

Federated Learning for Non-IID Data: From Theory to Algorithm

Bojian Wei, Jian Li*, Yong Liu, and Weiping Wang

Institute of Information Engineering, Chinese Academy of Sciences
University of Chinese Academy of Sciences
Renmin University of China

18th Pacific Rim International Conference on Artificial Intelligence (PRICAI 2021)

November 10, 2021

1. Introduction
2. Preliminaries and Notations
3. Generalization Analysis
4. Algorithm: FedAvgR
5. Experiments
6. Conclusion

Non-IID Data Partitioning in Federated Learning

- In Federated Learning (FL) [McMahan et al. \[2017\]](#), the raw data of each client is stored locally and can not be obtained by any third party → **Non-IID** problem.
- **Non-IID** (not independently or identically distributed) data partitioning:
 - 1 The distribution is different on different clients.
 - 2 The amount of data among different clients is unbalanced.
 - 3 Data on some clients may be relevant.
- The non-IID problem leads to the decline of the model's effectiveness compared to centralized learning (CL).

Lack of Generalization Analysis

- Many studies try to solve the non-IID problem by designing new algorithms empirically Wang et al. [2020a], Smith et al. [2017], Pustozero et al. [2021], McMahan et al. [2017], Li et al. [2020a], Karimireddy et al. [2020], Yu et al. [2019], Wang et al. [2020b], Li et al. [2020c], Briggs et al. [2020], while only a few studies have carried out generalization analysis Mohri et al. [2019] for FL.
 - FedAvg applies iterative model averaging to deal with non-IID data.
 - Using local momentum instead of local SGD.
 - Using clustering to FL.
 - Agnostic FL provides a generalization view of FL, but the target is to optimize the worst case in the hypothesis.
- Lack of generalization analysis for FL under the traditional framework.

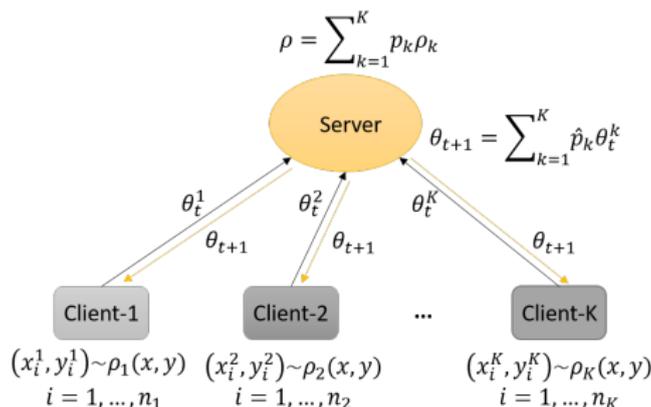
We analyze the **excess risk** for FL with non-IID data, which measures the gap between the global model trained by FL and the optimal model trained by CL.

- **Theoretically:** We give the *excess risk* bound between FL on non-IID data and CL for the first time and find out the factors that affect the accuracy decline. We give a reasonable explanation for the bound by decomposing the *excess risk* into three terms: *agnostic error*, *federated error* and *approximation error*.
- **Algorithmically:** We propose a novel algorithm FedAvgR (Federated Averaging with Regularization) to improve the performance of FL on non-IID data, which is regularized by Rademacher complexity and discrepancy distance.

1. Introduction
- 2. Preliminaries and Notations**
3. Generalization Analysis
4. Algorithm: FedAvgR
5. Experiments
6. Conclusion

Preliminaries and Notations

- Assume that there are K clients in a FL setting, where samples (x^k, y^k) on the k -th client with size of n_k are drawn i.i.d. from distribution ρ_k , data on different clients may not have the same distribution ($\rho_i \neq \rho_j$), and all clients participate in each communication round.
- The **global distribution** is assumed to be a mixture distribution of local distributions on all K clients: $\rho = \sum_{k=1}^K p_k \rho_k$, where p_k is the mixture weight ($\sum_{k=1}^K p_k = 1$). Actually, the mixture weight p_k is **unknown**, so an estimated weight \hat{p}_k will be applied in practice, which brings us the **estimated global distribution** $\tilde{\rho} = \sum_{k=1}^K \hat{p}_k \rho_k$.



Preliminaries and Notations

- The hypothesis space $\mathcal{H} = \{\mathbf{x} \rightarrow f(\mathbf{x})\}$ consists of labeling functions $f: \mathcal{X} \rightarrow \mathcal{Y}$, where $\mathcal{X} \subseteq \mathbb{R}^d$ represents the input space and $\mathcal{Y} \subseteq \mathbb{R}^C$ represents the label space. Let $\ell(f(\mathbf{x}), y)$ be the loss function, which is assumed to be upper bounded by M ($M > 0$), and $\mathcal{L} = \{\ell(f(\mathbf{x}), y) | f \in \mathcal{H}\}$ be the family of loss functions on \mathcal{H} , the expected loss of FL on ρ can be described as

$$\mathcal{E}_\rho(f) = \sum_{k=1}^K p_k \mathcal{E}_{\rho_k}(f) = \sum_{k=1}^K p_k \int_{\mathcal{X} \times \mathcal{Y}} \ell(f(\mathbf{x}), y) d\rho_k(\mathbf{x}, y),$$

and the corresponding empirical loss is

$$\hat{\mathcal{E}}_\rho(f) = \sum_{k=1}^K p_k \hat{\mathcal{E}}_{\rho_k}(f) = \sum_{k=1}^K p_k \frac{1}{n_k} \sum_{i=1}^{n_k} \ell(f(\mathbf{x}_i^k), y_i^k).$$

Contents

1. Introduction
2. Preliminaries and Notations
- 3. Generalization Analysis**
4. Algorithm: FedAvgR
5. Experiments
6. Conclusion

- **Excess risk:**

$$\mathcal{E}_\rho(\tilde{f}_{fl}) - \mathcal{E}_\rho(f^*) = \underbrace{\mathcal{E}_\rho(\tilde{f}_{fl}) - \mathcal{E}_\rho(\hat{f}_{fl})}_{A_1:=} + \underbrace{\mathcal{E}_\rho(\hat{f}_{fl}) - \mathcal{E}_\rho(\hat{f}_{cl})}_{A_2:=} + \underbrace{\mathcal{E}_\rho(\hat{f}_{cl}) - \mathcal{E}_\rho(f^*)}_{A_3:=}$$

- A_1 : *agnostic error*, A_2 : *federated error*, A_3 : *approximation error*.
- $\tilde{f}_{fl} = \arg \min_{f \in \mathcal{H}} \sum_{k=1}^K \hat{p}_k \hat{\mathcal{E}}_{\rho_k}(f)$ is the empirical learner of FL on $\tilde{\rho}$,
 $\hat{f}_{fl} = \arg \min_{f \in \mathcal{H}} \sum_{k=1}^K p_k \hat{\mathcal{E}}_{\rho_k}(f)$ is the empirical learner on ρ ,
 $\hat{f}_{cl} = \arg \min_{f \in \mathcal{H}} \frac{1}{n} \sum_{i=1}^n \ell(f(\mathbf{x}_i), y_i)$ is the empirical learner of CL, and
 $f^* = \arg \min_{f \in \mathcal{H}} \mathcal{E}_\rho(f)$ is the expected (optimal) learner in \mathcal{H} which minimizes the expected loss on ρ .
- The labeling function f is formed as $f(\mathbf{x}) = \mathbf{W}^T \phi(\mathbf{x})$, where $\mathbf{W} \in \mathbb{R}^{D \times C}$, $\phi(\mathbf{x}) \in \mathbb{R}^D$ and $\phi(\cdot)$ is the feature mapping with learnable parameters φ .

Bounds of Three Error Terms

To measure the performance gap of a model on different distributed data, we introduce the **discrepancy distance** Mansour et al. [2009] as follows:

$$disc_L(Q_1, Q_2) = \sup_{f \in \mathcal{H}} |\mathcal{E}_{Q_1}(f) - \mathcal{E}_{Q_2}(f)|,$$

where Q_1 and Q_2 are two different distributions.

Theorem (Agnostic Error Bound)

Assume that $\ell(f(\mathbf{x}), y)$ is λ -Lipschitz equipped with the 2-norm, that is $|\ell(f(\mathbf{x}), y) - \ell(f(\mathbf{x}'), y')| \leq \lambda \|\mathbf{x} - \mathbf{x}'\|_2$, $B = \sup_{f = \mathbf{W}^T \phi(\mathbf{x}) \in \mathcal{H}} \|\mathbf{W}\|_*$, where $\|\cdot\|_*$ denotes the trace norm. With probability at least $1 - \delta$ ($\delta > 0$):

$$A_1 \leq 2disc_L(\tilde{\rho}, \rho) + 4\sqrt{2}\lambda B \sum_{k=1}^K \frac{\hat{p}_k}{n_k} \sqrt{C} \|\phi(\mathbf{X}^k)\|_F + 6M \sqrt{\frac{\mathcal{S}(\hat{\mathbf{p}}|\bar{\mathbf{n}}) \log(2/\delta)}{2n}},$$

where $\|\phi(\mathbf{X}^k)\|_F = \sqrt{\sum_{i=1}^{n_k} \langle \phi(\mathbf{x}_i^k), \phi(\mathbf{x}_i^k) \rangle}$, $\mathcal{S}(\hat{\mathbf{p}}|\bar{\mathbf{n}}) = \chi^2(\hat{\mathbf{p}}|\bar{\mathbf{n}}) + 1$, χ^2 denotes the chi-squared divergence, $\hat{\mathbf{p}} = [\hat{p}_1, \dots, \hat{p}_K]$, and $\bar{\mathbf{n}} = \frac{1}{n} [n_1, \dots, n_K]$.

We first decompose A_1 into the following parts:

$$A_1 = \mathcal{E}_\rho(\tilde{f}_{fl}) - \mathcal{E}_{\tilde{\rho}}(\tilde{f}_{fl}) + \underbrace{\mathcal{E}_{\tilde{\rho}}(\tilde{f}_{fl}) - \mathcal{E}_{\tilde{\rho}}(\hat{f}_{fl})}_{A'_1 :=} + \mathcal{E}_{\tilde{\rho}}(\hat{f}_{fl}) - \mathcal{E}_\rho(\hat{f}_{fl})$$

We further decompose A'_1 as:

$$A'_1 = \underbrace{\mathcal{E}_{\tilde{\rho}}(\tilde{f}_{fl}) - \hat{\mathcal{E}}_{\tilde{\rho}}(\tilde{f}_{fl})}_{A_{11} :=} + \underbrace{\hat{\mathcal{E}}_{\tilde{\rho}}(\tilde{f}_{fl}) - \hat{\mathcal{E}}_{\tilde{\rho}}(\hat{f}_{fl})}_{A_{12} :=} + \underbrace{\hat{\mathcal{E}}_{\tilde{\rho}}(\hat{f}_{fl}) - \mathcal{E}_{\tilde{\rho}}(\hat{f}_{fl})}_{A_{13} :=}$$

- The remain parts of A_1 can be bounded by $2disc_L(\rho, \tilde{\rho})$.
- $\hat{\mathcal{E}}_{\tilde{\rho}}(\tilde{f}_{fl}) \leq \hat{\mathcal{E}}_{\tilde{\rho}}(\hat{f}_{fl}) \Rightarrow A_{12} \leq 0$.

Proof Sketch

Let \mathcal{H} be a hypothesis space of f defined over \mathcal{X} , \mathcal{L} be the family of loss functions associated to \mathcal{H} , $\mathbf{n} = [n_1, \dots, n_K]$ be the vector of sample sizes and $\mathbf{p} = [p_1, \dots, p_K]$ be the mixture weight vector, the empirical weighted Rademacher complexity of \mathcal{L} is

$$\widehat{\mathcal{R}}(\mathcal{L}, \mathbf{p}) = \mathbb{E}_\epsilon \left[\sup_{f \in \mathcal{H}} \sum_{k=1}^K \frac{p_k}{n_k} \sum_{i=1}^{n_k} \epsilon_i^k \ell(f(\mathbf{x}_i^k), y_i^k) \right],$$

and the empirical weighted Rademacher complexity of \mathcal{H} is

$$\widehat{\mathcal{R}}(\mathcal{H}, \mathbf{p}) = \mathbb{E}_\epsilon \left[\sup_{f \in \mathcal{H}} \sum_{k=1}^K \frac{p_k}{n_k} \sum_{i=1}^{n_k} \sum_{c=1}^C \epsilon_{ic}^k f_c(\mathbf{x}_i^k) \right],$$

where $f_c(\mathbf{x}_i^k)$ is the c -th value of $f(\mathbf{x}_i^k)$ corresponding to the C classes, ϵ_i^k s and ϵ_{ic}^k s are independent Rademacher variables, which are uniformly sampled from $\{-1, +1\}$, respectively.

A_{11} and A_{13} can be bounded by weighted Rademacher complexity:

- For any sample $S = \{S_1, \dots, S_n\}$ drawn from ρ , define $\Phi(S)$ by

$$\Phi(S) = \sup_{f \in \mathcal{H}} (\mathcal{E}_{\tilde{\rho}}(f) - \widehat{\mathcal{E}}_{\tilde{\rho}}(f)).$$

Let $S' = \{S'_1, \dots, S'_n\}$ be a sample differing from S only by point $x'_i{}^k$ in S'_k and x_i^k in S_k . Then, we have

$$|\Phi(S) - \Phi(S')| \leq \frac{\widehat{p}_k}{n_k} M.$$

Applying McDiarmid's inequality and Jensen's inequality, we get

$$\Phi(S) \leq 2\widehat{\mathcal{R}}(\mathcal{L}, \widehat{\mathbf{p}}) + 3M \sqrt{\frac{\chi^2(\widehat{\mathbf{p}} \parallel \bar{\mathbf{n}}) + 1}{2n} \log \frac{2}{\delta}}.$$

- $A_{11}, A_{13} \leq \Phi(S)$

Estimate weighted Rademacher complexity:

- Lipschitz assumption: $\widehat{\mathcal{R}}(\mathcal{L}, \widehat{\mathbf{p}}) \leq \sqrt{2\lambda} \widehat{\mathcal{R}}(\mathcal{H}, \widehat{\mathbf{p}})$.
- Rewriting $\widehat{\mathcal{R}}(\mathcal{H}, \widehat{\mathbf{p}})$:

$$\widehat{\mathcal{R}}(\mathcal{H}, \widehat{\mathbf{p}}) = \mathbb{E}_\epsilon \left[\sup_{f \in \mathcal{H}} \sum_{k=1}^K \frac{\widehat{p}_k}{n_k} \langle \mathbf{W}_k, \Phi_k \rangle \right],$$

where $\mathbf{W}_k, \Phi_k = [\sum_{i=1}^{n_k} \epsilon_{i1}^k \phi(\mathbf{x}_i^k), \dots, \sum_{i=1}^{n_k} \epsilon_{iC}^k \phi(\mathbf{x}_i^k)] \in \mathbb{R}^{D \times C}$ and $\langle \mathbf{W}_k, \Phi_k \rangle = \text{Tr}(\mathbf{W}_k^T \Phi_k)$.

- Hölder's inequality: $\widehat{\mathcal{R}}(\mathcal{H}, \widehat{\mathbf{p}}) \leq B \sum_{k=1}^K \frac{\widehat{p}_k}{n_k} \sqrt{\mathbb{E}_\epsilon [\|\Phi_k\|_F^2]}$.
- Jensen's inequality: $\mathbb{E}_\epsilon [\|\Phi_k\|_F^2] \leq C \|\phi(\mathbf{X}^k)\|_F^2$.

Bounds of Three Error Terms

Theorem (Federated Error Bound)

Under the same assumptions as the Theorem before, with probability at least $1 - \delta$ ($\delta > 0$), we have:

$$A_2 \leq \sum_{k=1}^K p_k \left(\text{disc}_L(\rho_k, \rho) + \frac{4\sqrt{2}\lambda B}{n_k} \sqrt{C} \|\phi(\mathbf{X}^k)\|_F \right) + \sum_{k=1}^K p_k \left(6M \sqrt{\frac{\log(2/\delta)}{2n_k}} \right).$$

Theorem (Approximation Error Bound)

Under the same assumptions as the Theorem before, with probability $1 - \delta$ ($\delta > 0$), we have:

$$A_3 \leq \frac{4\sqrt{2}\lambda B}{n} \sqrt{C} \|\phi(\mathbf{X})\|_F + 3M \sqrt{\frac{\log(2/\delta)}{2n}},$$

where $\|\phi(\mathbf{X})\|_F = \sqrt{\sum_{i=1}^n \langle \phi(\mathbf{x}_i), \phi(\mathbf{x}_i) \rangle}$.

Proof Sketch

Note that $A_2 = \sum_{k=1}^K p_k \underbrace{[\mathcal{E}_{\rho_k}(\hat{f}_{fl}) - \mathcal{E}_{\rho}(\hat{f}_{cl})]}_{A'_2}$, we decompose A'_2 as:

$$\underbrace{\mathcal{E}_{\rho_k}(\hat{f}_{fl}) - \hat{\mathcal{E}}_{\rho_k}(\hat{f}_{fl})}_{A_{21}} + \underbrace{\hat{\mathcal{E}}_{\rho_k}(\hat{f}_{fl}) - \hat{\mathcal{E}}_{\rho_k}(\hat{f}_{cl})}_{A_{22}} + \underbrace{\hat{\mathcal{E}}_{\rho_k}(\hat{f}_{cl}) - \mathcal{E}_{\rho_k}(\hat{f}_{cl})}_{A_{23}} + \underbrace{\mathcal{E}_{\rho_k}(\hat{f}_{cl}) - \mathcal{E}_{\rho}(\hat{f}_{cl})}_{A_{24}}$$

Substituting A_{22} into the equation of A_2 , due to the definition of \hat{f}_{fl} , we have

$$\sum_{k=1}^K p_k [\hat{\mathcal{E}}_{\rho_k}(\hat{f}_{fl}) - \hat{\mathcal{E}}_{\rho_k}(\hat{f}_{cl})] \leq 0$$

- A_{21} and A_{23} can be bounded by Rademacher complexity.
- A_{24} can be bounded by $disc_L(\rho_k, \rho)$.
- A_3 is a constant multiple of the generalization bound for CL, which can be bounded by Rademacher complexity, as well.

Theorem (Excess Risk Bound)

Under the same assumptions as the Theorem before, With probability at least $1 - \delta$ ($\delta > 0$), the excess risk bound of federated learning on non-IID data holds as follows:

$$\mathcal{E}_\rho(\widehat{f}_{fl}) - \mathcal{E}_\rho(f^*) \leq \mathcal{O}(G_1 + G_2 + G_3),$$

where $G_1 = \text{disc}_L(\tilde{\rho}, \rho) + \sum_{k=1}^K \frac{\widehat{p}_k B \sqrt{C}}{n_k} \|\phi(\mathbf{X}^k)\|_F + \sqrt{\frac{S(\widehat{\mathbf{p}}|\widehat{\mathbf{n}})}{n}}$,

$$G_2 = \frac{B\sqrt{C}}{n} \|\phi(\mathbf{X})\|_F + \sqrt{\frac{1}{n}} \text{ and}$$

$$G_3 = \sum_{k=1}^K p_k \left[\text{disc}_L(\rho_k, \rho) + \frac{B\sqrt{C}}{n_k} \|\phi(\mathbf{X}^k)\|_F + \sqrt{\frac{1}{n_k}} \right].$$

Excess Risk Bound

- According to the bound, we can lower the *excess risk* by reducing $disc_L(\rho_k, \rho)$, $\|\mathbf{W}\|_*$, $\|\phi(\mathbf{X}^k)\|_F$, and $disc_L(\tilde{\rho}, \rho)$.
- In non-IID condition, samples on different clients are drawn from different distributions, so the gap between ρ_k and ρ certainly exists. Furthermore, p_k is unknown, **how can we reduce** $disc_L(\tilde{\rho}, \rho)$?

$$disc_L(\rho_k, \rho) \downarrow \Rightarrow disc_L(\tilde{\rho}, \rho) \downarrow$$

Especially, when $\rho_k = \rho$, whatever value we choose for \hat{p}_k , it's not going to make big difference to the global distribution. Therefore, we are able to lower the *excess risk* by reducing $disc_L(\rho_k, \rho)$, $\|\mathbf{W}\|_*$ and $\|\phi(\mathbf{X}^k)\|_F$.

- **Corollaries** ($\phi(\cdot)$ is upper bounded by κ^2 and $\hat{p}_k = p_k$)
 - $disc_L(\rho_k, \rho) = 0$: $\mathcal{O}\left((\kappa B \sqrt{C} + 1) \sum_{k=1}^K \hat{p}_k \sqrt{\frac{1}{n_k}}\right)$
 - $n_k = n/K$: $\mathcal{O}\left(\kappa B \sqrt{KC/n}\right)$ (**distributed learning**)
 - $K = 1$: $\mathcal{O}(\kappa B \sqrt{C/n})$ (**centralized learning**)

Contents

1. Introduction
2. Preliminaries and Notations
3. Generalization Analysis
- 4. Algorithm: FedAvgR**
5. Experiments
6. Conclusion

FedAvgR: Federated Averaging with Regularization

- 1 We choose MMD (Maximum Mean Discrepancy) [Borgwardt et al. \[2006\]](#) to measure the distance between different distributions Q_1 and Q_2 , which is formed as

$$\text{MMD}[Q_1, Q_2] = \sup_{f \in \mathcal{H}} (\mathbb{E}_{Q_1}[f(\mathbf{x})] - \mathbb{E}_{Q_2}[f(\mathbf{x})]).$$

The local distribution ρ_k won't change during training, so we shall reduce the discrepancy after feature mapping. In other words, we can reduce $\text{disc}_L(\rho_k^\phi, \rho^\phi)$ instead of $\text{disc}_L(\rho_k, \rho)$, where ρ_k^ϕ and ρ^ϕ are respectively the local feature distribution on client k and global feature distribution.

- 2 Taking $\text{MMD}[\rho_k^\phi, \rho^\phi]$ as a regularizer with $\|\mathbf{W}\|_*$ and $\|\phi(\mathbf{X}^k)\|_F$, the objective function on the k -th client is

$$\min_{\mathbf{W}, \varphi} \frac{1}{n_k} \sum_{i=1}^{n_k} \ell(f(\mathbf{x}_i^k), y_i^k) + \alpha \|\mathbf{W}\|_* + \beta \|\phi(\mathbf{X}^k)\|_F + \gamma \text{MMD}[\rho_k^\phi, \rho^\phi].$$

FedAvgR: Federated Averaging with Regularization

Algorithm 1 FedAvgR. K clients are indexed by k , \mathcal{B} is the local mini-batch size, E is the number of local epochs, η is the learning rate, \mathbf{F} represents the objective function.

Server-Aggregate

```
1: initialize  $\mathbf{W}_0$  and  $\varphi_0$ 
2: for  $k = 1, \dots, K$  do
3:    $\hat{\rho}_k^\phi \leftarrow$  estimate the distribution of  $\phi(\mathbf{x})$ 
4:   upload the parameters of  $\hat{\rho}_k^\phi$  to the server
5: end for
6: get the global distribution  $\hat{\rho}^\phi = \sum_{k=1}^K \hat{p}_k \hat{\rho}_k^\phi$ 
7: for each round  $t = 1, 2, \dots$  do
8:   for each client  $k$  do
9:      $\mathbf{W}_{t+1}^k, \varphi_{t+1}^k, \hat{\rho}_k^\phi \leftarrow$  Client-Update( $k, \mathbf{W}_t, \varphi_t, \hat{\rho}^\phi$ )
10:  end for
11:  update the global distribution  $\hat{\rho}^\phi$ 
12:   $\mathbf{W}_{t+1} \leftarrow \sum_{k=1}^K \hat{p}_k \mathbf{W}_{t+1}^k, \varphi_{t+1} \leftarrow \sum_{k=1}^K \hat{p}_k \varphi_{t+1}^k$ 
13: end for
```

Client-Update($k, \mathbf{W}_t, \varphi_t, \hat{\rho}^\phi$)

```
1: draw samples  $Z_{\hat{\rho}^\phi}$  from  $\hat{\rho}^\phi$ 
2: for epoch = 1, ...,  $E$  do
3:   for  $(\mathbf{x}, y) \in \mathcal{B}$  do
4:     calculate  $\text{MMD}[\hat{\rho}_k^\phi, \hat{\rho}^\phi]$  by  $(\mathbf{x}, y)$  and  $Z_{\hat{\rho}^\phi}$ 
5:      $\mathbf{F} = \frac{1}{\mathcal{B}} \sum_{(\mathbf{x}, y) \in \mathcal{B}} \ell(f(\mathbf{x}), y) + \alpha \|\mathbf{W}\|_* + \beta \|\phi(\mathbf{X})\|_F + \gamma \text{MMD}[\hat{\rho}_k^\phi, \hat{\rho}^\phi]$ 
6:      $\mathbf{W}_{t+1}^k \leftarrow \mathbf{W}_t - \eta \nabla_{\mathbf{W}_t} \mathbf{F}, \varphi_{t+1}^k \leftarrow \varphi_t - \eta \nabla_{\varphi_t} \mathbf{F}$ 
7:   end for
8:    $\hat{\rho}_k^\phi \leftarrow$  estimate the distribution of  $\phi(\mathbf{x})$ 
9: end for
```

Contents

1. Introduction
2. Preliminaries and Notations
3. Generalization Analysis
4. Algorithm: FedAvgR
- 5. Experiments**
6. Conclusion

Experimental Setup

- **Environment:** All the experiments are trained on a Linux_x86_64 server (CPU: Intel(R) Xeon(R) Silver 4214 (RAM: 196 GB) / GPU: NVIDIA GeForce RTX-2080ti).
- **Datasets:**
 - ① The synthetic dataset in our experiment is generated related to the method in [Li et al. \[2020b\]](#), where the number of samples n_k on client k follows a power law.
 - ② We apply the partitioning method related to [McMahan et al. \[2017\]](#) to some LIBSVM [Chang and Lin \[20\]](#) datasets to get non-IID data. We sort each dataset by the label and divide it into N/N_s shards of size N_s , where N is the total number of samples, then we assign each client 2 shards.
- **Model:** We use Random Fourier Feature [Rahimi and Recht \[2007\]](#) to do the feature mapping $\phi(\cdot)$, which is formed as $\sqrt{2} \cos(\omega^T x + b)$, where ω is sampled from $\mathcal{N}(0, \sigma^2)$, σ is related to the corresponding Gaussian kernel, and b is uniformly sampled from $[0, 2\pi]$.

Table 1: Information of Different Datasets

Dataset	Class	Training Size	Testing Size	Features
a1a	2	1605	30956	123
svmguide1	2	3089	4000	4
splice	2	1000	2175	60
vehicle	4	500	446	18
dna	3	2000	1186	180
pendigits	10	7494	3498	16
satimage	6	4435	2000	36
usps	10	7291	2007	256
MNIST	10	60000	10000	28×28

Comparative Experiment

We compare FedAvgR with OneShot Zhang et al. [2015], FedAvg McMahan et al. [2017], FedProx Li et al. [2020a] and FL+HC Briggs et al. [2020].

- OneShot aggregates local models when local trainings converge.
- FedAvg iteratively averages local models by n_k/n .
- FedProx adds the last-round's global model to local training as regularization based on FedAvg.
- FL+HC uses hierarchical clustering to divide clients into several clusters and applies FedAvg separately.

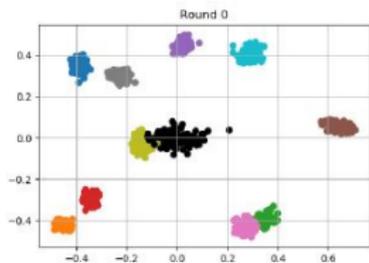
Comparative Experiment

Table 2: Test Accuracy on Real-World Datasets. We run methods on each dataset 10 times, each with 300 rounds. We bold the numbers of the best method and underline the numbers of other methods which are not significantly worse than the best one.

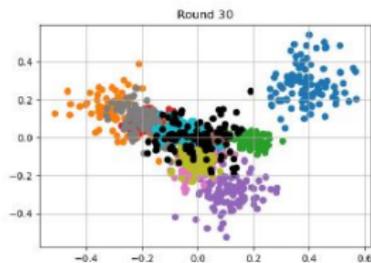
Dataset	OneShot	FedAvg	FedProx	FL+HC	FedAvgR
a1a	76.86±0.30	<u>84.29±0.06</u>	<u>84.27±0.06</u>	81.63±0.94	84.30±0.06
svmguidel	71.50±4.21	90.95±0.86	91.19±0.84	85.66±4.48	91.77±1.01
splice	75.95±4.56	<u>90.37±0.21</u>	<u>90.38±0.20</u>	85.12±2.14	90.40±0.26
vehicle	52.31±4.36	78.61±1.08	78.58±1.06	62.24±8.12	78.82±0.98
dna	63.73±1.02	95.23±0.17	95.18±0.21	92.09±3.25	95.59±0.23
pendigits	46.70±2.32	94.87±0.58	94.85±0.59	86.81±4.58	95.12±0.48
satimage	73.07±2.39	<u>88.83±0.41</u>	88.46±0.31	76.72±2.96	88.93±0.39
usps	56.83±4.06	94.57±0.15	94.53±0.13	88.03±3.62	94.80±0.19
MNIST	68.80±2.06	97.26±0.09	97.24±0.07	85.13±2.23	97.34±0.06

Impact of MMD $[\rho_k^\phi, \rho^\phi]$

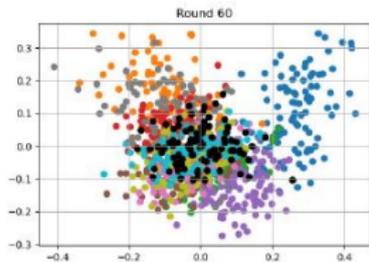
We run 100 rounds on the synthetic dataset with $(u, v) = (1, 1)$ and sample 100 points from each $\hat{\rho}_k^\phi$ and $\hat{\rho}^\phi$ to show the impact of MMD $[\rho_k^\phi, \rho^\phi]$.



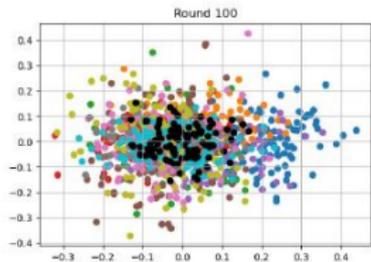
(a) Round 0: 0.766



(b) Round 30: 0.098



(c) Round 60: 0.034



(d) Round 100: 0.019

Impacts of Different Regularizers

We run 250 rounds on the synthetic dataset with $(u, v) = (0.5, 0.5)$ and some real-world datasets with non-IID partitioning.

Table 3: Test Accuracy of FedAvgR with Different Regularizers

Dataset	No Regularizer	$\ \mathbf{W}\ _*$	$\ \phi(\mathbf{X}^k)\ _F$	MMD	All Regularizers
svmguide1	89.05	89.20	89.45	89.61	90.70
vehicle	77.12	77.17	77.17	77.46	78.32
dna	95.33	95.52	95.36	95.45	95.70
pendigits	95.70	95.71	95.74	95.94	95.90
usps	94.57	94.72	94.82	94.80	94.82
synthetic	95.82	96.07	96.06	96.12	96.23

Contents

1. Introduction
2. Preliminaries and Notations
3. Generalization Analysis
4. Algorithm: FedAvgR
5. Experiments
- 6. Conclusion**

- In this paper, we give an *excess risk bound* for **federated learning on non-IID data** through **Rademacher complexity** and **discrepancy distance**, analyzing the error between it and the optimal centralized learning model. Based on our theory, we propose FedAvgR to improve the performance of federated learning in non-IID setting, where three regularizers are added to achieve a sharper bound. Experiments show that our algorithm outperforms the previous methods. As the first work to analyze the *excess risk* under a more general framework, our work will provide a reference for the future study of generalization properties in federated learning with non-IID data. Besides, the proof techniques in this paper are helpful to the research of error analysis related to the distributed framework.

References

- Karsten M. Borgwardt, Arthur Gretton, Malte J. Rasch, Hans-Peter Kriegel, Bernhard Schölkopf, and Alexander J. Smola. Integrating structured biological data by kernel maximum mean discrepancy. In *Proceedings of the 14th International Conference on Intelligent Systems for Molecular Biology*, pages 49–57, 2006.
- Christopher Briggs, Zhong Fan, and Peter Andras. Federated learning with hierarchical clustering of local updates to improve training on non-iid data. In *International Joint Conference on Neural Networks, IJCNN*, pages 1–9. IEEE, 2020.
- Chih-Chung Chang and Chin-Jen Lin. Libsvm: a library for support vector machines. *ACM Transactions on Intelligent Systems and Technology*, pages 2:27:1–27:27, 20.
- Sai Praneeth Karimireddy, Satyen Kale, Mehryar Mohri, Sashank J. Reddi, Sebastian U. Stich, and Ananda Theertha Suresh. SCAFFOLD: stochastic controlled averaging for federated learning. In *ICML*, volume 119, pages 5132–5143, 2020.
- Tian Li, Anit Kumar Sahu, Manzil Zaheer, Maziar Sanjabi, Ameet Talwalkar, and Virginia Smith. Federated optimization in heterogeneous networks. In *MLSys*, 2020a.
- Xiang Li, Kaixuan Huang, Wenhao Yang, Shusen Wang, and Zhihua Zhang. On the convergence of fedavg on non-iid data. In *ICLR*, 2020b.
- Zhize Li, Dmitry Kovalev, Xun Qian, and Peter Richtárik. Acceleration for compressed gradient descent in distributed and federated optimization. In *ICML*, volume 119, pages 5895–5904, 2020c.
- Yishay Mansour, Mehryar Mohri, and Afshin Rostamizadeh. Domain adaptation: Learning bounds and algorithms. In *COLT*, 2009.
- Brendan McMahan, Eider Moore, Daniel Ramage, Seth Hampson, and Blaise Agüera y Arcas. Communication-efficient learning of deep networks from decentralized data. In *AISTATS*, volume 54, pages 1273–1282, 2017.
- Mehryar Mohri, Gary Sivek, and Ananda Theertha Suresh. Agnostic federated learning. In *ICML*, volume 97, pages 4615–4625, 2019.
- Anastasia Pustozherova, Andreas Rauber, and Rudolf Mayer. Training effective neural networks on structured data with federated learning. In *AINA*, volume 226, pages 394–406, 2021.
- Ali Rahimi and Benjamin Recht. Random features for large-scale kernel machines. In *NIPS*, pages 1177–1184, 2007.
- Virginia Smith, Chao-Kai Chiang, Maziar Sanjabi, and Ameet S. Talwalkar. Federated multi-task learning. In *NIPS*, pages 4424–4434, 2017.
- Hongyi Wang, Mikhail Yurochkin, Yuekai Sun, Dimitris S. Papailiopoulos, and Yasaman Khazaeni. Federated learning with matched averaging. In *ICLR*, 2020a.
- Jiayu Wang, Vinayak Tantia, Nicolas Ballas, and Michael G. Rabbat. Slowmo: Improving communication-efficient distributed SGD with slow momentum. In *ICLR*, 2020b.
- Hao Yu, Rong Jin, and Sen Yang. On the linear speedup analysis of communication efficient momentum SGD for distributed non-convex optimization. In *ICML*, volume 97, pages 7184–7193, 2019.
- Yuchen Zhang, John C. Duchi, and Martin J. Wainwright. Divide and conquer kernel ridge regression: a distributed algorithm with minimax optimal rates. *J. Mach. Learn. Res.*, 16:3299–3340, 2015.



中国科学院信息工程研究所
INSTITUTE OF INFORMATION ENGINEERING, CAS



中国科学院大学
University of Chinese Academy of Sciences



中國人民大學
RENMIN UNIVERSITY OF CHINA

Thank You