 5 consecutive validation rounds are conducted without improvements on the validation loss. The latter is evaluated on a held out set of clients, consisting

of 10% of the total number. Validation is performed every 5 training rounds, thus the process terminates after 25 rounds without the model's performance improvement. The optimization process aims to minimize either the RMSE loss or the Cross Entropy loss (Zhang and Sabuncu, 2018) between model's predictions and the target class.