# Novel clustered federated learning based on local loss

Endong Gu, Yongxin Chen, Hao Wen, Xingju Cai, Deren Han

*Abstract*—This paper proposes `LCFL`, a novel clustering metric for evaluating clients' data distributions in federated learning. `LCFL` aligns with federated learning requirements, accurately assessing client-to-client variations in data distribution. It offers advantages over existing clustered federated learning methods, addressing privacy concerns, improving applicability to non-convex models, and providing more accurate classification results. `LCFL` does not require prior knowledge of clients' data distributions. We provide a rigorous mathematical analysis, demonstrating the correctness and feasibility of our framework. Numerical experiments with neural network instances highlight the superior performance of `LCFL` over baselines on several clustered federated learning benchmarks.

*Index Terms*—federated learning, clustering, local loss, LCFL

## I. INTRODUCTION

The rapid development of the Internet and the mass popularization of private devices have generated vast amounts of personal data. Benefiting from the fast growth of storage and computational capacities, companies can fully utilize this data to predict market demand through machine learning methods and make profits. As a result, machine learning research has gradually shifted from model-driven to data-driven. Although the full mining of user data has brought more personalized services, people's concerns about their data being freely sold and privacy leakage have also followed. Various nations have progressively passed relevant laws and regulations on privacy protection [1]. In this context, a new machine learning paradigm known as "federated learning" has been proposed in the field of machine learning [2], [3]. It investigates how to use data from all parties (which we refer to as clients) to complete machine learning while upholding user privacy.

The term "federated learning", abbreviated as `FL`, was introduced in 2016 by Google's academic team [2] which is a distributed machine learning model training strategy that focuses more on privacy protection and communication efficiency. A standard `FL` progress includes an iterative three-step protocol [4] (illustrated in Figure 1). In each communication round, denoted as $t$, the clients first download the latest global

Endong Gu is with the LMIB, School of Mathematical Sciences, Beihang University, Beijing 100191, P.R. China, e-mail: endonggu@buaa.edu.cn

Yongxin Chen is with the LMIB, School of Mathematical Sciences, Beihang University, Beijing 100191, P.R. China, e-mail: chenyongxin@buaa.edu.cn

Hao Wen is is with the College of Science, China Agricultural University, Beijing 100083, P.R. China, e-mail: wenh06@cau.edu.cn

Xingju Cai is with the School of Mathematical Sciences, Nanjing Normal University, Nanjing 210023, P.R. China, e-mail: caixingju@njnu.edu.cn

Deren Han is the corresponding author and is with the School of Mathematical Sciences, Beihang University, LMIB and NSLSCS, Beijing 100191, China, e-mail: handr@buaa.edu.cn

model, denoted as $w_t$, from the central server. This ensures that all clients start with the same initial model. Then each client independently improves the downloaded model by using their respective local data, denoted as $D_i$, to train the model on their devices. The local training process typically involves performing multiple iterations of an optimization algorithm, such as stochastic gradient descent, on mini-batches of data to prevent overfitting and enhance generalization. This process allows the model to learn from the unique characteristics of each client's data and valuable work focusing on this process has been done [5]. After local model training, all clients upload their updated model parameters to the central server. Finally, the server aggregates the received model updates. The aggregation step can involve simple averaging or more sophisticated techniques, such as secure aggregation or weighted averaging considering the importance of each client's update. This step ensures that the aggregated model reflects a collective representation of all clients' contributions while accounting for variations in the size and quality of local datasets.
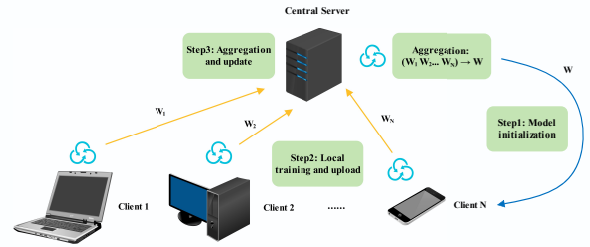


Fig. 1. `FL` three-step protocol illustration

By iteratively executing these steps across multiple rounds, `FL` enables the collaborative training of a shared model while preserving data privacy. The protocol ensures that clients can benefit from the aggregated knowledge of the entire network without directly sharing their raw data with the server or other clients.

Two key challenges distinguish `FL` from traditional distributed optimization: high degrees of systems and statistical heterogeneity [6], [7], [8]. System heterogeneity in `FL` encompasses significant variations in system characteristics among

devices involved in the learning process, including disparities in hardware specifications (such as CPU and memory), network connectivity options (such as 3G, 4G, 5G, and Wi-Fi), and power availability (e.g., battery level). These diverse system-level attributes introduce additional challenges, particularly in the areas of straggler mitigation and fault tolerance within the federated network [9]. On the other hand, statistical heterogeneity in FL refers to the existence of non-identically distributed data across the network, where each device or client possesses a distinct local dataset influenced by factors such as user preferences, geographical locations, or device-specific characteristics. This statistical heterogeneity gives rise to challenges concerning varying model performance and potential biases and fairness issues.

Federated Averaging (FedAvg) is a fundamental algorithm in FL that enables the training of machine learning models on decentralized data while preserving privacy. It involves iterative rounds of training on local devices, followed by aggregation of model updates through a central server using weighted averaging. However, FedAvg has drawbacks, including high communication and bandwidth requirements, challenges with straggler devices that can slow down the training, and difficulties in handling non-IID data distributions. Ongoing research aims to address these limitations to enhance the efficiency and effectiveness of FL [5], [9], [10], [11], [12], [13].

These works focus on improving efficiency and effectiveness with non-IID data [2]. Federated Proximal (FedProx) is an algorithm for federated learning that tackles this non-IID data challenge [9]. It extended FedAvg by introducing a proximal term to the optimization objective to encourage consistent model updates across devices. It improved global model convergence and generalization performance. However, FedProx has drawbacks including hyperparameter sensitivity, communication overhead, and lack of explicit handling of straggler devices. Splitting Scheme for Solving Federated Problems (FedSplit) is an algorithmic framework introduced in [10] for efficient distributed convex minimization in a hub-and-spoke model inspired by FL. [10] examined previous procedures, FedAvg and FedProx, revealing that their fixed points do not necessarily correspond to stationary points of the original optimization problem, even in simple convex scenarios with deterministic updates. To address this issue, FedSplit applied the Peaceman-Rachford splitting technique [14] into FL, ensuring that fixed points align with the optima of the original problem. However, it should be noted that FedSplit does not explicitly account for system heterogeneity, which limits its practicality in FL settings. [11] also noticed the problem that FedAvg and FedProx cannot converge to the global stationary solution and proposed a novel algorithmic framework based on primal-dual optimization called Federated Primal-Dual Algorithm (FedPD) to address this problem. Nevertheless, FedPD's requirement for all clients to participate in each communication round makes it less practical and applicable in FL settings as well. To overcome this limitation, [12] combined a nonconvex Douglas-Rachford splitting method [15], randomized block-coordinate strategies, and asynchronous implementation and proposed Federated Douglas-Rachford (FedDR). Unlike previous methods such as FedSplit and

FedPD, FedDR updates only a subset of users at each communication round, potentially in an asynchronous manner, making it more practical for real-world implementations. It also achieves the best-known $O(\epsilon^{-2})$ communication complexity for finding a stationary point under standard assumptions, where $\epsilon$ is a given accuracy.

Moreover, many researchers pay more attention to designing robust models [16], [17], reducing expensive communication [18] and preserving the privacy of user data [19], [20]. It is necessary to not only consider the similarity between data and integrate similar data to improve model performance but also identify the data with differences and develop model personalization ability. Personalized FL that meets the two requirements at the same time has become one of the hot topics in FL research [6], [7]. To address statistical heterogeneity, clustered FL is also a hot topic that is studied in this paper. For more details, the readers are referred to the survey papers [6], [7], [8], [21].

Before introducing clustered FL, it is essential to provide a concise overview of clustering methods, a classical unsupervised machine learning technique. Unsupervised machine learning aims to explore and uncover the inherent structure and patterns within unlabeled training data, without relying on class label information. Clustering, as a fundamental task in unsupervised learning, partitions a dataset into disjoint subsets, called clusters, to identify potential categories. Given the vast array of clustering algorithms available, we will highlight some classical algorithms categorized into four main groups: partitioning-based, density-based, hierarchical-based, and model-based [22].

Partitioning-based algorithms rapidly determine clusters by initially assigning samples and iteratively reallocating them to appropriate groupings. One well-known algorithm in this category is $k$-means, which has been reinvented many times in history by scholars in various fields such as Steinhaur (1956), Lloyd (1957), and McQueen (1967) [23], [24]. Variants such as $k$-medoids [25] enforce cluster centers to be actual training samples, while the number of clusters, $k$, is typically predetermined, although heuristics exist for automatic determination [26], [27]. Density-based clustering algorithms characterize the sample distribution density to discover clusters of arbitrary shapes. Noteworthy methods in this category include Density-Based Spatial Clustering of Applications with Noise (DBSCAN) [28], Ordering Points To Identify the Clustering Structure (OPTICS) [29] and Density-Based Clustering of Applications with Noise using Local Updates (DENCLUE) [30]. Hierarchical-based algorithms organize data hierarchically based on proximity. Agglomerative Nesting (AGNES) [31] adopts a bottom-up strategy, while Divisive Analysis (DIANA) [31] follows a top-down approach. Hierarchical clustering algorithms like Balanced Iterative Reducing and Clustering using Hierarchies (BIRCH) [32] and Robust Clustering using Links (ROCK) [33] improve upon AGNES and DIANA by allowing backtracking adjustments on merged or split clusters. Model-based methods optimize the fit between the given data and a predefined mathematical model, assuming the data is generated by a mixture of underlying probability distributions. Model-based Clustering (MCLUST) [34] is a well-known algorithm in this

category. For a more comprehensive understanding of clustering methods, we recommend referring to the papers [23], [35], [36], [37].

Given that clustering can uncover potential patterns in data, it is a natural approach to cluster users with similar data distributions into cohesive groups and train a FL model for each cluster [38], [39]. However, the challenge of clustering in FL arises from the fact that the classification goal is based on the distribution of user data, which cannot be shared in real-world scenarios. Consequently, it becomes necessary to employ alternative information, rather than the distribution itself, to discern the differences among users.

In early research on FL, the prevailing assumption was that all users shared a single global model. However, the heterogeneity of data distribution often led to sub-optimal performance of the global model on individual user data [38]. To address this statistical heterogeneity, personalized FL approaches have emerged, including the notable clustered FL. Clustering users based on data distribution has gained significant attention in various applications such as recommendation systems and precision medicines.

To estimate the clustering structure of data distribution, [38] employed the gradient cosine similarity $\alpha_{i,j}$ as a metric for user classification using the bisection method and hierarchical structures to select the optimal cluster.

$$\alpha_{i,j} = \frac{\langle \nabla L_i(w^*), \nabla L_j(w^*) \rangle}{\|\nabla L_i(w^*)\| \|\nabla L_j(w^*)\|}$$

Subsequently, a customized FL model is provided for each cluster. However, it is important to note that [38] relies on gradient information, which poses inherent privacy risks [40]. Moreover, their clustering procedure is centralized, resulting in high computational costs, especially when dealing with a large number of devices.

To overcome these limitations, [41] and [42] utilize local model parameters trained by users instead of relying on gradient information as metrics (1).

$$d_{i,j} = \|w_i^* - w_j^*\| \tag{1}$$

[42] introduced non-convex paired penalty functions and proposed the Fusion Penalized Federated Clustering(FPFC) algorithm, which performs the clustering process and model updates simultaneously, eliminating the need for central clustering. However, we observe that employing model parameters as metrics may lead to incorrect classification and deviate from the original intent of clustering in federated learning.

Drawing inspiration from the classical $k$-means algorithm, Mansor et al. and Ghosh et al. proposed similar algorithms, namely HYPCLUSTER [43] and the Iterative Federated Clustering Algorithm (IFCA) [39], respectively. These algorithms alternate between the clustering and model update processes. The cluster assignment for each user is estimated by minimizing the loss function. Specifically, the cluster identity $\hat{j}$ of client $i$ is estimated by

$$\hat{j} = \arg\min_{j \in [k]} L_i(w_j).$$

[39] established the convergence rate of the population loss function under favorable initialization, ensuring both convergence of the training loss and generalization to test data.

[43] provided guarantees only for generalization. While these algorithms eliminate the need for centralized clustering, they still require the specification of the number of clusters $k$, and the clustering results are typically sensitive to the initial values.

In our proposed clustering framework, we adopt a loss function-based metric (2), which avoids the privacy leakage caused by using gradient information.

$$d(i,j) = |L_i(w_i) - L_i(w_j)| + |L_j(w_i) - L_j(w_j)|. \tag{2}$$

Unlike [39] and [43], our framework does not assume the cluster number $k$ in advance. Furthermore, our framework is flexible and compatible with various clustering methods, such as binary hierarchical clustering or density-based clustering algorithms like DBSCAN. If the value of $k$ is known, the $k$-medoids clustering algorithm can also be utilized.

*a) Contributions:* In this paper, we present a novel clustering framework based on loss function metrics, called "Loss based Clustered Federated Learning (LCFL)". The primary objective of LCFL is to enable the use of clustering methods in FL while ensuring compliance with data protection regulations. Our main contributions can be summarized as follows:

1) Flexibility in clustering methods: LCFL does not impose any specific clustering algorithm, but instead supports partitioning-based clustering, density-based clustering, and hierarchical-based clustering methods. This flexibility allows practitioners to choose the most suitable clustering method based on their practical experience and specific requirements.

2) Efficient and privacy-preserving solution: By leveraging loss function metrics, LCFL provides an efficient and privacy-preserving solution for clustering in FL. Instead of sharing raw data, our framework makes use of the loss function metrics, addressing the challenges of data privacy and security, allowing multiple clients to collaborate and contribute their local data without compromising individual data privacy.

3) Valuable contribution to FL: The LCFL framework offers a valuable contribution to the field of FL by facilitating the application of clustering algorithms in a privacy-preserving manner. Its compatibility with different clustering methods and FL methods enhances its practicality and usefulness in various FL scenarios.

By combining these contributions, LCFL empowers practitioners to utilize clustering methods effectively in FL while upholding data protection norms and ensuring privacy preservation.

*b) Notation:* We use [M] to denote the set of integers $\{1, 2, \ldots, M\}$. We denote by $X$ the set of all possible examples or instances which is also referred to as the input space. The set of all possible labels or target values is denoted by $Y$. The data stored on the client $i$ is denoted by $X_i$, the number of which is denoted by $m_i$. We use $x_i$ to refer to a sample in the dataset $X_i$. We assume that any sample $x_i \in X_i$ is independently and identically distributed according to a fixed but unknown distribution $D_i$. Let $f(x; w) : X \to Y$ be the machine learning model associated with $w$, where $w \in \mathcal{W}$ is

the model parameter to be learned from data. The loss function of a model $f(\cdot; w)$ on some example $x \in X$ also referred to as the risk or true error of $f$ is denoted by $l(f(x; w), y)$ (sometimes dismissing $w$ for simplicity). We use $L_i(w)$ to denote the empirical loss of model $f(\cdot; w)$ on data $X_i$, that is $L_i(w) = \frac{1}{m_i} \sum_{k=1}^{m_i} l(f(x_{i,k}; w), y_{i,k})$. And the optimal value of $L_i(w)$ is denoted by $L_i^*$, that is $L_i^* = \min_{w \in \mathcal{W}} L_i(w)$. If $(X, y) \sim D$, $\mu_X, \mu_y$ denote the expectation of $X$, $y$ separately, $\Sigma_{XX}$ denotes the variance of $X$, $\Sigma_{Xy}$ denotes the covariance of $X$ and $y$, and so on.

The subsequent sections of this paper are organized as follows. Section II presents a description of our clustered FL framework, namely LCFL, accompanied by the introduction of our main theorem. Additionally, implementation considerations regarding our framework are discussed. Supplementary explanations of our theorem are provided in Section III. In Section IV, we outline the limitations associated with employing model parameters. Section V presents the experimental results we obtained, followed by the conclusion of our study in Section VI.

## II. CLUSTERED FEDERATED LEARNING BASED ON LOCAL LOSS

Due to privacy concerns, the propagation of the user's data distribution is not feasible in clustered FL. Therefore, an alternative approach is required to substitute for the data distribution information. Zhu et al. [40] have highlighted the privacy risks associated with using model gradient information, making gradient cosine similarity an unsuitable clustering metric. In Section IV, we elaborate on the limitations of using model parameters as a metric, particularly when the best model parameter is not unique, leading to incorrect classification results. While [39] demonstrated that it is possible to achieve satisfactory clustering performance without a metric function, it is necessary to have a good initial setup. Consequently, a new metric function is needed to classify clients and address the aforementioned scenario effectively.

### A. Algorithm and Main Results

In this subsection, we propose our algorithm framework and some implementation considerations about our framework. The motivation of clustered FL is to reduce the impact of non-IID data by clustering users with similar data distribution into one cluster. For users with very different local data distributions, model parameters trained using only local data will not work on other user data [43]. When $w_i^* \in \arg\min_{w \in \mathcal{W}} L_i(w)$, $w_j^* \in \arg\min_{w \in \mathcal{W}} L_j(w)$, according to optimality, we have

$$\forall w \in \mathcal{W}, \ L_i^* = L_i(w_i^*) \leq L_i(w).$$

If two users' datasets $X_i, X_j$ come from two distributions $D_i, D_j$ which we cannot check whether they are different, there will be two cases. The first we expect is

$$L_i(w_j^*) > L_i(w_i^*). \tag{3}$$

Otherwise the case $L_i(w_j^*) = L_i(w_i^*) = L_i^*$ indicates that

$$w_j^* \in \arg\min_{w \in \mathcal{W}} L_i(w)$$

which means datasets $X_i, X_j$ share the same best model parameter $w_j$. Then, if the two users collaborate on training one model $L_{i+j}(w) = \frac{m_i}{m_i + m_j} L_i(w) + \frac{m_j}{m_i + m_j} L_j(w)$ it will be:

$$L_{i+j}(w_j^*) = \frac{m_i}{m_i + m_j} L_i^* + \frac{m_j}{m_i + m_j} L_j^*$$
$$\leq \frac{m_i}{m_i + m_j} L_i(w) + \frac{m_j}{m_i + m_j} L_j(w), \ \forall w \in \mathcal{W}.$$

which means that $w_j^*$ is still the best model parameter. Under this case, both users $i, j$ will not get a worse model after cooperation. Moreover, because more data are used in the training process, the generalization performance of the model $L_{i+j}(w)$ will be improved. This result shows assigning user $i, j$ in the same cluster is appropriate.

When considering the first case mentioned in equation (3), we assert that the two distributions, $D_i$ and $D_j$, exhibit distinct characteristics. By "different", we mean that these distributions can be discerned by the model hypothesis set $\mathcal{W}$. In the subsequent context of this section, we will further demonstrate and provide evidence for this claim. Drawing inspiration from the aforementioned analyses and the inherent connection between the loss function and distribution, we propose utilizing equation (2) as the clustering metric. In Section III, we present an illustrative example to illustrate the correlation between the loss function and distribution.

We utilize equation (2) as a measure of the "distance" between the data distributions of users. This metric function, by definition, is symmetric and satisfies the triangle inequality for any loss function $L$. While the gradient cosine similarity $\alpha_{i,j}$ exhibits symmetry, the model parameter metric $d_{i,j}$ possesses both symmetry and satisfies the triangle inequality. However, it is important to note that the expression (2) may not always qualify as a true distance. This is because there may exist cases where $d(i, j) = 0$ for $D_i \neq D_j$, which violates the requirement of a distance metric. The same applies to the other two metrics; they are not actual distances either.

We provide a companion algorithmic framework for using this metric function in FL, taking into account both communication and practical considerations. The main idea of the algorithm is to cluster the users before conducting FL and apply the FL algorithm to each cluster separately. We do not require each user to communicate after achieving the optimal solution locally. Instead, we use a parameter called the local iteration count $T$. The local training process ends when the local iteration count reaches the desired value. This preliminary step can be seen as a warm-up phase for FL. The algorithm is formally presented in Algorithm 1.

In the framework of the above clustering distance function $d(i, j)$, the $k$-means algorithm cannot be used, because in this framework, there are only distances between any two participants, and there is no way to choose the "cluster center", that is, the clustering method that requires the identity information of participants is not applicable in our clustering framework. However, a distance metric alone is sufficient for most algorithms, such as improvements to $k$-means. This $k$-medoids method selects the medoid based only on distance, and most hierarchical and density-based clustering methods.

---

**Algorithm 1:** Loss based Clustered Federated Learning(LCFL)

**Input:** initialization $w_i^0, i \in [M]$, number of local iterations $T$, local step size $\gamma$, clustering method $C(\cdot)$, federated learning method FedOpt

1 **for** *Client $i \in [M]$ in parallel* **do**
2    **for** $t = 0$ to $T - 1$ **do**
3      $w_i^{t+1} = w_i^t - \gamma \nabla L_i(w_i^t)$;
4    **end**
5    Sends $w_i^T$ to Server;
6    **for** $j \in [M]/i$ **do**
7      Client $i$ receives $w_j^T$ from Server;
8      $d(i,j)_i = \|L_i(w_j^T) - L_i(w_i^T)\|$;
9      Sends $d(i,j)_i$ to Server;
10    **end**
11 **end**
12 **for** $i, j \in [M], i \neq j$, *Server* **do**
13    $d(i,j) = d(i,j)_i + d(j,i)_j$
14 **end**
15 Server uses $C(d(i,j))$ to obtain clustering structure $C_1, \ldots, C_k$;
16 **for** *Client $k$ in clustering structure $C_k$* **do**
17    $w_k^* =$ FedOpt($C_k$).
18 **end**

**Output:** $C_1, \ldots, C_k, w_1^*, \ldots, w_k^*$

---

Now we turn to prove that if $d(i,j) \geq R$, distributions $D_i, D_j$ can be distinguished by the model hypothesis set $\mathcal{W}$. We will prove this property starting with the measure of the distribution. The divergence between distributions is usually measured by Kullback-Leibler divergence (KL divergence) [44].

**Definition 1** (KL divergence [44]). *For discrete probability distributions $P$ and $Q$ defined on the same sample space $X$ the KL divergence from $Q$ to $P$ is defined as:*

$$D_{KL}(P||Q) = \sum_{x \in X} P(x) \log(\frac{P(x)}{Q(x)}).$$

However, the KL divergence is only related to the distributions and does not consider the machine learning task that this paper cares about most. To measure the difference in data distributions under a machine learning model, we use the concept of "label-discrepancy" between distributions in [45]:

**Definition 2** (Label-discrepancy [45]). *Given a loss function $L$, the discrepancy $disc$ between two distributions $D_i$ and $D_j$ over $X \times Y$ is defined by:*

$$disc_{\mathcal{W}}(D_i, D_j) = \sup_{w \in \mathcal{W}} |L_{D_i}(w) - L_{D_j}(w)|,$$

*where $D_i, D_j$ is the distribution of the feature-label data, $\mathcal{W}$ is the hypothesis set for machine learning models, $L_{D_i}(w)$ denotes the expected loss function when the model is $w$ and the data distribution is $D_i$, that is $L_{D_i}(w) = \mathbb{E}_{(x,y)\sim D_i} l(f(x; w), y)$.*

Under this definition, two distributions $D_i$ and $D_j$ are distinguishable under $\mathcal{W}$ if and only if the label discrepancy between them is 0. If $D_i = D_j$, it is evident that $disc_{\mathcal{W}}(D_i, D_j) = 0$. Conversely, if $disc_{\mathcal{W}}(D_i, D_j) = 0$, it implies that for every $w \in \mathcal{W}$, $L_{D_i}(w) = L_{D_j}(w)$. In this case, the two distributions are not distinguishable under $\mathcal{W}$. This means that the loss of all models in the hypothesis set is the same under both $D_i$ and $D_j$, indicating that models trained on $D_i$ generalize well on $D_j$ and vice versa. Moreover, it is advisable to assign users $i$ and $j$ to the same cluster based on this observation.

We have observed that the expected form of $d(i,j)$, denoted as $\hat{d}(i,j)$, is less than $2 \cdot disc_{\mathcal{W}}(D_i, D_j)$:

$$\hat{d}(i,j) = |L_{D_i}(w_i^*) - L_{D_i}(w_j^*)| + |L_{D_j}(w_i^*) - L_{D_j}(w_j^*)|$$
$$\leq 2 \sup_{w \in \mathcal{W}} |L_{D_i}(w) - L_{D_j}(w)|$$
$$= 2 \cdot disc_{\mathcal{W}}(D_i, D_j).$$

If there exists $\epsilon > 0$ such that $\hat{d}(i,j) \geq \epsilon$, we can conclude

$$0 < \frac{\epsilon}{2} \leq \frac{1}{2}\hat{d}(i,j) \leq disc_{\mathcal{W}}(D_i, D_j). \quad (4)$$

Based on the previous analysis, we can deduce that under (4), models from $\mathcal{W}$ cannot distinguish between the two distributions $D_i$ and $D_j$. Furthermore, the two distributions should be grouped into the same cluster.

However, $d(i,j)$ and $\hat{d}(i,j)$ are not equal in general. To analyze the bound of the two terms, we first introduce the definitions of empirical and average Rademacher complexity [46]. The formal definitions of the empirical and average Rademacher complexity are stated as follows.

**Definition 3** (Empirical Rademacher complexity [46]). *Let $\mathcal{G}$ be a family of functions mapping from $Z$ to $[a, b]$ and $S = (z_1, \ldots, z_m)$ be a fixed sample of size $m$ with elements in $Z$. Then, the empirical Rademacher complexity of $\mathcal{G}$ with respect to the sample $S$ is defined as:*

$$\hat{\mathcal{R}}_S(\mathcal{G}) = \mathbb{E}_{\sigma}\left[\sup_{g \in \mathcal{G}} \frac{1}{m} \sum_{i=1}^{m} \sigma_i g(z_i)\right],$$

*where $\sigma = (\sigma_1, \ldots, \sigma_m)^\top$, with $\sigma_i$ is independent uniform random variables taking values in $\{-1, +1\}$. The random variables $\sigma_i$ are called Rademacher variables.*

**Definition 4** (Rademacher complexity [46]). *Let $D$ denote the distribution according to which samples are drawn. For any integer $m \geq 1$, the Rademacher complexity of $\mathcal{G}$ is the expectation of the empirical Rademacher complexity over all samples of size $m$ drawn according to $D$:*

$$\mathcal{R}_m(\mathcal{G}) = \mathbb{E}_{S \sim D^m}\left[\hat{\mathcal{R}}_S(\mathcal{G})\right].$$

The Rademacher complexity measures the richness of a family of functions by quantifying the extent to which the hypothesis set can fit random noise. For further details on the Rademacher complexity, we refer readers to the book by Mohri et al. [46].

These definitions are provided in the general case of a family of functions $\mathcal{G}$ mapping from an arbitrary input space $Z$ to $\mathbb{R}$.

However, under Assumption 1, we can treat an arbitrary loss function $l : Y \times Y \to \mathbb{R}$ as $\mathcal{G}$, where $Z = X \times Y$:

$$\mathcal{G} = \{g : (x,y) \to l(f(x;w),y) | w \in \mathcal{W}\}.$$

**Assumption 1.** *The loss function $l$ is bounded by some $M > 0$, that is $l(y,y') \leq M$ for all $y, y' \in Y$ or, more strictly, $l(f(x;w),y) \leq M$ for all $w \in \mathcal{W}$ and $(x,y) \in X \times Y$, the problem is referred to as a bounded problem.*

Additionally, we assume that the range of the loss function $l$ is $[0,1]$, so we can omit the explicit mention of $M$ in the following theorem for brevity.

Now, armed with the definitions of empirical and average Rademacher complexity, we can proceed with the analysis of the term to be bounded and give the following theorem.

**Theorem 1.** *For any $\delta > 0$, with probability at least $(1 - \delta)^4$ over the draw of i.i.d samples $X_i, X_j$ of sizes $m_i, m_j$ respectively, the following holds:*

$$|d(i,j) - \hat{d}(i,j)| \leq 2\sqrt{\frac{\log(2/\delta)}{2m_i}} + 2\sqrt{\frac{\log(2/\delta)}{2m_j}} \quad (5)$$
$$+ \mathcal{R}_{m_i}(\mathcal{W}) + \mathcal{R}_{m_j}(\mathcal{W}).$$

We denote

$$C_\delta(X_i, X_j; \mathcal{W}) = 2\sqrt{\frac{\log(2/\delta)}{2m_i}} + 2\sqrt{\frac{\log(2/\delta)}{2m_j}}$$
$$+ \mathcal{R}_{m_i}(\mathcal{W}) + \mathcal{R}_{m_j}(\mathcal{W}).$$

In the context of the Theorem 1, if the empirical discrepancy $d(i,j)$ between users $i$ and $j$ is greater than the bound $C_\delta(X_i, X_j; \mathcal{W})$, then with probability at least $(1 - \delta)^4$, we obtain:

$$\hat{d}(i,j) \geq d(i,j) - C_\delta(X_i, X_j; \mathcal{W}) > 0$$

which indicates that users $i$ and $j$ should be assigned to the same cluster.

In simpler terms, if the observed difference between two users exceeds a certain threshold, the theorem guarantees that the expected difference will also be sufficiently large. This ensures that the clustering algorithm will correctly identify these users as belonging to the same cluster with a high level of confidence.

This result provides a statistical assurance for the performance of the clustering algorithm and the accuracy of its clustering decisions based on pairwise user distances.

Calculating the exact value of the Rademacher complexity can pose significant challenges. In fact, for certain hypothesis sets, computing the empirical Rademacher complexity is known to be NP-hard [46]. As a viable alternative, the VC dimension (Vapnik-Chervonenkis dimension [47]) is frequently employed to provide an upper bound on the Rademacher complexity. While we won't delve into the specific definition of the VC dimension in this article, it is widely utilized in machine learning theory to estimate the complexity of hypothesis classes and analyze the behavior of learning algorithms.

### B. Implementation considerations

In this subsection, we delve into the practical implementation details of our method, highlighting its seamless integration with the existing communication protocol of FL. Furthermore, we will emphasize that our method maintains the privacy of the participating clients throughout the process.

The choice of the metric function $d(i,j)$ in clustering for FL also takes into account the communication cost involved. In this context, it is important to minimize the amount of communication required between the participating clients.

One possible metric function is given by (6):

$$\tilde{d}(i,j) = |L_i(w_i) - L_j(w_i)| + |L_i(w_j) - L_j(w_j)|, \quad (6)$$

where $\tilde{d}(i,j)$ is defined as the difference between the local loss functions of clients $i$ and $j$ under their respective local models. Since exchanging local data is generally prohibited in FL due to privacy concerns, calculating $|L_i(w_i) - L_j(w_i)|$ would require an additional communication step to transmit $L_i(w_i)$.

To address this issue and reduce the communication cost, the metric function $d(i,j)$ (as defined in our algorithm) is preferable. It relies on exchanging model parameters rather than local data. By exchanging the model parameters $w_i$ and $w_j$ between clients, the computation of $d(i,j)$ can be performed locally at each client without the need for additional data exchange. This approach reduces the communication overhead while still allowing the calculation of the pairwise distances necessary for clustering. Considering the goal of minimizing communication cost in FL, the original metric function $d(i,j)$ is a more suitable choice compared to $\tilde{d}(i,j)$, as the former avoids the need for additional data communication.

Our algorithm acts as a warm-up phase before starting the FL process, seamlessly integrating with existing FL frameworks and systems without requiring significant modifications. It follows the standard FL communication protocol, ensuring compatibility and easy adoption in various FL setups. During the warm-up phase, our algorithm establishes a clustering structure among the clients. The local model parameters used for clustering training can also serve as initialization parameters for subsequent FL. This utilization of pre-trained parameters speeds up the FL process, as it provides a well-suited starting point. Thus, any communication costs during the warm-up stage are offset by the subsequent acceleration achieved through the use of pre-trained parameters.

Moreover, privacy preservation holds paramount importance in distributed learning scenarios. With our approach, we place a strong emphasis on protecting the privacy of the participating clients and ensuring the security and confidentiality of their sensitive data. One key aspect of our method is that it avoids the use of gradient information, which can be prone to information leakage.

In our algorithm, we do not explicitly check whether the discrepancy $d(\cdot,\cdot)$ between users $i$ and $j$ exceeds the specified threshold $C_\delta(i,j;\mathcal{W})$ to avoid the calculation of Rademacher complexity. Instead, we utilize the discrepancy matrix obtained from the clustering algorithm to derive the clustering structure. By employing the discrepancy matrix, we leverage the pairwise

discrepancies between users to determine their similarities in the clustering process. This approach provides flexibility in capturing the inherent structure and relationships among the users, allowing the clustering algorithm to autonomously identify and group users based on the discrepancy matrix.

## III. EXAMPLES AND MATHEMATICAL ANALYSIS

In this section, we provide an example to showcase the effectiveness of our algorithm when dealing with user data that exhibits different data distributions. By presenting this example, we aim to illustrate how our algorithm can successfully handle scenarios where the data distribution among users varies. This capability is crucial in clustered FL, as it allows us to leverage diverse and distributed data sources while maintaining robust and accurate model performance.

We analyze our proposed algorithm in a concrete linear model. As mentioned earlier, let us denote the input space as $X$, and the set of target values as $y$. We assume that the data on a certain client are generated from a distribution of $D$. Furthermore, we use the expected squared loss function $L(f(X; w), y) = \mathbb{E}_{(X,y) \sim D} \left( w^\top X - y \right)^2$ instead of the empirical form for the sake of theoretical analysis. The optimizer of this linear model is

$$w^* = (\Sigma_{XX} + \mu_X \mu_X^\top)^{-1}(\Sigma_{Xy} + \mu_y \mu_X),$$

and thus

$$
\begin{aligned}
&L(f(X; \hat{w}), y) - L(f(X; w^*), y) \\
=&(\hat{w} - w^*)^\top \left[ (\Sigma_{XX} + \mu_X \mu_X^\top)(\hat{w} + w^*) - 2(\Sigma_{Xy} + \mu_y \mu_X) \right] \\
\geq& \frac{\lambda_{\min}}{(\hat{\lambda}_{\max} + \hat{\mu}_X^\top \hat{\mu}_X)^2} \|\hat{\Sigma}_{Xy} + \hat{\mu}_X \hat{\mu}_y\|^2 \\
&+ \frac{1}{\lambda_{\max} + \mu_X^\top \mu_X} \|\Sigma_{Xy} + \mu_X \mu_y\|^2 \\
&- \frac{2}{\hat{\lambda}_{\max} + \hat{\mu}_X^\top \hat{\mu}_X} \|\Sigma_{Xy} + \mu_X \mu_y\| \|\hat{\Sigma}_{Xy} + \hat{\mu}_X \hat{\mu}_y\|,
\end{aligned}
\tag{7}
$$

where $\hat{w}$ is the optimizer of this linear model on another data distribution $(X, y) \sim \hat{D}$, $\lambda_{\min}$ denotes the minimum eigenvalue of $\Sigma_{XX}$, $\lambda_{\max}$ denotes the maximum eigenvalue of $\Sigma_{XX}$ and $\hat{\lambda}_{\max}$ denotes the maximum eigenvalue of $\hat{\Sigma}_{XX}$.

If the condition

$$\|\Sigma_{Xy} + \mu_X \mu_y\| \leq \frac{\|\hat{\Sigma}_{Xy} + \hat{\mu}_X \hat{\mu}_y\|}{2(\hat{\lambda}_{\max} + \hat{\mu}_X^\top \hat{\mu}_X)}$$

is satisfied, then

$$
\begin{aligned}
&\frac{\lambda_{\min}}{(\hat{\lambda}_{\max} + \hat{\mu}_X^\top \hat{\mu}_X)^2} \|\hat{\Sigma}_{Xy} + \hat{\mu}_X \hat{\mu}_y\|^2 \\
&+ \frac{1}{\lambda_{\max} + \mu_X^\top \mu_X} \|\Sigma_{Xy} + \mu_X \mu_y\|^2 \\
>& \frac{2}{\hat{\lambda}_{\max} + \hat{\mu}_X^\top \hat{\mu}_X} \|\Sigma_{Xy} + \mu_X \mu_y\| \|\hat{\Sigma}_{Xy} + \hat{\mu}_X \hat{\mu}_y\|
\end{aligned}
$$

holds. It means that the most right-hand side of (7) is greater than zero, hence the gap between $L(f(X; \hat{w}), y)$ and $L(f(X; w^*), y)$ exits. Accordingly using our proposed algorithm, one can obtain a clear distance matrix that is easy for clustering.

Moreover, we assume that all clients' data are normalized, such that $\mu_X = \hat{\mu}_X = \mathbf{0}$, $\mu_y = \hat{\mu}_y = 0$, and $\Sigma_{XX} = \hat{\Sigma}_{XX}$. This implies that the difference between distributions $D$ and $\hat{D}$ only exists in $\Sigma_{Xy}$. From this, we can derive the following expression:

$$
\begin{aligned}
&L(f(X; \hat{w}), y) - L(f(X; w^*), y) \\
=&(\hat{\Sigma}_{Xy} - \Sigma_{Xy})^\top \Sigma_{XX}^{-1} (\hat{\Sigma}_{Xy} - \Sigma_{Xy}).
\end{aligned}
$$

Consequently, we can bound this term as follows:

$$
\begin{aligned}
\frac{1}{\lambda_{\max}} \|\hat{\Sigma}_{Xy} - \Sigma_{Xy}\|^2 &\leq L(f(X; \hat{w}), y) - L(f(X; w^*), y) \\
&\leq \frac{1}{\lambda_{\min}} \|\hat{\Sigma}_{Xy} - \Sigma_{Xy}\|^2.
\end{aligned}
\tag{8}
$$

The above inequality (8) demonstrates the close connection between the loss function and the data distribution. This motivates us to consider equation (2) as a suitable metric.

## IV. DRAWBACKS OF THE MODEL PARAMETER METRIC

In this section, we elaborate that using model parameters as a metric will lead to wrong classification while using model loss functions will not.

The literature [41] states that using the model parameter $w^*$ as the metric of clustering needs to satisfy the following assumption:

**Assumption 2.** *Model parameters of different clusters $w_i^*, w_j^*$ satisfy:*
$$\|w_i^* - w_j^*\| \geq R.$$

However, we illustrate that models trained with the same data distribution can also satisfy the above separation assumption, which will result in users with the same data distribution being assigned to different clusters, which goes against the original intention of integrating similar data to improve model performance.

Consider using the Softmax model in FL. Changing the Softmax model parameter $w \in \mathbb{R}^K$ to $w - \varphi = (w_1 - \varphi, \ldots, w_K - \varphi)^\top$ does not affect the model, since:

$$
\begin{aligned}
p(y_i = j | x_i; w - \varphi) &= \frac{e^{(w_j - \varphi)^T x_i}}{\sum_{l=1}^K e^{(w_l - \varphi)^T x_i}} \\
&= \frac{e^{w_j^T x_i} e^{-\varphi^T x_i}}{\sum_{l=1}^K e^{w_l^T x_i} e^{-\varphi^T x_i}} \\
&= \frac{e^{w_j^T x_i}}{\sum_{l=1}^K e^{w_l^T x_i}} \\
&= p(y_i = j | x_i; w).
\end{aligned}
\tag{9}
$$

(9) means the Softmax model will get the same output under a family of model parameters: $\{w - \varphi | \varphi \in \mathbb{R}\}$. In this situation, even though the model parameters $w^*$ and $\hat{w}^* = w^* + \varphi$ are trained with the same data, they still satisfy the Assumption 2, for any $\varphi$ that satisfies the condition (10).

$$\|w^* - \hat{w}^*\| = \|(\varphi, \ldots, \varphi)^\top\| \geq R. \tag{10}$$

This means the model parameter metric divides users belonging to one cluster into two different clusters, resulting in the loss of data volume under this cluster.

Moreover, considering the Softmax model with regularization term $\|\cdot\|_2$ which is more commonly used in the field of machine learning, if

$$\varphi = \frac{2\sum_{k=1}^{K} w_k^*}{K},$$

such that both $w^*$ and $\hat{w}^* = w^* + \varphi$ are optimal solutions to the regularized Softmax model. In this situation, to cluster users with similar data distribution, it should have at least:

$$R > \|\varphi\| = \|\frac{2\sum_{k=1}^{K} w_k^*}{K}\|.$$

This implies that the choice of the $R$ value significantly impacts the clustering performance. In a broader context, when dealing with models that do not have a unique optimal solution (as is often the case with deep neural network models due to their combinatorial symmetry), it is not suitable to rely on the norm of model parameter differences as the clustering metric.

However, using our proposed metric function, this situation does not occur. The loss function commonly used for training a softmax model is the categorical cross-entropy loss. Mathematically, if we denote the predicted probabilities as $p$ and the true label as $y$ (both vectors of length equal to the number of classes $K$), the categorical cross-entropy loss can be computed as follows:

$$L = -\sum_{i=1}^{K} (y_i \log(p_i)). \tag{11}$$

When the same dataset is used, according to (9), it is obvious that the value of (11) will be the same. It is precisely this property that motivates us to use (2) as a metric function.

## V. Numerical experiments

In this section, we present the experimental results that validate our theoretical analysis. We demonstrate the performance of our proposed metric, as well as other clustering metrics such as the difference of model parameters and the cosine of gradients. Additionally, we compare the performance of our algorithm with IFCA and several baseline methods.

(**Machine**) We carried out the simulation experiments on a commodity workstation with one Intel® Xeon® Gold 5218R CPU with 64 GB RAM and 2 NVidia® 3090 GPUs.

(**Implementation and Hyperparameters**) In order to conduct fair comparisons, Stochastic Gradient Descent (SGD) was adopted as the local problem solver for all FL algorithms involved in this paper. For all numerical experiments, the learning rates on the clients decay every global iteration by a factor of 0.99. The rest of the primary hyperparameters are summarized in Table I. In this table, "participation rate" refers to the proportion of clients that participate in the training process every global iteration, and "init. learning rate" means the initial learning rates of the clients which decay with the pattern stated previously.

The code, including the re-implementations of the FedAvg and IFCA algorithm, was built on top of a simulation framework fl-sim for FL based on PyTorch [48]. More details can be found in https://github.com/wenh06/fl-sim and https://github.com/wenh06/LCFL.

TABLE I
PRIMARY HYPERPARAMETERS OF THE NUMERICAL EXPERIMENTS
CONDUCTED IN THIS PAPER

| hyperparameter | FEMNIST | Rotated MNIST | Rotated CIFAR10 |
|---|---|---|---|
| participation rate | 30% | 100% | 10% |
| init. learning rate | 0.03 | 0.02 | 0.02 |
| batch size | 20 | 20 | 20 |
| local epochs | 5 | 10 | 5 |
| global iterations | 100 | 60 | 200 |

We use three federated datasets to verify the performance: Federated Extended MNIST(FEMNIST) [49], Rotated MNIST and Rotated CIFAR10 [39].

In the case of the FEMNIST dataset, we utilized a sub-sampled and repartitioned version of the original full dataset to introduce additional statistical heterogeneity. This dataset has also been employed in previous studies such as [9], [12]. Specifically, it followed the procedure of selecting 10 lowercase characters ('a'-'j') from the EMNIST dataset [50] and distributing only 3 classes to each device. The dataset consists of 200 devices, and we used a convolutional neural network model to classify the inputs, which were gray-scale images with 784 pixels and flattened to 1-dimensional.

The Rotated MNIST and Rotated CIFAR10 datasets were created based on the MNIST [51] and CIFAR10 datasets [52] by [39]. These datasets were designed to simulate a FL scenario where the data distribution varies among different clients. For Rotated MNIST, the data pictures were rotated by four degrees: 0, 90, 180, and 270, resulting in the creation of $k = 4$ clusters with a relatively straightforward cluster structure. Considering that the original MNIST dataset contains 60,000 training images and 10,000 test images, we randomly partitioned the $60,000k$ training images and $10,000k$ test images among $m$ clients, ensuring each client had the same rotation. Rotated CIFAR10 was created similarly to Rotated MNIST, except that the rotation degrees were 0 and 180, resulting in $k = 2$ clusters, and we set the number of clients to $m = 200$, as in [39].

### A. Performances of the metrics

In this subsection, we compared the performance of three different metrics which use model loss functions, model parameters, and the cosine of gradients. We utilized our proposed framework, using the same federated optimizer FedAvg for fairness. However, the metrics used in the clustering stage are different. We performed experiments on two datasets: FEMNIST and Rotated MNIST.

The results are shown in Figure 2 and Figure 3. In these figures, "LCFL" is our algorithm. "DiffNorm" and "Gradcos" refer to the metrics using model parameters and the cosine of gradients respectively. We conducted experiments for $m = 1200$ clients on Rotated MNIST. We repeated the experiment five times using five different random seeds and computed the standard deviation of the output accuracy with the "±STD" legend. We use average accuracy to measure the effect and "Global Iter." stands for the global iteration number.

In Figure 2, these three models ultimately achieved an average accuracy of over 90%, but the number of iterations
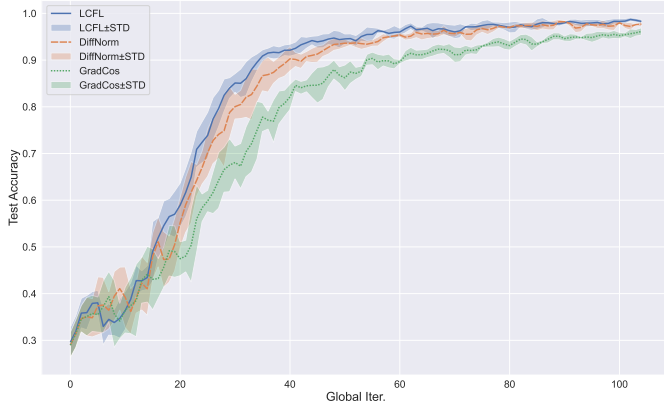
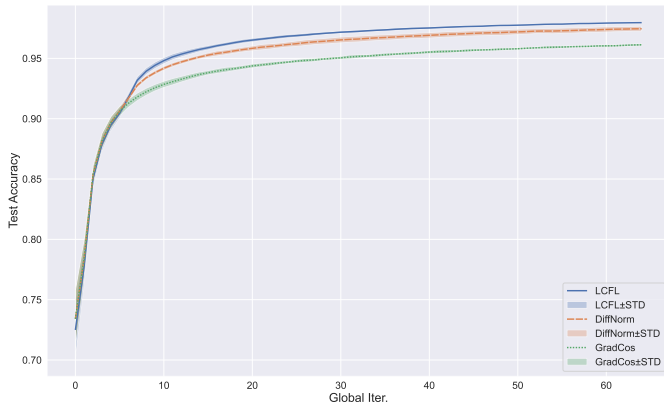Fig. 2. Performances of different metrics on FEMNIST



Fig. 3. Performances of different metrics on Rotated MNIST

data, resulting in improved clustering performance.

### B. Performances of algorithms

In this subsection, we compare our proposed algorithm LCFL with IFCA and FedAvg. Moreover, we conduct local training to verify the necessity of FL. We conducted experiments on three datasets: FEMNIST, Rotated CIFAR10, and Rotated MNIST. For the experiments on Rotated MNIST, we conducted for $m = 1200$ and $m = 2400$ clients, respectively.

For LCFL, we use the $k$-medoids to cluster users and set $k = 10, T = 10$. For IFCA, we also set $k = 10$. For all algorithms, we run the experiments with 5 random seeds and report the average test accuracy and standard deviation in Figure 4, Figure 5, Table II, and Table III. The "Local" legend refers to the local training, calculating the average test set accuracy of all local models.
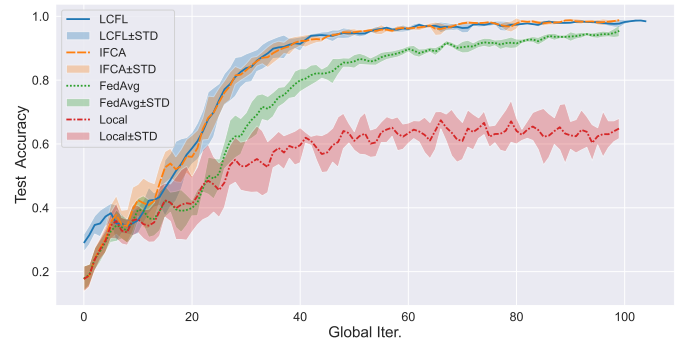


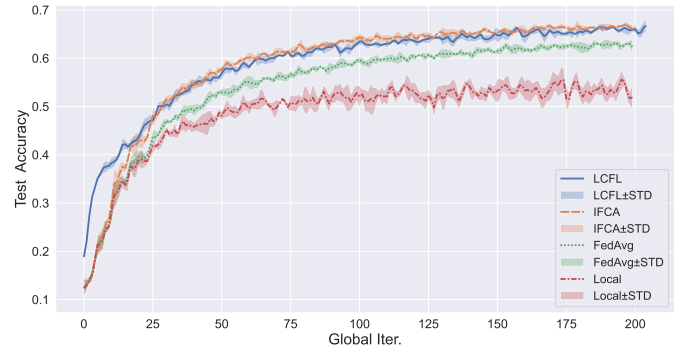Fig. 4. Test Accuracy comparison on FEMNIST

they took to achieve 90% accuracy varied greatly. Our model achieved this result around the 37th iteration (considering standard deviation), while the other two measurement criteria were around the 40th iteration and the 60th iteration, respectively. On the results after 100 iterations, there is an approximate result of using the model loss functions and model parameters as metrics, which is about 98%. However, using the cosine of gradients as the metric result is slightly worse, at around 95%. The standard deviation of accuracy gradually decreases as the global iteration process progresses.

In Figure 3, the results are similar to the previous experiment shown in Figure 2, and using our proposed metric is superior to the other two metrics. In this experiment, the standard deviation of the average accuracy is very small. The clustering federated learning algorithm using three different metrics ultimately achieved an average accuracy of over 95%, achieved in the 11th, 14th, and 27th iterations, respectively. After 60 iterations, the three metrics achieved results of approximately 96%, 97%, and 97.5%, respectively.

The metric we proposed and utilized in the LCFL algorithm is unambiguous, leading to the best performance compared to other metrics. This clear advantage is highlighted in the two figures. The unambiguous nature of our metric ensures that the clustering process in the LCFL algorithm is robust and effective. By leveraging this metric, we can accurately capture the underlying patterns and similarities among the distributed



Fig. 5. Test Accuracy comparison on Rotated CIFAR10

TABLE II
TEST ACCURACIES(%) $\pm$ STD(%) ON ROTATED MNIST($m = 1200$)

| Iteration | LCFL | IFCA | FedAvg | Local |
|---|---|---|---|---|
| 5 | $89.66 \pm 0.57$ | $91.13 \pm 0.51$ | $89.66 \pm 0.57$ | $74.82 \pm 1.86$ |
| 10 | $94.45 \pm 0.31$ | $94.04 \pm 0.23$ | $92.52 \pm 0.31$ | $79.54 \pm 2.89$ |
| 15 | $95.75 \pm 0.22$ | $95.14 \pm 0.19$ | $93.64 \pm 0.24$ | $84.33 \pm 3.22$ |
| 30 | $97.11 \pm 0.12$ | $96.48 \pm 0.09$ | $95.03 \pm 0.17$ | $84.27 \pm 3.60$ |
| 60 | $97.90 \pm 0.08$ | $97.42 \pm 0.09$ | $96.12 \pm 0.14$ | $84.33 \pm 3.75$ |

From the results of these experiments, our algorithm has demonstrated superior performance compared to the two baseline approaches, FedAvg and local training, and is slightly

| Iteration | LCFL | IFCA | FedAvg | Local |
|-----------|------|------|--------|-------|
| 5 | $81.05 \pm 1.13$ | $85.48 \pm 1.39$ | $81.05 \pm 1.13$ | $78.66 \pm 1.34$ |
| 10 | $91.05 \pm 1.11$ | $90.98 \pm 0.96$ | $87.51 \pm 0.50$ | $82.67 \pm 1.55$ |
| 15 | $93.19 \pm 1.05$ | $92.77 \pm 0.82$ | $89.82 \pm 0.38$ | $84.12 \pm 1.67$ |
| 30 | $95.30 \pm 0.84$ | $94.95 \pm 0.61$ | $92.53 \pm 0.27$ | $85.85 \pm 1.86$ |
| 60 | $96.65 \pm 0.64$ | $96.41 \pm 0.50$ | $94.30 \pm 0.20$ | $87.01 \pm 1.96$ |

better than `IFCA`. By implementing the `LCFL` algorithm, we can proactively uncover the underlying cluster identities of the clients. Once the correct cluster identity is discovered, each model is trained and tested using data from a similar distribution, resulting in improved accuracy. In contrast, the `FedAvg` baseline attempts to fit all data from different distributions, which limits its ability to make personalized predictions, leading to inferior performance. The local model baseline algorithm is prone to overfitting, as proved by the gradually increasing standard deviation.

We can also see the advantages of clustering federated learning algorithms from our experimental results. As the test accuracy of `FedAvg` after 60 global iterations is less than clustered `FL` methods(`LCFL` and `IFCA`).

Furthermore, the results provide compelling evidence for the necessity of conducting `FL` tasks. They clearly illustrate that `FL` models outperform local models by a significant margin. This emphasizes the importance of collaborative learning and leveraging the collective knowledge of distributed devices to enhance the performance and accuracy of machine learning models, demonstrating the superiority of `FL` models over local models.

In summary, the clarity and effectiveness of our proposed metric, as well as the compelling evidence presented in the figures, support the superiority of `FL` models over local models and emphasize the importance of conducting `FL` tasks.

## VI. CONCLUSION

In this article, we presented `LCFL`, a novel framework for clustered `FL` that aims to improve existing `FL` methods by enabling the participating clients to learn more specialized models through clustering. By leveraging our finding that the similarity between clients' loss functions is highly indicative of the similarity of their data distributions, `LCFL` overcomes the limitations associated with using the similarity between weight updates.

To recap, `LCFL` offers several advantages over previous clustered `FL` approaches [38], [41], [39], [43], [42]. Firstly, it does not require the transfer of gradient data, making it particularly suitable for privacy-preserving situations. Secondly, `LCFL` can effectively differentiate between scenarios where learning a single model from all clients' data is feasible and situations where it is not practical, and only segregates clients in the latter scenario. Moreover, the clustering step in `LCFL` is performed as a warm-up, incorporating a classic `FL` process.

Our experiment verified that using the value of loss functions as a metric has a better effect on clustering compared to the similarity between gradient cosine and that between weight update. In addition, our experiments on convolutional deep neural networks demonstrate that `LCFL` can achieve significant advancements over the `FL` baseline in terms of classification accuracy. This is particularly evident when all clients' data exhibit a clustering structure.

While this work has made significant strides in assessing the relationship between model performance and data distribution in `FL`, there are still areas for further investigation. Future research should delve into broader scenarios, conduct detailed analyses, and address potential limitations to provide a more comprehensive understanding of this field.

In conclusion, the `LCFL` framework presents a promising approach for clustered `FL`, with potential applications in various domains. By addressing privacy concerns, enhancing clustering accuracy, and improving classification performance, `LCFL` contributes to the advancement of `FL` methods and paves the way for more efficient and specialized learning in data-driven environments.

## REFERENCES

[1] J. Albrecht, "How the GDPR will change the world," *European Data Protection Law Review*, vol. 2, no. 3, pp. 287–289, 2016.

[2] B. McMahan, E. Moore, D. Ramage *et al.*, "Communication-efficient learning of deep networks from decentralized data," in *Proceedings of the 20th International Conference on Artificial Intelligence and Statistics*, vol. 54, 2017, pp. 1273–1282.

[3] J. Konečnỳ, H. B. McMahan, D. Ramage *et al.*, "Federated optimization: Distributed machine learning for on-device intelligence," *arXiv preprint arXiv:1610.02527*, 2016.

[4] W. Y. B. Lim, N. C. Luong, D. T. Hoang *et al.*, "Federated learning in mobile edge networks: A comprehensive survey," *IEEE Communications Surveys & Tutorials*, vol. 22, no. 3, pp. 2031–2063, 2020.

[5] T. Li, S. Hu, A. Beirami *et al.*, "Ditto: Fair and robust federated learning through personalization," in *International Conference on Machine Learning*, 2021, pp. 6357–6368.

[6] P. Kairouz, H. B. McMahan, B. Avent *et al.*, "Advances and open problems in federated learning," *Foundations and Trends in Machine Learning*, vol. 14, no. 1–2, pp. 1–210, 2021.

[7] T. Li, A. K. Sahu, A. Talwalkar *et al.*, "Federated learning: Challenges, methods, and future directions," *IEEE Signal Processing Magazine*, vol. 37, no. 3, pp. 50–60, 2020.

[8] Q. Yang, Y. Liu, T. Chen *et al.*, "Federated machine learning: Concept and applications," *ACM Transactions on Intelligent Systems and Technology*, vol. 10, no. 12, pp. 1–19, 2019.

[9] T. Li, A. K. Sahu, M. Zaheer *et al.*, "Federated optimization in heterogeneous networks," in *Proceedings of Machine Learning and Systems*, vol. 2, 2020, pp. 429–450.

[10] R. Pathak and M. J. Wainwright., "Fedsplit: An algorithmic framework for fast federated optimization," in *Advances in Neural Information Processing Systems*, vol. 33, 2020, pp. 7057–7066.

[11] X. Zhang, M. Hong, S. Dhople *et al.*, "A federated learning framework with adaptivity to non-iid data," *IEEE Transactions on Signal Processing*, vol. 69, pp. 6055–6070, 2021.

[12] Q. T. Dinh, D. P. N. Pham, and L. Nguyen, "FedDR-randomized Douglas-Rachford splitting algorithms for nonconvex federated composite optimization," in *Advances in Neural Information Processing Systems*, vol. 34, 2021, pp. 30 328–30 338.

[13] F. Hanzely and P. Richtárik, "Federated learning of a mixture of global and local models," *arXiv preprint arXiv:2002.05516*, 2020.

[14] D. W. Peaceman and H. H. Rachford, Jr, "The numerical solution of parabolic and elliptic differential equations," *Journal of the Society for Industrial and Applied Mathematics*, vol. 3, no. 1, pp. 28–41, 1955.

[15] J. Douglas and H. H. Rachford, "On the numerical solution of heat conduction problems in two and three space variables," *Transactions of the American Mathematical Society*, vol. 82, no. 2, pp. 421–439, 1956.

[16] E. Bagdasaryan, A. Veit, Y. Hua *et al.*, "How to backdoor federated learning," in *Proceedings of the 23rd International Conference on Artificial Intelligence and Statistics*, vol. 108, 2020, pp. 2938–2948.

[17] A. N. Bhagoji, S. Chakraborty, P. Mittal *et al.*, "Analyzing federated learning through an adversarial lens," in *Proceedings of the 36th International Conference on Machine Learning*, vol. 97, 2019, pp. 634–643.

[18] L. P. Barnes, Y. Han, and A. Özgür, "Lower bounds for learning distributions under communication constraints via fisher information," *The Journal of Machine Learning Research*, vol. 21, no. 1, pp. 9583–9612, 2020.

[19] S. Patel, G. Persiano, and K. Yeo, "Private stateful information retrieval," in *Proceedings of the 2018 ACM SIGSAC Conference on Computer and Communications Security*, 2018, p. 1002–1019.

[20] S. Augenstein, H. B. McMahan, D. Ramage *et al.*, "Generative models for effective ML on private, decentralized datasets," *arXiv preprint arXiv:1911.06679*, 2019.

[21] Q. Li, Z. Wen, Z. Wu *et al.*, "A survey on federated learning systems: Vision, hype and reality for data privacy and protection," *IEEE Transactions on Knowledge and Data Engineering*, vol. 35, no. 4, pp. 3347–3366, 2023.

[22] Z.-H. Zhou, *Machine learning*. Springer Nature, 2021.

[23] A. K. Jain and R. C. Dubes, *Algorithms for clustering data*. Prentice-Hall, Incorporated, 1988.

[24] A. K. Jain, "Data clustering: 50 years beyond $k$-means," *Pattern Recognition Letters*, vol. 31, no. 8, pp. 651–666, 2010.

[25] L. Kaufman and P. J. Rousseeuw, "Clustering by means of medoids," in *Data Analysis Based on $L_1$ Norm and Related Methods*, vol. 31, 1987, pp. 405–416.

[26] D. Pelleg and A. W. Moore, "X-means: Extending k-means with efficient estimation of the number of clusters," in *International Conference on Machine Learning*, vol. 1, 2000, pp. 727–734.

[27] R. Tibshirani, G. Walther, and T. Hastie, "Estimating the number of clusters in a data set via the gap statistic," *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, vol. 63, no. 2, pp. 411–423, 2001.

[28] M. Ester, H. P. Kriegel, J. Sander *et al.*, "A density-based algorithm for discovering clusters in large spatial databases with noise," in *Proceedings of the 2nd International Conference on Knowledge Discovery and Data Mining*, vol. 96, no. 34, 1996, pp. 226–231.

[29] M. Ankerst, M. M. Breunig, H. P. Kriegel *et al.*, "Optics: Ordering points to identify the clustering structure," *Association for Computing Machinery Special Interest Group on Management of Data Record*, vol. 28, no. 2, pp. 49–60, 1999.

[30] A. Hinneburg and D. A. Keim, "An efficient approach to clustering in large multimedia databases with noise," in *Proceedings of the 4th International Conference on Knowledge Discovery and Data Mining*, 1998, pp. 58–65.

[31] L. Kaufman and P. J. Rousseeuw, *Finding groups in data: an introduction to cluster analysis*. John Wiley & Sons, 2009.

[32] T. Zhang, R. Ramakrishnan, and M. Livny, "Birch: an efficient data clustering method for very large databases," *Association for Computing Machinery Special Interest Group on Management of Data Record*, vol. 25, no. 2, pp. 103–114, 1996.

[33] S. Guha, R. Rastogi, and K. Shim, "Rock: A robust clustering algorithm for categorical attributes," *Information Systems*, vol. 25, no. 5, pp. 345–366, 2000.

[34] C. Fraley and A. E. Raftery, "Model-based clustering, discriminant analysis, and density estimation," *Journal of the American Statistical Association*, vol. 97, no. 458, pp. 611–631, 2002.

[35] R. Xu and D. Wunsch, "Survey of clustering algorithms," *IEEE Transactions on Neural Networks*, vol. 16, no. 3, pp. 645–678, 2005.

[36] G. Gan, C. Ma, and J. Wu, *Data clustering: theory, algorithms, and applications*. SIAM, 2020.

[37] A. Fahad, N. Alshatri, Z. Tari *et al.*, "A survey of clustering algorithms for big data: Taxonomy and empirical analysis," *IEEE Transactions on Emerging Topics in Computing*, vol. 2, no. 3, pp. 267–279, 2014.

[38] F. Sattler, K. R. Muller, and W. Samek, "Clustered federated learning: Model-agnostic distributed multitask optimization under privacy constraints," *IEEE Transactions on Neural Networks and Learning Systems*, vol. 32, no. 8, pp. 3710–3722, 2021.

[39] A. Ghosh, J. Chung, D. Yin *et al.*, "An efficient framework for clustered federated learning," in *Advances in Neural Information Processing Systems*, vol. 33, 2020, pp. 19 586–19 597.

[40] L. Zhu, Z. Liu, and S. Han, "Deep leakage from gradients," in *Advances in Neural Information Processing Systems*, vol. 32, 2019, pp. 3710–3722.

[41] A. Ghosh, J. Hong, D. Yin *et al.*, "Robust federated learning in a heterogeneous environment," *arXiv preprint arXiv:1906.06629*, 2019.

[42] Y. Xu, Z. Liu, Y. Sun *et al.*, "Clustered federated learning based on nonconvex pairwise fusion," *arXiv preprint arXiv:2211.04218*, 2022.

[43] Y. Mansour, M. Mohri, J. Ro *et al.*, "Three approaches for personalization with applications to federated learning," *arXiv preprint arXiv:2002.10619*, 2020.

[44] S. Kullback and R. A. Leibler, "On information and sufficiency," *The Annals of Mathematical Statistics*, vol. 22, no. 1, pp. 79–86, 1951.

[45] M. Mohri and A. Muñoz Medina, "New analysis and algorithm for learning with drifting distributions," in *Algorithmic Learning Theory*, 2012, pp. 124–138.

[46] M. Mohri, A. Rostamizadeh, and A. Talwalkar, *Foundations of machine learning*. MIT Press, 2018.

[47] V. N. Vapnik and A. Y. Chervonenkis, "Uniform convergence of relative frequencies of events to their probabilities theory of pobility and its applications," *Theory of Probability & Its Applications*, vol. 16, no. 2, pp. 264–280, 1971.

[48] A. Paszke, S. Gross, F. Massa *et al.*, "Pytorch: An imperative style, high-performance deep learning library," *Advances in Neural Information Processing Systems*, vol. 32, pp. 8024–8035, 2019.

[49] S. Caldas, S. M. K. Duddu, P. Wu *et al.*, "Leaf: A benchmark for federated settings," *arXiv preprint arXiv:1812.01097v3*, 2018.

[50] G. Cohen, S. Afshar, J. Tapson *et al.*, "EMNIST: Extending MNIST to Handwritten Letters," in *2017 International Joint Conference on Neural Networks*, 2017, pp. 2921–2926.

[51] Y. Lecun, L. Bottou, Y. Bengio *et al.*, "Gradient-Based Learning Applied to Document Recognition," *Proceedings of the IEEE*, vol. 86, no. 11, pp. 2278–2324, 1998.

[52] A. Krizhevsky, G. Hinton *et al.*, "Learning multiple layers of features from tiny images," University of Toronto, Tech. Rep., 2009.

## APPENDIX

**Theorem 2** (Hoeffding's Inequality). *Let $X_1, X_2, \ldots, X_n$ be independent and identically distributed (i.i.d.) random variables, where $X_i \in [a_i, b_i]$ for all $i$. Define the sample mean as $\bar{X}_n = \frac{1}{n}\sum_{i=1}^{n} X_i$. Then, for any $\epsilon > 0$, the following inequality holds:*

$$P\left(\left|\bar{X}_n - \mathbb{E}(\bar{X}_n)\right| \geq \epsilon\right) \leq 2e^{-2n\epsilon^2 / \sum_{i=1}^{n}(b_i-a_i)^2}$$

*where $\mathbb{E}(\bar{X}_n)$ denotes the expected value of $\bar{X}_n$.*

**Proof of Theorem 1**

*Proof.* The part to be bounded, denoted as $|d(i,j) - \hat{d}(i,j)|$, is decomposed into four terms by breaking down the absolute value:

$$\begin{aligned}
&|d(i,j) - \hat{d}(i,j)| \\
&\leq |L_{D_i}(w_j^*) - L_i(w_j^*)| + |L_i(w_i^*) - L_{D_i}(w_i^*)| \quad (12)\\
&+ |L_{D_j}(w_i^*) - L_j(w_i^*)| + |L_j(w_j^*) - L_{D_j}(w_j^*)|.
\end{aligned}$$

For the first part of the right hand of (12), the inequality is established using Hoeffding's inequality. By applying Hoeffding's inequality to the difference between the empirical and expected loss on samples drawn from $D_i$, a lower bound is obtained. This bound is given by (13).

$$\mathbb{P}_{X_i \sim D_i^{m_i}}\left[|L_{D_i}(w_j^*) - L_i(w_j^*)| \geq \epsilon\right] \leq 2\exp^{-2m_i\epsilon^2}. \quad (13)$$

By setting the right-hand side of 13 to be equal to $\delta$ and solving for $\epsilon$, a new inequality is derived in the following:

$$\mathbb{P}_{X_i \sim D_i^{m_i}}\left[|L_{D_i}(w_j^*) - L_i(w_j^*)| \leq \sqrt{\frac{\log(2/\delta)}{2m_i}}\right] \geq 1 - \delta, \quad (14)$$

Similarly, Hoeffding's inequality is applied to the difference between the empirical and expected loss on samples drawn from $D_j$, leading to the following two inequalities.

$$\mathop{\mathbb{P}}_{X_j \sim D_j^{m_j}} \left[ |L_{D_j}(w_i^*) - L_j(w_i^*)| \geq \epsilon \right] \leq 2 \exp^{-2m_j \epsilon^2},$$

and

$$\mathop{\mathbb{P}}_{X_j \sim D_j^{m_j}} \left[ |L_{D_j}(w_i^*) - L_j(w_i^*)| \leq \sqrt{\frac{\log(2/\delta)}{2m_j}} \right] \geq 1 - \delta. \tag{15}$$

According to Theorem 3.3 in [46], the remaining two terms can be bounded using the generalization bounds in terms of Rademacher complexity. The probabilities of the differences between the true loss and expected loss being larger than a certain value are upper-bounded.

$$\mathop{\mathbb{P}}_{X_i \sim D_i^{m_i}} \left[ |L_{D_i}(w_i^*) - L_i(w_i^*)| \geq \epsilon \right] \leq 2 \exp^{-2m_i[\epsilon - \mathcal{R}_{m_i}(\mathcal{W})]^2},$$

$$\mathop{\mathbb{P}}_{X_j \sim D_j^{m_j}} \left[ |L_{D_j}(w_j^*) - L_j(w_j^*)| \geq \epsilon \right] \leq 2 \exp^{-2m_j[\epsilon - \mathcal{R}_{m_j}(\mathcal{W})]^2},$$

We do the same transformation to get

$$\mathop{\mathbb{P}}_{X_i \sim D_i^{m_i}} \left[ |L_{D_i}(w_i^*) - L_i(w_i^*)| \leq \sqrt{\frac{\log(2/\delta)}{2m_i}} + \mathcal{R}_{m_i}(\mathcal{W}) \right]$$
$$\geq 1 - \delta \tag{16}$$

and

$$\mathop{\mathbb{P}}_{X_j \sim D_j^{m_j}} \left[ |L_{D_j}(w_j^*) - L_j(w_j^*)| \leq \sqrt{\frac{\log(2/\delta)}{2m_j}} + \mathcal{R}_{m_j}(\mathcal{W}) \right]$$
$$\geq 1 - \delta. \tag{17}$$

By combining the inequalities (12), (14),(15),(16) and (17), the desired bound in (5) is obtained. The probabilities of all the terms being within their respective bounds are multiplied, resulting in the final inequality.

$$\mathop{\mathbb{P}}_{\substack{X_i \sim D_i^{m_i} \\ X_j \sim D_j^{m_j}}} \left[ |d(i,j) - \hat{d}(i,j)| \leq 2\sqrt{\frac{\log(2/\delta)}{2m_i}} + 2\sqrt{\frac{\log(2/\delta)}{2m_j}} \right.$$
$$\left. + \mathcal{R}_{m_i}(\mathcal{W}) + \mathcal{R}_{m_j}(\mathcal{W}) \right] \geq (1 - \delta)^4$$

The final inequality guarantees that with probability at least $(1 - \delta)^4$, the bound in 5 holds. □

Overall, this proof establishes a probabilistic bound on the difference between the target expected discrepancy $\hat{d}(i,j)$ and its empirical estimate $d(i,j)$, considering the individual terms and their respective probabilities.

**Endong Gu** received his B.S. degree in Statistics from Beihang University, Beijing, China, in 2021. Currently, he is pursuing his M.S. degree in Applied Mathematics at the School of Mathematical Sciences, Beihang University, Beijing, China. Endong's research interests primarily revolve around federated learning and the application of optimization problems in the field of machine learning.

**Yongxin Chen** received his master's degree from the School of Mathematical Sciences, Nanjing Normal University. He is currently a Ph.D. candidate at the School of Mathematical Sciences, Beihang University, pursuing advanced studies in optimization theory and methods. His research interests include sparse optimization, distributed optimization, and nonconvex optimization.

**Hao Wen** received his Ph.D. degree in Mathematics in 2018 from Tsinghua University, Beijing, China. He worked as a post-doctoral fellow at Beihang University under the supervision of Professor Deren Han. He is currently a lecturer in the Department of Applied Mathematics in the College of Science at China Agricultural University. His primary research interests include theory, algorithm, and application of optimization; machine learning, especially federated learning, and its applications in medical and agricultural sciences.

**Xingju Cai** is a Professor, and doctoral supervisor at the School of Mathematical Sciences, Nanjing Normal University. She is currently the deputy secretary-general of the Chinese Operations Research Society and the chairman of the Jiangsu Operations Research Society. She has served as a visiting editor of the Asia-Pacific Journal of Operational Research and Operations Research Transactions. She received her Ph.D. degree in Computational Mathematics from Nanjing University in 2013. Her research mainly focuses on machine learning, first-order algorithms for convex and non-convex optimization problems, and their computational complexity and applications. She has published over 30 academic papers in related fields and has made outstanding achievements in structural optimization and variational inequality decomposition algorithms. She was awarded the first prize for Science and Technology Progress of Jiangsu Province in 2021.

**Deren Han** is a Professor, doctoral supervisor, the dean of the School of Mathematical Sciences, Beihang University. He is also the secretary general of the Mathematics Education Steering Committee of the Ministry of Education. He received his PhD degree in computational mathematics from Nanjing University in 2002. His research mainly focuses on numerical methods for large-scale optimization and variational inequality problems, as well as the application of optimization and variational inequality problems in transportation planning and magnetic resonance imaging. He has published many academic papers. He has won the Youth Operations Research Award of the Operations Research Society of China, the Second prize for Science and Technology Progress of Jiangsu Province, and other awards. He has presided over the National Science Fund for Distinguished Young Scholars and other projects. He serves as the standing director of the Operations Research Society of China and the president of the Operations Research Society of Jiangsu Province. Han also is an editorial board member of Numerical Computation and Computer Applications, the Journal of the Operations Research Society of China, Journal of Global Optimization.