

Quantifying the Impact of Label Noise on Federated Learning

Shuqi Ke¹ Chao Huang² Xin Liu²

¹ The Chinese University of Hong Kong, Shenzhen ² University of California, Davis
shuqi@link.cuhk.edu.cn, fchhuang@ucdavis.edu, xinliu@ucdavis.edu

Abstract

Federated Learning (FL) is a distributed machine learning paradigm where clients collaboratively train a model using their local (human-generated) datasets. While existing studies focus on FL algorithm development to tackle data heterogeneity across clients, the important issue of data quality (e.g., label noise) in FL is overlooked. This paper aims to fill this gap by providing a quantitative study on the impact of label noise on FL. We derive an upper bound for the generalization error that is *linear* in the clients' label noise level. Then we conduct experiments on MNIST and CIFAR-10 datasets using various FL algorithms. Our empirical results show that the global model accuracy linearly decreases as the noise level increases, which is consistent with our theoretical analysis. We further find that label noise slows down the convergence of FL training, and the global model tends to overfit when the noise level is high.

Introduction

Federated Learning (FL) is a distributed machine learning paradigm where clients (e.g., distributed devices or organizations) collaboratively train a global model (Kairouz et al. 2021). The local data of the clients are often human-generated and have critical privacy concerns. An FL process consists of some communication rounds. In each round, each client trains its local model with its local data and then uploads the model updates to a central server (Hsu, Qi, and Brown 2019). The central server aggregates the local updates from clients and sends back an aggregated global model to all clients. After that, clients update their local models according to the information from the central server (Collins et al. 2022). The client-server interaction stops when the global model converges.

There has been an increasing volume of research studies on FL over the last few years (Kairouz et al. 2021; Wang et al. 2021; Li et al. 2021b; Jiang et al. 2022). Among these studies, a critical bottleneck, which without appropriate algorithmic treatment usually fails FL, is data heterogeneity (non-IID). For example, in a classification task, some clients may collect more data for class *A* while others may collect more data for class *B*. Previous studies among this line focused on two categories of non-IID: attribute skew and label skew (Zhu et al. 2021). Attribute skew refers to the case

where the feature distribution of each client is different from one another. For example, attribute skew could occur in a handwritten digit classification task as users may write the same digit with different font styles, sizes, and stroke widths (Kairouz et al. 2021). Label skew refers to the case where the label distribution of each client is different from one another. Label skew, for example, could occur in an animal recognition task. Label distributions are different because clients are in different geo-regions and different animal habitats — dolphins only live near coastal regions, or aquariums (Kairouz et al. 2021).

While existing studies focus on tackling the non-IIDness, some implicitly assume that the data are clean, i.e., the data are correctly labeled. In practical applications, however, clients' datasets usually contain noisy labels (Northcutt, Athalye, and Mueller 2021). Label noise has been identified in many widely used FL datasets, including MNIST (Northcutt, Jiang, and Chuang 2021; LeCun and Cortes 2005), EMNIST (Al-Rawi and Karatzas 2018; Cohen et al. 2017), CIFAR-10 (Krizhevsky 2009; Al-Rawi and Karatzas 2018), ImageNet (Northcutt, Jiang, and Chuang 2021; Russakovsky et al. 2014), and Clothing1M (Xiao et al. 2015). The causes of label noise can be human error, subjective labeling tasks, non-exact data labeling processes, and malfunctioning data collection infrastructure (Johnson and Khoshgoftaar 2022; Chen et al. 2021). Moreover, in an FL setting, as clients collect and label local data in a distributed and private fashion, their labels are likely to be noisy and have different noise patterns (Xu et al. 2022). For example, wearable devices can access various human-generated data, such as heart rate, sleep patterns, medication records, and mental health logs. Such data could contain different levels of label noise due to various sensor precision issues and human bias (Kim, Jo, and Choi 2021).

Label noise is known to lessen model performance (Johnson and Khoshgoftaar 2022). This paper focuses on the issue of label noise in FL, and we are particularly interested in answering the following two key questions:

- **Question 1:** *How does label noise affect FL convergence?*
- **Question 2:** *How does label noise affect FL generalization?*

To answer Question 1, we conduct numerical experiments

and show that the training loss converges slower with a higher noise level. To answer Question 2, we proceed from both theoretical and empirical perspectives. First, under minor assumptions, we prove that, for any distributed learning algorithm, the generalization error of the global model is linearly bounded above by a multiple of the system noise level. Then we conduct experiments using MNIST and CIFAR-10, showing that the results are consistent with the assumptions and theoretical results. We further show that the global model’s accuracy decreases linearly in the clients’ label noise level.

The key contributions of this paper are summarized below.

- To the best of our knowledge, this is the first quantitative study that analyzes the impact of label noise on FL. Our study bears practical significance for its use in different applications, e.g., incentive design (Huang et al. 2022).
- We provide a generic upper bound on the FL generalization error that applies to any FL algorithms. We further obtain a tighter upper bound considering the widely adopted ReLU networks in clients’ local models.
- We run experiments under various algorithms and different settings in FL. Our numerical results justify our theoretic assumption. We also observe that label noise linearly degrades FL performance by reducing the test accuracy of the global model.
- Our study reveals several important and interesting insights. (1) Label noise slows down FL convergence; (2) label noise induces overfitting to the global model; (3) among three benchmark FL algorithms, SCAFFOLD (Karimireddy et al. 2020) achieves the best test accuracy than other algorithms with minor label noise, while FedNova (Wang et al. 2020) achieves the best test accuracy with more extensive label noise.

Related Work

Label noise

Label noise has been an active topic in FL over the last few years. We classify the existing methods into three categories: (1) Some methods apply *noise-tolerant loss functions* to achieve robust performance (Sharma et al. 2022).

(2) Some methods *distill confident training sample by selection or a weighting scheme* (Chen et al. 2021; Yang et al. 2022b,a; Ma et al. 2021; Chen et al. 2020; Fang and Ye 2022; Duan et al. 2022; Han and Zhang 2020; Li et al. 2021a; Kim et al. 2022; Li, Pei, and Huang 2022; Tuor et al. 2021). Li et al. discovered that label noise might cause overfitting for FedAvg algorithm. However, they did not analytically characterize the hidden linear relation between noise level and the global model’s performance.

(3) Based on (2), some methods further *correct noisy samples* (Xu et al. 2022; Zeng et al. 2022; Wang et al. 2022; Tsouvalas et al. 2022). Tsouvalas et al. proposed FedLN that estimates per-client noise level and corrects noisy labels. However, their definition of label noise is limited because they only considered the engineering method to generate label noise as the definition of label noise. They considered a

case where conditional distributions¹ $\Pr(\text{label}|\text{feature})$ are the same across clients (Kairouz et al. 2021). But in practice, the conditional distributions could be different for different clients. We provide a more general definition in this work and fill this gap. Xu et al. studied an FL scenario where different clients have different levels of label noise (Xu et al. 2022). They introduced local intrinsic dimension (LID), a measure of the dimension of the data manifold. They discovered a strong linear relation between cumulative LID score and local noise level. However, their work did not provide either empirical observation or theoretical results on the relation between the global model’s performance and local noise level. Moreover, there is no systematic study on how label noise affects FL in terms of convergence and generalization. We bridge this research gap in this work.

Path-norm

This work uses path-norm to measure the global model’s generalization ability under label noise. People introduced different measures to explain the generalization ability of neural networks (Zheng et al. 2019; Jiang* et al. 2020). Behnam Neyshabur et al. proposed path-norm as a capacity measure for ReLU networks (Neyshabur, Salakhutdinov, and Srebro 2015; Neyshabur et al. 2017). Empirical studies showed that path-norm positively correlates with generalization in all categories of hyper-parameter (Jiang* et al. 2020).

The value of path-norm increases throughout the learning process. E et al. showed that the path-norm increases at most polynomially under centralized training (Weinan et al. 2020). In this work, we conduct the first formal study on the evolution of path-norm in FL. This is also the first work that analyzes the generalization ability of models in FL with path-norm proxy. We introduce path-norm proxy to the FL context because this proxy does not require unrealistic assumptions and allows us to characterize a large class of FL algorithms. For example, the assumptions on convexity, smoothness, etc., are no longer necessary in our analysis. Moreover, we have empirically verified our theory based on the definition of path-norm proxy.

Preliminaries and Problem Statement

Federated Learning

In this subsection, we briefly introduce the problem formulation and algorithmic framework of FL.

Consider a typical FL task (Kairouz et al. 2021), where N clients collaboratively train a global model under the coordination of a central server through R communication rounds. FL aims to solve a distributed optimization problem with distributed data. Here we first introduce the objective of the distributed optimization problem and then define the relevant notations. The objective is

$$\min_{W \in \mathcal{W}} \frac{1}{N} \sum_{k=1}^N \left[\frac{1}{n_k} \sum_{i=1}^{n_k} \ell(f(x_{k,i}; W), y_{k,i}) \right] \quad (1)$$

where we define

¹In some work, the conditional distribution is also referred to as “feature-to-label mapping”.

- Hypothesis space: $\mathcal{W} \subset \mathbb{R}^{d_w}$ denotes the hypothesis space of all feasible parameters of learning models, and $d_w \in \mathbb{N}$ is the dimension of the hypothesis space.
- Local data: Each client has a local dataset S_k . We assume that in the k -th dataset S_k , each data point is drawn from a distribution π_k over $\mathcal{S} \subset \mathbb{R}^{d_x+d_y}$ where d_x denotes the dimension of feature space and d_y denotes the dimension of the label space. A data point $(x, y) \in \mathbb{R}^{d_x+d_y}$ is a real-valued vector where $x \in \mathbb{R}^{d_x}$ denotes its feature and $y \in \mathbb{R}^{d_y}$ denotes its label. There are in total n_k data points in client k 's local dataset

$$S_k = \{(x_{k,1}, y_{k,1}), (x_{k,2}, y_{k,2}), \dots, (x_{k,n_k}, y_{k,n_k})\}$$

Let μ_k denote the ground truth distribution (i.e., clean labels) and π_k denote client k 's possibly noisy data distribution. There exists label noise in the local dataset of client k if there exists $x \in \mathbb{R}^{d_x}, y \in \mathbb{R}^{d_y}$ such that

$$\Pr_{\mu_k}(y|x) \neq \Pr_{\pi_k}(y|x) \quad (2)$$

where \Pr represents a probability mass/density function with a given distribution and an event. One can consider the data points sampled from π_k as training data and those sampled from μ_k as test data.

- Global parameter and local parameter: We denote the global model's parameter as a real-valued vector $W \in \mathcal{W}$. Each client has a local model with parameter $w_k \in \mathcal{W}$.
- Meta model: We define the meta model $f : \mathbb{R}^{d_x} \times \mathcal{W} \rightarrow \mathbb{R}^{d_y}$ as a function that maps the data feature and model parameter to an estimated label. For example, a meta model could be a neural network with variable parameters. We obtain a model by substituting the variable parameters with real number values.
- Loss function: We denote the loss function as

$$\ell : \mathbb{R}^{d_y} \times \mathbb{R}^{d_y} \rightarrow \mathbb{R}_{\geq 0}$$

For example, a squared loss function is defined as $\ell : (y, \hat{y}) \mapsto \|y - \hat{y}\|^2$.

In each communication round, a client trains its local model for E epochs to minimize the local training loss $\frac{1}{n_k} \sum_{i=1}^{n_k} \ell(f(x_{k,i}; W), y_{k,i})$ over its local dataset S_k . After local model training, the clients upload their local model parameters w_k to a central server. The central server aggregates the uploaded parameters and updates the global model's parameter W . After that, the central server sends the global model's new parameter back to each client. We provide a general FL framework in Algorithm 1.

Different FL algorithms use different aggregation mechanisms. We use FedAvg as an example to explain the aggregation step in Algorithm 1. In FedAvg, the aggregation is defined as

$$\phi : (w_1, \dots, w_N, W) \mapsto (1 - \eta_{gl})W + \eta_{gl} \frac{\sum_{k=1}^N w_k}{N} \quad (3)$$

where η_{gl} denotes the global learning rate. Note that in a realistic setting, there could be limitations on computation and

Algorithm 1: A General FL Framework

Initialization: Local datasets $\{S_1, S_2, \dots, S_N\}$, aggregation function ϕ

Output: Global model parameter vector W and local model parameter vectors $\{w_1, w_2, \dots, w_N\}$ after the R -th communication round

```

1: for t ← 1 to R do
2:   Parallel for k ← 1 to N do
3:     for i ← 1 to E do                                     ▷ local training
4:       Update local model parameter  $w_k$ 
5:     end for
6:     Send  $w_k$  to the central server
7:   end for
8:    $W \leftarrow \phi(w_1, \dots, w_N, W)$                        ▷ aggregation
9:   for k ← 1 to N do                                       ▷ broadcast
10:    Send  $W$  to client  $k$ 
11:    Update local model parameter  $w_k$  according to
       $W$ 
12:   end for
13: end for

```

communication, including computational efficiency, communication bandwidth, and network robustness (Ghosh et al. 2020). For example, some clients may fail to communicate with the central server due to network issues. Therefore, the server only samples a subset of available clients. Since we focus on data noise, we ignore these realistic considerations and assume that all clients participate in all communication rounds.

Model performance

This subsection introduces the theoretical tools to measure a learning algorithm's performance. Here we inherit most notations from the last part with some revisions. We consider fixed data points for a FL process in the previous part. But in this part, we consider each data point and each local dataset S_k as random variables to investigate the generalization performance of an algorithm given an arbitrary training dataset. The pair (x, y) in lowercase represents a deterministic data point, and the pair (X, Y) in uppercase represents pair of random variables. We re-write a local dataset S_k as

$$S_k = \{(X_{k,1}, Y_{k,1}), (X_{k,2}, Y_{k,2}), \dots, (X_{k,n_k}, Y_{k,n_k})\}$$

where $(X_{k,i}, Y_{k,i}) \sim \pi_k$. We define the empirical risk $L : \mathcal{W} \rightarrow \mathbb{R}_{\geq 0}$ of the global model as

$$L(W) = \sum_{k=1}^N \frac{n_k}{n} \mathbb{E}_{\pi_k} [\ell(f(X; W), Y)] \quad (4)$$

where $n := \sum_{k=1}^N n_k$ and W denotes the parameter of the global model. Given the ground truth distribution μ_k of each client, we further define the ground-truth risk $L^\dagger : \mathcal{W} \rightarrow \mathbb{R}_{\geq 0}$ of the global model as

$$L^\dagger(W) = \sum_{k=1}^N \frac{n_k}{n} \mathbb{E}_{\mu_k} [\ell(f(X; W), Y)] \quad (5)$$

Then we define the generalization error of the global model as (Yagli, Dytso, and Vincent Poor 2020)

$$G(W) := |L^\dagger(W) - L(W)| \quad (6)$$

Path-norm proxy

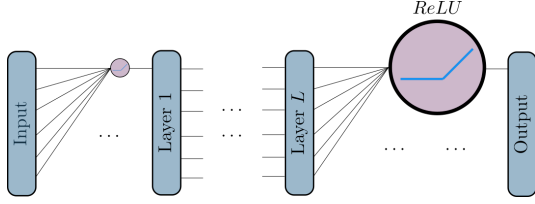


Figure 1: An L -layer ReLU network.

This paper uses ReLU network and path-norm proxy for a case study of the generalization error. The authors in (Weinan et al. 2020) provided a mathematical description of ReLU networks. Based on their definitions, we define the path-norm proxy below.

Definition 1 (Path-norm proxy (Weinan et al. 2020)). *The path-norm proxy of an L -layer ReLU network is defined as*

$$\|f(\cdot; \theta)\|_{\text{pnp}} = \sum_{(i_0, \dots, i_{L+1})} \prod_{l=0}^L |\theta_l(i_l, i_{l+1})| \quad (7)$$

where θ denotes the parameter vector of the ReLU network; $\theta_l(i_l, i_{l+1})$ refers to the weight of the edge connecting the i_l -th node in layer l and the i_{l+1} -th node in layer $l+1$.

The authors in (Weinan et al. 2020) also proved that the path norm proxy controls the generalization error in a centralized learning setting. Next, we will show that the path norm proxy controls the generalization error in FL.

Theoretical Results

In this section, we provide a theoretical analysis of the generalization error of the global model in FL. In particular, we give proof of the upper bound of the global model's generalization error.

In practical FL applications, local data distributions are complicated as we cannot explicitly find the distribution functions. To simplify our theoretical analysis, we make the following assumption:

Assumption 2 (Simplified label noise condition). *For any client i and client j , we assume*

$$\forall(x, y) \in \mathbb{R}^{d_x+d_y}, \Pr(x; \pi_i) = \Pr(x; \pi_j) \quad (8)$$

This assumption assures the feature distributions to be identical for all clients, which is a standard setting in studies about concept drift (Jothimurugesan et al. 2022). Although it is difficult to show that Assumption 2 holds in our experiment settings, the numerical results are still consistent with our theoretical results.

We first provide a general result on the upper bound of generalization error in Theorem 3. Then we extend this general bound by studying some specific cases with more assumptions in Corollary 8.

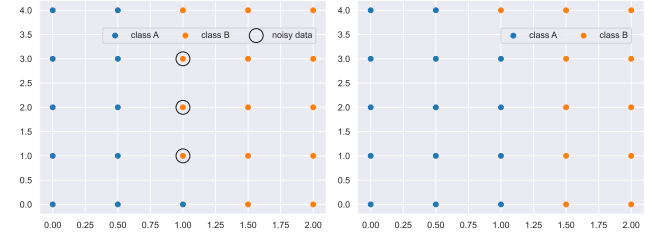
Theorem 3 (Bound the evolution of generalization error). *Consider any Federated Learning algorithm with a neural network with an arbitrary structure for a classification task of C classes under label noise and use the cross-entropy function for loss computation, then under Assumption 2*

$$G(W) \leq \Omega \cdot \mathbb{E}_X \left[\sum_{i=1}^C \sum_{k=1}^N \frac{n_k}{n} |\Pr_\mu(Y = i|X) - \Pr_{\pi_k}(Y = i|X)| \right] \quad (9)$$

where Ω is the upper bound of f .

Interpretation of Theorem 3: This theorem implies that the generalization error of global model in FL is linearly bounded by the degree of label noise in the distributed system. The theorem quantitatively characterizes the impact of label noise. This linear bound is also consistent with our empirical findings. When $N = 1$, this linear bound applies to centralized learning.

We can interpret the expectation term in the upper bound with an example. In this example, we set $N = 2$, i.e., two clients. The input space consists of 25 discrete grid points and two classes. Client 2's local data distribution is identical to the ground truth. Client 1 has label noise in its local data where three circled data points in class A are mislabelled as class B.



(a) Data distribution of client 1 (b) Data distribution of client 2

Figure 2: An example of label noise.

If the two clients has the same number of data samples, i.e., $n_1 = n_2$, then

$$\begin{aligned} & \mathbb{E}_X \left[\sum_{i=1}^C \sum_{k=1}^N \frac{n_k}{n} |\Pr_\mu(Y = i|X) - \Pr_{\pi_k}(Y = i|X)| \right] \\ &= \mathbb{E}_X \left[\sum_{i=1}^C \frac{1}{2} |\Pr_\mu(Y = i|X) - \Pr_{\pi_1}(Y = i|X)| \right] \\ &= \frac{1}{2} \left(\frac{1}{5} \cdot \left| \frac{1}{5} - \frac{4}{5} \right| + \frac{1}{5} \left| \frac{4}{5} - \frac{1}{5} \right| \right) = \frac{3}{25} \end{aligned} \quad (10)$$

This expectation represents the expected **percentage of noisy data points** in a dataset, e.g. there are in total 25 grid points and 3 noisy data points in Fig 2a.

Now we consider a slightly different example where client 2 also has label noise as shown in Figure 3. Then the expectation is $\frac{4}{25}$.

Before we prove Theorem 3, we need a lemma on cross-entropy.

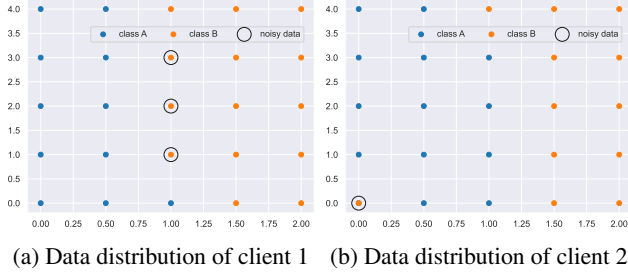


Figure 3: An example of label noise.

Lemma 4. Consider a classification problem of C classes. Given a data distribution π such that $(x, y) \sim \pi, y \in [1 : C]$, a neural network f and a probability measure \Pr , then the expectation of cross-entropy loss is

$$-\sum_{i=1}^C \Pr_{\pi}(Y = i) \mathbb{E}_{X|Y=i} \left[f_i(X) - \log \left(\sum_{r=1}^C \exp(f_r(X)) \right) \right] \quad (11)$$

In most machine learning tasks, it is reasonable to assume that the input and output of the model are bounded, which we formalize in Assumptions 5 and 6.

Assumption 5 (Bounded input space). The input space \mathcal{X} is bounded in $[0, 1]^{d_x} \subset \mathbb{R}^{d_x}$.

Assumption 6 (Bounded model output). Consider a neural network $f : \mathbb{R}^{d_x} \times \mathcal{W} \rightarrow \mathbb{R}^{d_y}$. We assume that its range $f(\mathbb{R}^{d_x}; \mathcal{W})$ is bounded in \mathbb{R}^{d_y} . That is, $\exists C_f \geq 0$ such that $\forall x \in \mathbb{R}^{d_x}, \forall \theta \in \mathcal{W}, \forall i \in \{1, \dots, C\}, C_f \geq |f_i(x; \theta)|$.

Note that the upper bound of model output could change as we train the model for more epochs. To model the evolution of the output upper bound, we can relax Assumption 6 and study a specific family of classifiers: ReLU networks. Later we can bound the generalization error evolution given the growth of path-norm proxy through iterations.

Proposition 7 (Polynomial growth of path-norm proxy). Consider an FL process with a L -layer neural network $f : \mathbb{R}^{d_x} \times \mathcal{W} \rightarrow \mathbb{R}^{d_y}$ as its global model, assume that its path-norm increases at most polynomially,

$$\|f(\cdot; \theta(t))\|_{\text{pnp}} = \mathcal{O}(t^{L+1} E^{(L+1)/2}) \quad (12)$$

where $t \leq R$ denotes the number of communication rounds and E denotes the local training time.

If we consider a generic decentralized algorithm, we have

$$\|f(\cdot; \theta(t))\|_{\text{pnp}} = \mathcal{O}(e^{C't(L+1)} E^{(L+1)/2}) \quad (13)$$

where C' is a constant independent of t, L, E .

Corollary 8. We can specify Ω in Theorem 3 with various assumptions:

1. By Assumption 6, $\Omega = C_f$.
2. If we use a ReLU network as our model in the FL task, then $\Omega = \|f(\cdot; \theta(t))\|_{\text{pnp}}$.

3. By Assumption 5 and Proposition 7,

$$\Omega = C_0 t^{L+1} E^{(L+1)/2}$$

where C_0 is a constant independent of t, E, L .

There are some important implications behind Corollary 8.

- Since the first two statements of the corollary do not rely on the aggregation mechanism of the algorithm, they could also be extended from FL to a decentralized learning scenario, e.g. Swarm learning in decentralized clinical ML (Warnat-Herresthal et al. 2021), decentralized optimization algorithms (Zhang, Ahmad, and Wang 2019; Luo and Ye 2022), ML on blockchain (Liu et al. 2020).
- Theorem 3 does not characterize the upper bound with communication rounds and local epochs in its general form. But it is a symbolic and concise term that helps us understand the impact of label noise. Nonetheless, case 3 in Corollary 8 provides the interplay between the label noise, communication rounds, and local epochs.

Numerical Results

We present three numerical experiments to validate our theoretical results and draw new insights. We first verify our theoretical work on the path-norm proxy. Then we show experiments of 2-client, 4-client, and 15-client FL settings.

Our main findings are 1) the growth of path-norm proxy empirically increases in a **polynomial** order in FL; 2) there exists an approximate negative **linear** relation between the test accuracy of global model and the number of incorrectly labeled data; 3) label noise slows down the **convergence** of FL algorithms and induces **over-fitting** to the global model.

Path-norm Proxy

In this subsection, we study the path norm proxy and observe its relation with the number of layers and communication rounds.

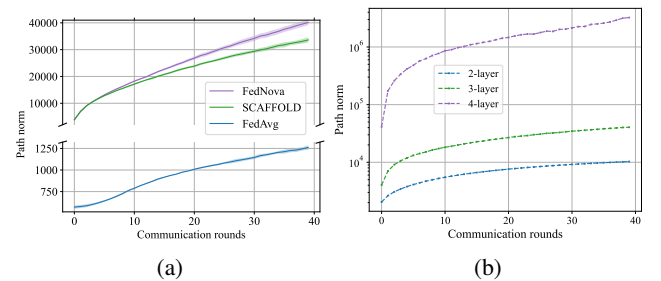


Figure 4: A case study on MNIST dataset. (a) The path-norm generated by different FL algorithms with 3-layer ReLU network. (b) The path-norm generated by FedAvg algorithm and ReLU networks with different numbers of layers.

Compare different FL algorithms. We use the same neural network structure and consider $N = 4$ clients. We study three FL algorithms, FedAvg (McMahan et al. 2017), SCAFFOLD, and FedNova on MNIST dataset. We observe a concave growth of the global model’s path-norm in Figure 4a.

We can sort the three algorithms according to the scale of their path-norm values

Path norm value: FedNova > SCAFFOLD >> FedAvg

which corresponds to their different usage of local gradient grad and local parameter change Δw

Algorithm	Aggregation input
FedNova	grad
SCAFFOLD	grad, Δw
FedAvg	Δw

These results empirically show that the path norm increases polynomially in FL.

Compare ReLU networks with different numbers of layers. We train ReLU networks using FedAvg on MNIST. We use different numbers of layers in $\{2, 3, 4\}$. The result is illustrated in Figure 4b. To verify the polynomial rule, we use a logarithmic scale on the y-axis. The three path-norm curves have similar shapes and almost differ up to a constant factor. This result is consistent with Proposition 7.

Pilot experiments

We run 2-client experiments with FedAvg algorithm on MNIST dataset. We study a 2-client setting for multiple considerations.

- 2-client setting exists in practice. In cross-silo FL, clients could be enterprises, and each client could provide abundant data, so the total number of clients is relatively small. For example, since 2019, two insurance companies, Swiss Re and WeBank, have collaborated on federated learning (Huang, Huang, and Liu 2022).
- This experiment serves as a starting point and gives us a thorough pedagogical understanding of the impact of label noise. We will study the 4-client and 15-client cases in the next subsection.

We generate the local datasets for two clients by dividing the whole dataset into two equally-sized parts. We add label noise to local datasets by uniformly flipping some instances' labels to other class labels. Each client has different noise levels. Denote the noise level of client i as wp_i , then $(wp_1, wp_2) \in \{0\%, 10\%, 20\%, \dots, 80\%, 90\%\}^2$. Pathological noise levels (greater than 50%) have been studied in supervised learning settings (Luo et al. 2022). We illustrate the test accuracy of global models under different degrees of label noise as bar charts in Fig 5. To verify the linear trend of test accuracy, we perform linear regression and visualize the result in Fig 6.

Negative bilinear trend by label noise. Figure 5 shows a negative bilinear relation between the test accuracy of the global model and noise label. When we apply linear regression on the test accuracy of the global model and the proportion of wrongly labeled data, we obtain a coefficient of determination of 0.98 in Figure 6. That means the relation between the test accuracy and label noise has a strong linear relation.

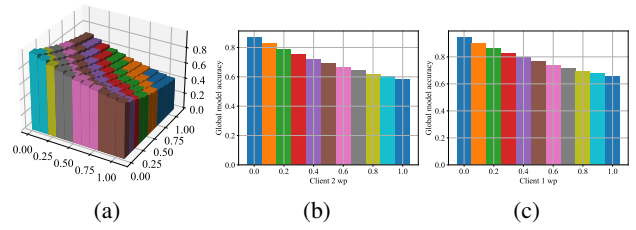


Figure 5: (a) Bar plot of global model accuracy. x,y axes control the levels of label noise of each client. z axis represents the test accuracy of the global model; (b) Slice of bar plot when client 1 has 30% of wrongly labelled data; (c) Slice of bar plot when client 2 has 10% of wrongly labelled data.

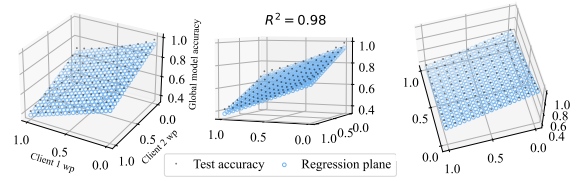


Figure 6: Linear regression on the global model accuracy.

Experiments with larger cohort size

We run experiments on CIFAR-10 dataset respectively with four clients and fifteen clients. Local datasets are generated by dividing the whole dataset into equally-sized parts. We add label noise to local datasets by uniformly flipping some instances' labels to other class labels. In a case study by Gu et al., the real human annotation has a rater error rate of around 4.8% (Gu et al. 2022). Therefore it is reasonable to study the error rate within a relatively small range that contains 4.8%, i.e., from 0% to 10%. We set the same proportion of wrongly labeled data for each client (0%, 2%, 4%, 8%).

Slow Convergence by label noise. In Figure 7 and Figure 8, we plot how client 1's local model loss depends on the communication rounds at different percentages of wrongly labeled data. The training loss decreases slower with a larger proportion of wrongly labeled data, i.e., the algorithm converges slower with a larger proportion of wrongly labeled data.

Overfitting by label noise. We observe in Figure 9 that for all three algorithms, the global model's test accuracy de-

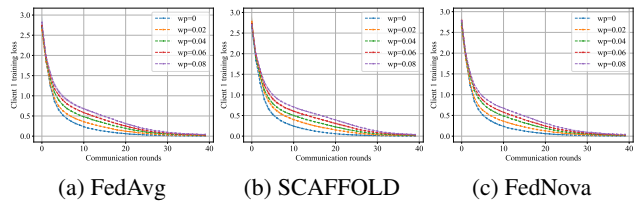


Figure 7: Training loss of Client 1 for 0%, 2%, 4%, 6%, 8% percentages of wrongly labelled data (4-client setting).

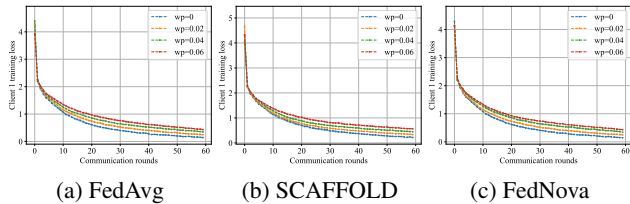


Figure 8: Training loss of Client 1 for 0%, 2%, 4%, 6%, 8% percentages of wrongly labelled data (15-client setting).

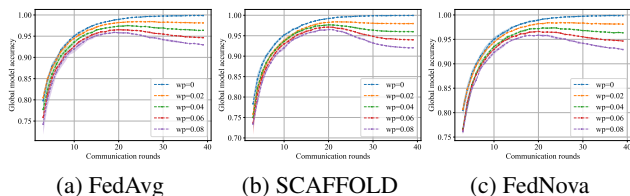


Figure 9: Test accuracy of global model for 0%, 2%, 4%, 6%, 8% percentages of wrongly labelled data (4-client setting).

creases after 20 communication rounds. The global model is more over-fitted with a larger percentage of wrongly labeled data. This result provides an engineering insight in FL that the over-fitting of the global model could result from some wrongly labeled data in the local datasets. It also motivates the study of mitigating label noise in FL (Li et al. 2021a).

Negative linear trend by label noise. In Figure 10, all three algorithms show a negative linear relation between the test accuracy of the global model and the proportion of wrongly labeled data. This is consistent with our theoretical analysis.

Discussions

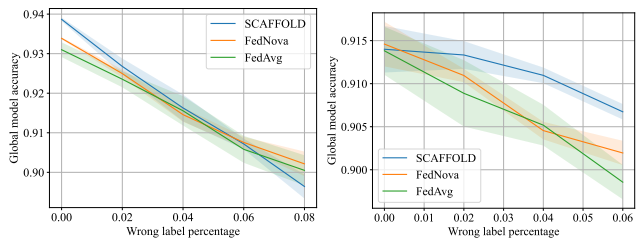
In this section, we discuss the limitation and potential application of our work.

Improving theoretical bounds: We prove a linear upper bound for the generalization error, and the bound is consistent with numerical results. However, the upper bound can be loose. One can provide a lower bound or improve the upper bound by making more restrictive assumptions. For example, one can consider a regression task with MSE loss function that provides nicer theoretical properties (Damian, Ma, and Lee 2021).

More comprehensive experiments: Our experiments use a small number of clients, which applies to cross-silo FL. In future research, we plan to study the impact of label noise with a larger number of clients (e.g., as in cross-device FL).

Application: Our results potentially serve as “domain knowledge” to improve FL algorithm design. Our work could also be used in designing incentive mechanisms in FL systems (Huang et al. 2022). In particular, the qualitative relation in this paper helps model the performance of global model under label noise.

Methodology: The emergence of large machine learning



(a) 4-client (b) 15-client

Figure 10: Test accuracy of global model by different FL algorithms under different label error rates.

models has shifted the nature of AI research from an engineering science (iteratively improving models) to a natural science (probing capabilities of the models we designed) (Kambhampati 2022). Researchers have been proposing hundreds of new models/algorithms for different AI problems. However, more must be done to understand how and why a proposed model/algorithm performs in a certain way. We must build theories based on observation and experiments to understand these artificial black boxes. In this way, we can transform AI research from engineering alchemy to white-box chemistry. Our work follows this scientific paradigm shift and conducts a case study for different FL models under label noise.

Concluding Remarks

This paper takes the first step to quantify the impact of label noise on the global model in FL. The critical challenge is that we have little knowledge of the underlying information related to local data distributions and we do not have an explicit expression of the outcome of an FL algorithm. We show with both empirical evidence and theoretical proof that 1) label noise linearly degrades the global model’s performance in FL; 2) label noise slows down the convergence of the global model; 3) label noise induces overfitting to the global model. Our results could provide insights into the design of a noise-robust algorithm and the design of an incentive mechanism.

Acknowledgments

The work was partially supported through grant US-DA/NIFA 2020-67021-32855, and by NSF through IIS-1838207, CNS 1901218, OIA-2134901.

References

- Al-Rawi, M. S.; and Karatzas, D. 2018. On the Labeling Correctness in Computer Vision Datasets. In *IAL@PKDD/ECML*.
- Chen, C.; Zheng, S.; Chen, X.; Dong, E.; Liu, X. S.; Liu, H.; and Dou, D. 2021. Generalized DataWeighting via Class-Level Gradient Manipulation. In Ranzato, M.; Beygelzimer, A.; Dauphin, Y.; Liang, P.; and Vaughan, J. W., eds., *Advances in Neural Information Processing Systems*, volume 34, 14097–14109. Curran Associates, Inc.

- Chen, Y.; Yang, X.; Qin, X.; Yu, H.; Chan, P.; and Shen, Z. 2020. *Dealing with Label Quality Disparity in Federated Learning*, 108–121. Cham: Springer International Publishing. ISBN 978-3-030-63076-8.
- 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)*, 2921–2926.
- Collins, L.; Hassani, H.; Mokhtari, A.; and Shakkottai, S. 2022. FedAvg with Fine Tuning: Local Updates Lead to Representation Learning.
- Damian, A.; Ma, T.; and Lee, J. D. 2021. Label Noise SGD Provably Prefers Flat Global Minimizers. In *NeurIPS*.
- Duan, S.; Liu, C.; Cao, Z.; Jin, X.; and Han, P. 2022. Fed-DR-Filter: Using Global Data Representation to Reduce the Impact of Noisy Labels on the Performance of Federated Learning. *Future Gener. Comput. Syst.*, 137(C): 336–348.
- Fang, X.; and Ye, M. 2022. Robust Federated Learning with Noisy and Heterogeneous Clients. In *2022 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 10062–10071.
- Ghosh, A.; Chung, J.; Yin, D.; and Ramchandran, K. 2020. An Efficient Framework for Clustered Federated Learning. In Larochelle, H.; Ranzato, M.; Hadsell, R.; Balcan, M.; and Lin, H., eds., *Advances in Neural Information Processing Systems*, volume 33, 19586–19597. Curran Associates, Inc.
- Gu, K.; Masotto, X.; Bachani, V.; Lakshminarayanan, B.; Nikodem, J.; and Yin, D. 2022. An instance-dependent simulation framework for learning with label noise. *Machine Learning*.
- Han, Y.; and Zhang, X. 2020. Robust Federated Learning via Collaborative Machine Teaching. *Proceedings of the AAAI Conference on Artificial Intelligence*, 34(04): 4075–4082.
- Hsu, T. H.; Qi, H.; and Brown, M. 2019. Measuring the Effects of Non-Identical Data Distribution for Federated Visual Classification. *CoRR*, abs/1909.06335.
- Huang, C.; Huang, J.; and Liu, X. 2022. Cross-Silo Federated Learning: Challenges and Opportunities.
- Huang, C.; Ke, S.; Kamhoua, C.; Mohapatra, P.; and Liu, X. 2022. Incentivizing Data Contribution in Cross-Silo Federated Learning.
- Jiang*, Y.; Neyshabur*, B.; Mobahi, H.; Krishnan, D.; and Bengio, S. 2020. Fantastic Generalization Measures and Where to Find Them. In *International Conference on Learning Representations*.
- Jiang, Z.; Wang, W.; Li, B.; and Yang, Q. 2022. Towards Efficient Synchronous Federated Training: A Survey on System Optimization Strategies. *IEEE Transactions on Big Data*, 1–1.
- Johnson, J. M.; and Khoshgoftaar, T. M. 2022. A Survey on Classifying Big Data with Label Noise. *J. Data and Information Quality*. Just Accepted.
- Jothimurugesan, E.; Hsieh, K.; Wang, J.; Joshi, G.; and Gibbons, P. 2022. Federated Learning under Distributed Concept Drift. In *NeurIPS 2022 Workshop on Distribution Shifts: Connecting Methods and Applications*.
- Kairouz, P.; McMahan, H. B.; Avent, B.; Bellet, A.; Bennis, M.; Bhagoji, A. N.; Bonawitz, K.; Charles, Z.; Cormode, G.; Cummings, R.; D’Oliveira, R. G. L.; Eichner, H.; Rouayheb, S. E.; Evans, D.; Gardner, J.; Garrett, Z.; Gascón, A.; Ghazi, B.; Gibbons, P. B.; Gruteser, M.; Harchaoui, Z.; He, C.; He, L.; Huo, Z.; Hutchinson, B.; Hsu, J.; Jaggi, M.; Javidi, T.; Joshi, G.; Khodak, M.; Konečný, J.; Korolova, A.; Koushanfar, F.; Koyejo, S.; Lepoint, T.; Liu, Y.; Mittal, P.; Mohri, M.; Nock, R.; Özgür, A.; Pagh, R.; Raykova, M.; Qi, H.; Ramage, D.; Raskar, R.; Song, D.; Song, W.; Stich, S. U.; Sun, Z.; Suresh, A. T.; Tramèr, F.; Vepakomma, P.; Wang, J.; Xiong, L.; Xu, Z.; Yang, Q.; Yu, F. X.; Yu, H.; and Zhao, S. 2021. *Advances and Open Problems in Federated Learning. Foundations and Trends in Machine Learning*.
- Kambhampati, S. 2022. Changing the Nature of AI Research. *Commun. ACM*, 65(9): 8–9.
- Karimireddy, S. P.; Kale, S.; Mohri, M.; Reddi, S.; Stich, S.; and Suresh, A. T. 2020. SCAFFOLD: Stochastic Controlled Averaging for Federated Learning. In III, H. D.; and Singh, A., eds., *Proceedings of the 37th International Conference on Machine Learning Research*, 5132–5143. PMLR.
- Kim, B. H.; Jo, S.; and Choi, S. 2021. ALIS: Learning Affective Causality Behind Daily Activities From a Wearable Life-Log System. *IEEE Transactions on Cybernetics*, 1–13.
- Kim, S.; Shin, W.; Jang, S.; Song, H.; and Yun, S.-Y. 2022. FedRN. In *Proceedings of the 31st ACM International Conference on Information & Knowledge Management*. ACM.
- Krizhevsky, A. 2009. Learning Multiple Layers of Features from Tiny Images.
- LeCun, Y.; and Cortes, C. 2005. The mnist database of handwritten digits.
- Li, J.; Pei, J.; and Huang, H. 2022. Communication-Efficient Robust Federated Learning with Noisy Labels. In *Proceedings of the 28th ACM SIGKDD Conference on Knowledge Discovery and Data Mining, KDD ’22*, 914–924. New York, NY, USA: Association for Computing Machinery. ISBN 9781450393850.
- Li, L.; Gao, L.; Fu, H.; Han, B.; Xu, C.-Z.; and Shao, L. 2021a. Federated Noisy Client Learning.
- Li, Q.; Wen, Z.; Wu, Z.; Hu, S.; Wang, N.; Li, Y.; Liu, X.; and He, B. 2021b. A Survey on Federated Learning Systems: Vision, Hype and Reality for Data Privacy and Protection. *IEEE Transactions on Knowledge and Data Engineering*, 1–1.
- Liu, J.; and Morse, A. S. 2011. Accelerated linear iterations for distributed averaging. *Annual Reviews in Control*, 35(2): 160–165.
- Liu, Y.; Yu, F. R.; Li, X.; Ji, H.; and Leung, V. C. M. 2020. Blockchain and Machine Learning for Communications and Networking Systems. *IEEE Communications Surveys & Tutorials*, 22(2): 1392–1431.
- Luo, L.; and Ye, H. 2022. Decentralized Stochastic Variance Reduced Extragradient Method.
- Luo, Y.; Liu, G.; Guo, Y.; and Yang, G. 2022. Deep Neural Networks Learn Meta-Structures from Noisy Labels in Semantic Segmentation. In *AAAI*.

- Ma, J.; Sun, X.; Xia, W.; Wang, X.; Chen, X.; and Zhu, H. 2021. Client Selection Based on Label Quantity Information for Federated Learning. In *2021 IEEE 32nd Annual International Symposium on Personal, Indoor and Mobile Radio Communications (PIMRC)*, 1–6.
- McMahan, H. B.; Moore, E.; Ramage, D.; Hampson, S.; and y Arcas, B. A. 2017. Communication-Efficient Learning of Deep Networks from Decentralized Data. In *AISTATS*.
- Neyshabur, B.; Bhojanapalli, S.; Mcallester, D.; and Srebro, N. 2017. Exploring Generalization in Deep Learning. In Guyon, I.; Luxburg, U. V.; Bengio, S.; Wallach, H.; Fergus, R.; Vishwanathan, S.; and Garnett, R., eds., *Advances in Neural Information Processing Systems*, volume 30. Curran Associates, Inc.
- Neyshabur, B.; Salakhutdinov, R. R.; and Srebro, N. 2015. Path-SGD: Path-Normalized Optimization in Deep Neural Networks. In Cortes, C.; Lawrence, N.; Lee, D.; Sugiyama, M.; and Garnett, R., eds., *Advances in Neural Information Processing Systems*, volume 28. Curran Associates, Inc.
- Northcutt, C.; Jiang, L.; and Chuang, I. 2021. Confident Learning: Estimating Uncertainty in Dataset Labels. *J. Artif. Int. Res.*, 70: 1373–1411.
- Northcutt, C. G.; Athalye, A.; and Mueller, J. 2021. Pervasive Label Errors in Test Sets Destabilize Machine Learning Benchmarks. In *Thirty-fifth Conference on Neural Information Processing Systems Datasets and Benchmarks Track (Round 1)*.
- Russakovsky, O.; Deng, J.; Su, H.; Krause, J.; Satheesh, S.; Ma, S.; Huang, Z.; Karpathy, A.; Khosla, A.; Bernstein, M.; Berg, A.; and Fei-Fei, L. 2014. ImageNet Large Scale Visual Recognition Challenge. *International Journal of Computer Vision*, 115.
- Sharma, R.; Ramakrishna, A.; MacLaughlin, A.; Rumshisky, A.; Majumdar, J.; Chung, C.; Avestimehr, S.; and Gupta, R. 2022. Federated Learning with Noisy User Feedback. In *Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, 2726–2739. Seattle, United States: Association for Computational Linguistics.
- Tsouvalas, V.; Saeed, A.; Ozelebi, T.; and Meratnia, N. 2022. Federated Learning with Noisy Labels.
- Tuor, T.; Wang, S.; Ko, B.; Liu, C.; and Leung, K. K. 2021. Overcoming Noisy and Irrelevant Data in Federated Learning. *2020 25th International Conference on Pattern Recognition (ICPR)*, 5020–5027.
- Wang, J.; Charles, Z. B.; Xu, Z.; Joshi, G.; McMahan, H. B.; y Arcas, B. A.; Al-Shedivat, M.; Andrew, G.; Avestimehr, S.; Daly, K.; Data, D.; Diggavi, S. N.; Eichner, H.; Gadhikar, A.; Garrett, Z.; Girgis, A. M.; Hanzely, F.; Hard, A.; He, C.; Horváth, S.; Huo, Z.; Ingerman, A.; Jaggi, M.; Javid, T.; Kairouz, P.; Kale, S.; Karimireddy, S. P.; Konečný, J.; Koyejo, S.; Li, T.; Liu, L.; Mohri, M.; Qi, H.; Reddi, S. J.; Richtárik, P.; Singhal, K.; Smith, V.; Soltanolkotabi, M.; Song, W.; Suresh, A. T.; Stich, S. U.; Talwalkar, A. S.; Wang, H.; Woodworth, B. E.; Wu, S.; Yu, F. X.; Yuan, H.; Zaheer, M.; Zhang, M.; Zhang, T.; Zheng, C.; Zhu, C.; and Zhu, W. 2021. A Field Guide to Federated Optimization. *ArXiv*, abs/2107.06917.
- Wang, J.; Liu, Q.; Liang, H.; Joshi, G.; and Poor, H. V. 2020. Tackling the Objective Inconsistency Problem in Heterogeneous Federated Optimization. In *Proceedings of the 34th International Conference on Neural Information Processing Systems, NIPS’20*. Red Hook, NY, USA: Curran Associates Inc. ISBN 9781713829546.
- Wang, Z.; Zhou, T.; Long, G.; Han, B.; and Jiang, J. 2022. FedNoiL: A Simple Two-Level Sampling Method for Federated Learning with Noisy Labels.
- Warnat-Herresthal, S.; Schultze, H.; Shastry, K.; Manamohan, S.; Mukherjee, S.; Garg, V.; Sarveswara, R.; Händler, K.; Pickkers, P.; Aziz, N. A.; Ktena, S.; Tran, F.; Bitzer, M.; Ossowski, S.; Casadei, N.; Herr, C.; Petersheim, D.; Behrends, U.; Kern, F.; and Velavan, T. 2021. Swarm Learning for decentralized and confidential clinical machine learning. *Nature*, 594.
- Weinan, E.; ; 9212; ; E, W.; Stephan; Wojtowysch; ; 9213; ; and Wojtowysch, S. 2020. On the Banach Spaces Associated with Multi-Layer ReLU Networks: Function Representation, Approximation Theory and Gradient Descent Dynamics. *SIAM Transactions on Applied Mathematics*, 1(3): 387–440.
- Xiao, T.; Xia, T.; Yang, Y.; Huang, C.; and Wang, X. 2015. Learning from massive noisy labeled data for image classification. In *2015 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2691–2699.
- Xu, J.; Chen, Z.; Quek, T. Q.; and Chong, K. F. E. 2022. FedCorr: Multi-Stage Federated Learning for Label Noise Correction. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*.
- Yagli, S.; Dytso, A.; and Vincent Poor, H. 2020. Information-Theoretic Bounds on the Generalization Error and Privacy Leakage in Federated Learning. In *2020 IEEE 21st International Workshop on Signal Processing Advances in Wireless Communications (SPAWC)*, 1–5.
- Yang, M.; Qian, H.; Wang, X.; Zhou, Y.; and Zhu, H. 2022a. Client Selection for Federated Learning With Label Noise. *IEEE Transactions on Vehicular Technology*, 71(2): 2193–2197.
- Yang, S.; Park, H.; Byun, J.; and Kim, C. 2022b. Robust Federated Learning With Noisy Labels. *IEEE Intelligent Systems*, 37(2): 35–43.
- Zeng, B.; Yang, X.; Chen, Y.; Yu, H.; and Zhang, Y. 2022. CLC: A Consensus-Based Label Correction Approach in Federated Learning. *ACM Trans. Intell. Syst. Technol.*, 13(5).
- Zhang, C.; Ahmad, M.; and Wang, Y. 2019. ADMM Based Privacy-Preserving Decentralized Optimization. *IEEE Transactions on Information Forensics and Security*, 14(3): 565–580.
- Zheng, S.; Meng, Q.; Zhang, H.; Chen, W.; Yu, N.; and Liu, T.-Y. 2019. Capacity Control of ReLU Neural Networks by Basis-Path Norm. In *Proceedings of the Thirty-Third AAAI Conference on Artificial Intelligence and Thirty-First Innovative Applications of Artificial Intelligence Conference and*

Zhu, H.; Xu, J.; Liu, S.; and Jin, Y. 2021. Federated Learning on Non-IID Data: A Survey. *Neurocomput.*, 465(C): 371–390.

Appendix I

In Federated Learning, the “non-IID” issue is defined as the statistical difference or statistical dependence of different local datasets from different clients (Kairouz et al. 2021; Zhu et al. 2021). In this work, we consider that for different clients i and j , the distributions π_i, π_j of their local datasets are different

$$\pi_i \neq \pi_j$$

There are further two major types of non-IID: feature distribution skew and label distribution skew.

- For label distribution skew, all clients share the conditional probability $\pi_i(x|y)$, that is,

$$\pi_i(x|y) = \pi_j(x|y), \forall i, j \in [1, 2, \dots, N], \forall (x, y) \in \mathbb{R}^{d_x+d_y}$$

while clients have different label distributions $\pi_i(y)$.

- For feature distribution skew, clients share the same conditional probability $\pi_i(y|x)$ while clients have different feature distributions $\pi_i(x)$.

Appendix II

The ReLU network is a powerful prototype model among various types of neural network for its successful performance in different fields, including image classification and natural language processing (Zheng et al. 2019). In this work, we study ReLU network as a sub-case in generalization error analysis.

We represent a L -layer neural network as a map $f : \theta \mapsto f(\cdot; \theta)$ where θ denotes the weight of the network and $f(\cdot; \theta) : x \mapsto f(x; \theta)$ is a function that maps an input $x \in \mathbb{R}^{d_x}$ to an output $y \in \mathbb{R}^{d_y}$ of the network.

Denote the width of the l -th layer as d_l where $1 \leq l \leq L$ and let $d_0 = d_x + 1, d_{L+1} = d_y$, i.e. there are d_l nodes in the l -th layer. Here we simplify the notation by converting the affine map to a linear one: identify $x \in \mathbb{R}^{d_x}$ with $(x, 1) \in \mathbb{R}^{d_x+1}$ (Weinan et al. 2020).

Definition 9 (Rectified linear unit). *The rectified linear unit function is defined to be*

$$\sigma : \mathbb{R} \rightarrow \mathbb{R}_{\geq 0}, \quad \sigma(x) = \max(0, x) \quad (14)$$

Definition 10 (Rectifier activation function). *By an abuse of notation, we define the rectifier activation function by applying the rectified linear unit function element-wise*

$$\sigma : \mathbb{R}^d \rightarrow \mathbb{R}_{\geq 0}^d, \quad \sigma \left(\begin{bmatrix} x_1 \\ \vdots \\ x_d \end{bmatrix} \right) = \begin{bmatrix} \max(0, x_1) \\ \vdots \\ \max(0, x_d) \end{bmatrix} \quad (15)$$

The k -th entry of output $f(x; \theta)$ of a L -layer ReLU network is (Zheng et al. 2019; Weinan et al. 2020)

$$\begin{aligned} f_k(x; \theta) &= \sum_{i_L=1}^{d_L} \theta_L(i_L, k) \sigma \left(\sum_{i_{L-1}=1}^{d_{L-1}} \theta(i_{L-1}, i_L) \sigma \left(\sum_{i_{L-2}=1}^{d_{L-2}} \cdots \sigma \left(\sum_{i_1=1}^{d_1} \theta_1(i_1, i_2) \sigma \left(\sum_{i_0=1}^{d_0} \theta_0(i_0, i_1) x_{i_0} \right) \right) \right) \right) \\ &= \sum_{(i_0, \dots, i_L)} \theta_L(i_L, k) \cdot \prod_{l=1}^L \theta_l(i_{l-1}, i_l) \cdot \prod_{l=1}^L \mathbf{1}\{g_{i_l}^l(\theta, x) > 0\} \cdot x_{i_0} \end{aligned} \quad (16)$$

where parameter $\theta_l(i_l, i_{l+1})$ refers to the weight of the edge connecting the i_l -th node in layer l and the i_{l+1} -th node in layer $l+1$; $g_{i_l}^l$ denotes the output of i_l -th node in layer l .

Appendix III

Definition 11 (Cross-entropy loss). *Given a neural network f , an input vector x , and the output vector $y = \text{Softmax}(f(x))$, we define the cross-entropy loss as*

$$-\sum_i y_i \log \left(\frac{\exp(f_i(x))}{\sum_j \exp(f_j(x))} \right) \quad (17)$$

where subscript i denotes the i -th entry of a vector.

Proof of Theorem 3. Without loss of generality, we prove the theorem given that we use a neural network as our classifier. By Lemma 4, we expand the formula of generalization error as follows

$$\begin{aligned} G(W) &= |L(W) - L^\dagger(W)| \\ &= \left| \sum_{k=1}^N \frac{n_k}{n} \sum_{i=1}^C \int_{\mathcal{X}} f_i(x) (d\Pr_{\pi_k}(x, y) - d\Pr_{\mu}(x, y)) \right| \\ &\text{by Assumption 2} \\ &= \left| \sum_{k=1}^N \frac{n_k}{n} \mathbb{E}_X \left[\sum_{i=1}^C f_i(X) (\Pr_{\pi_k}(Y = i|X) - \Pr_{\mu}(Y = i|X)) \right] \right| \\ &\leq \sum_{k=1}^N \frac{n_k}{n} \mathbb{E}_X \left[\sum_{i=1}^C f_i(X) |\Pr_{\pi_k}(Y = i|X) - \Pr_{\mu}(Y = i|X)| \right] \\ &\leq \sum_{k=1}^N \frac{n_k}{n} \Omega \cdot \mathbb{E}_X \left[\sum_{i=1}^C |\Pr_{\pi_k}(y = i|X) - \Pr_{\mu}(y = i|X)| \right] \end{aligned} \quad (18)$$

□

Theorem 12 (Growth of path-norm proxy (Weinan et al. 2020, Corollary 5)). *Consider an arbitrarily wide L -layer ReLU neural network. If the network’s weight evolves under continuous gradient flow dynamics, then the network’s path-norm increases at most polynomially*

$$\|f(\cdot; \theta(t))\|_{\text{pnp}} \leq \left(C_0 + \sqrt{\mathcal{R}(f(\cdot; \theta(t=0)))} t^{1/2} \right)^{L+1} \quad (19)$$

where $\theta(t)$ denotes the weight of the network at time t , C_0 is a constant and \mathcal{R} is the expectation of a sufficiently smooth loss function ℓ

$$\mathcal{R}(\theta; f) = \int_{\mathbb{R}^{d_x+d_y}} \ell(f(x; \theta), y) \mathbb{P}(dx \otimes dy) \quad (20)$$

Proof sketch of Proposition 7. This proof assumes the gradient-flow evolution and an arbitrarily wide neural network based on the proofs in (Weinan et al. 2020).

We first consider an FL setting. Let $\theta^{(i)}(t_k)$ denote the collection of all the parameters of client i 's neural network uploaded to the central server at the k -th communication round before FL aggregation. Let $\theta_j^{(i)}(t_k)$ denote the collection of j -th layer parameters of client i 's neural network uploaded to the central server at the k -th communication round before FL aggregation. Let $\tilde{\theta}_j(t_k)$ denote the collection j -th layer parameters of the global neural network at the k -th communication round after FL aggregation, i.e.

$$\begin{aligned}\tilde{\theta}_j(t_k) &= \text{aggregation}(\theta_j^{(1)}(t_k), \dots, \theta_j^{(N)}(t_k)) \\ &= \phi(\theta_j^{(1)}(t_k), \dots, \theta_j^{(N)}(t_k))\end{aligned}$$

Define the risk functional of client i :

$$\mathcal{R}^{(i)}(\cdot) := \int_{\mathcal{X}} \int_{\mathcal{Y}} \ell(f(x; \cdot), y) d_{\pi_i} \Pr(x, y)$$

where \mathcal{X} denotes the feature/input space and \mathcal{Y} denotes the label/output space.

We first consider the path-norm evolution during local update of client i . Let E denote the time range of local update under gradient flow and $t_k = t_{k+1} + E, \forall k \leq R$. By Theorem 12, we have

$$\begin{aligned}\|\theta_j^{(i)}(t_{k+1})\|_{L^2(\pi^{j+1} \otimes \pi^j)} &\leq \|\tilde{\theta}_j(t_k)\|_{L^2(\pi^{j+1} \otimes \pi^j)} \\ &\quad + \sqrt{\mathcal{R}^{(i)}(\tilde{\theta}(t_k))} E^{1/2}\end{aligned}$$

Without loss of generality, consider the FedAvg aggregation scheme, i.e. for $1 \leq k \leq R$,

$$\begin{aligned}\tilde{\theta}(t_k) &= \phi(\theta^{(1)}(t_k), \dots, \theta^{(N)}(t_k)) \\ &= \frac{1}{N} \sum_{i=1}^N \theta^{(i)}(t_k)\end{aligned}$$

Let Ω_j denote the index space of the j -th layer of a neural network, and $\theta_j^{(i)}(w_{j+1}, w_j, t_{k+1})$ denote the weight of the neural network given indices w_{j+1}, w_j , then

$$\begin{aligned}&\|\tilde{\theta}_j(t_{k+1})\|_{L^2(\pi^{j+1} \otimes \pi^j)} \\ &= \left\| \frac{1}{N} \sum_{i=1}^N \theta_j^{(i)}(t_{k+1}) \right\|_{L^2(\pi^{j+1} \otimes \pi^j)} \\ &= \int_{\Omega_{j+1} \times \Omega_j} \left(\frac{1}{N} \sum_{i=1}^N \theta_j^{(i)}(w_{j+1}, w_j, t_{k+1}) \right)^2 dw_{j+1} w_j \\ &\leq \frac{1}{N} \sum_{i=1}^N \|\theta_j^{(i)}(t_{k+1})\|_{L^2(\pi^{j+1} \otimes \pi^j)} \\ &\leq \frac{1}{N} \sum_{i=1}^N \left(\|\tilde{\theta}_j(t_k)\|_{L^2(\pi^{j+1} \otimes \pi^j)} + \sqrt{\mathcal{R}^{(i)}(\tilde{\theta}(t_k))} E^{1/2} \right) \\ &\leq \|\tilde{\theta}_j(t_k)\|_{L^2(\pi^{j+1} \otimes \pi^j)} + \max_{1 \leq i \leq N} \sqrt{\mathcal{R}^{(i)}(\tilde{\theta}(t_k))} E^{1/2}\end{aligned}$$

$$\begin{aligned}&\leq \|\tilde{\theta}_j(t_k)\|_{L^2(\pi^{j+1} \otimes \pi^j)} + \max_{1 \leq i \leq N} \sqrt{\mathcal{R}^{(i)}(\tilde{\theta}(t_0))} E^{1/2} \\ &\leq \|\tilde{\theta}_j(t_0)\|_{L^2(\pi^{j+1} \otimes \pi^j)} + (k+1) \max_{1 \leq i \leq N} \sqrt{\mathcal{R}^{(i)}(\tilde{\theta}(t_0))} E^{1/2}\end{aligned}$$

Now we have an upper bound of the j -th layer parameters, we can then derive the upper bound of the path-norm. By Lemma 4.6 in (Weinan et al. 2020),

$$\|f\|_{\text{pnp}} \leq \left(C' + k \max_{1 \leq i \leq N} \sqrt{\mathcal{R}^{(i)}(\tilde{\theta}(t_0))} E^{1/2} \right)^{L+1}$$

where C' is a constant and $C' \geq \|\tilde{\theta}_j(t_0)\|_{L^2(\pi^{j+1} \otimes \pi^j)}$ for all $1 \leq j \leq L$.

Note that the above result also applies for SCAFFOLD with a new risk functional of client i :

$$\mathcal{R}^{(i)}(\cdot; f) = \int_{\mathbb{R}^{d_x + d_y}} (\ell(f(x; \cdot), y) + \langle \cdot, c_i \rangle) \mathbb{P}(dx \otimes dy) + M$$

where c_i denotes the client control variate and M is a constant. We introduce M to ensure that the risk functional is always non-negative. The gradient-flow dynamics does not depend on the choice of M .

Next, we consider a decentralized learning setting. Here we take the multi-consensus stochastic variance reduced extragradient algorithm as an example (Luo and Ye 2022). The local update analysis is similar to FL. As for the aggregation step in decentralized optimization, the communication step is typically written as the matrix multiplication

$$\tilde{\theta}(t_k) = W\theta(t_k)$$

where θ denotes the matrix which collects all clients' parameter vector $\theta^{(i)}$. We assume that

- $W_{ij} \neq 0$ if client i and j can exchange information;
- W is a symmetric matrix;
- $0 \preceq W \preceq I, W\mathbf{1} = \mathbf{1}, \text{null}(I - W) = \text{span}(\mathbf{1})$.

Let $\theta^{(i)}(t_k)$ denote the collection of all the parameters of client i 's neural network uploaded to the central server at the k -th communication round before decentralized communication. Let $\theta_j^{(i)}(t_k)$ denote the collection of j -th layer parameters of client i 's neural network uploaded to the central server at the k -th communication round before decentralized communication. Let $\tilde{\theta}_j^{(i)}(t_k)$ denote the collection j -th layer parameters of client i 's neural network at the k -th communication round after the decentralized communication.

By Lemma 2 in (Liu and Morse 2011) and Lemma 2.1 in (Luo and Ye 2022), we obtain a bound on the mixing rate of parameters,

$$\begin{aligned}&\left\| \tilde{\theta}_j^{(i)}(t_k) - \frac{1}{N} \sum_{I=1}^N \theta_j^{(I)}(t_k) \right\|_{L^2(\pi^{j+1} \otimes \pi^j)} \\ &\leq \lambda_2(W) \left\| \theta_j^{(i)}(t_k) - \frac{1}{N} \sum_{I=1}^N \theta_j^{(I)}(t_k) \right\|_{L^2(\pi^{j+1} \otimes \pi^j)}\end{aligned}$$

where $\lambda_2(W)$ denotes the second largest eigenvalue of W .

Then

$$\begin{aligned}
& \|\tilde{\theta}_j^{(i)}(t_{k+1})\|_{L^2(\pi^{j+1} \otimes \pi^j)} \\
& \leq (1 + 2\lambda_2(W)) \|\tilde{\theta}_j(t_k)\|_{L^2(\pi^{j+1} \otimes \pi^j)} \\
& \quad (1 + \lambda_2(W)) \max_{1 \leq i \leq N} \sqrt{\mathcal{R}^{(i)}(\tilde{\theta}(t_k))} E^{1/2} \\
& \leq (1 + 2\lambda_2(W)) \|\tilde{\theta}_j(t_k)\|_{L^2(\pi^{j+1} \otimes \pi^j)} \\
& \quad (1 + \lambda_2(W)) \max_{1 \leq i \leq N} \sqrt{\mathcal{R}^{(i)}(\tilde{\theta}(t_0))} E^{1/2} \\
& \leq (1 + 2\lambda_2(W))^{k+1} \|\tilde{\theta}_j(t_0)\|_{L^2(\pi^{j+1} \otimes \pi^j)} \\
& \quad + \frac{1 + \lambda_2(W)}{2\lambda_2(W)} ((1 + 2\lambda_2(W))^{k+1} - 1) \max_{1 \leq i \leq N} \sqrt{\mathcal{R}^{(i)}(\tilde{\theta}(t_0))} E^{1/2}
\end{aligned}$$

again, by Lemma 4.6 in (Weinan et al. 2020), we have

$$\|f\|_{\text{pnp}} = \mathcal{O} \left(e^{C''k} \max_{1 \leq i \leq N} \sqrt{\mathcal{R}^{(i)}(\tilde{\theta}(t_0))} E^{1/2} \right)^{L+1}$$

□