

FedCos: A Scene-adaptive Federated Optimization Enhancement for Performance Improvement

Hao Zhang*, Tingting Wu*, Siyao Cheng*, and Jie Liu†

*Harbin Institute of Technology, Harbin, China

†Harbin Institute of Technology (Shenzhen), Shenzhen, China

Email: zhh1000@hit.edu.cn, ttwu@ir.hit.edu.cn, {csy, jieliu}@hit.edu.cn

Abstract—As an emerging technology, federated learning (FL) involves training machine learning models over distributed edge devices, which attracts sustained attention and has been extensively studied. However, the heterogeneity of client data severely degrades the performance of FL compared with that in centralized training. It causes the locally trained models of clients to move in different directions. On the one hand, it slows down or even stalls the global updates, leading to inefficient communication. On the other hand, it enlarges the distances between local models, resulting in an aggregated global model with poor performance. Fortunately, these shortcomings can be mitigated by reducing the angle between the directions that local models move in. Based on this fact, we propose FedCos, which reduces the directional inconsistency of local models by introducing a cosine-similarity penalty. It promotes the local model iterations towards an auxiliary global direction. Moreover, our approach is auto-adapt to various non-IID settings without an elaborate selection of hyperparameters. The experimental results show that FedCos outperforms the well-known baselines and can enhance them under a variety of FL scenes, including varying degrees of data heterogeneity, different number of participants, and cross-silo and cross-device settings. Besides, FedCos improves communication efficiency by 2 to 5 times. With the help of FedCos, multiple FL methods require significantly fewer communication rounds than before to obtain a model with comparable performance.

Index Terms—Edge computing, federated learning, data heterogeneity, performance improvement, communication efficiency

I. INTRODUCTION

With the proliferation of sensing devices, a new era of Internet of Things (IoT) is sparked. These distributed devices generate significant amounts of data all the time, which promotes artificial intelligence into our life, such as smart healthcare system, automatic driving, smart city, etc. Traditionally, to reap the benefits of data of edge devices, the predominant approach is to collect all the data to the remote central cloud for processing and modeling. However, with the rapid development of IoT applications, transmitting the generated data results in high communication cost. Moreover, uploading the data may impose great privacy leakage risk. Under the limitation of legislation such as General Data Protection Regulation (GDPR) [1], training the deep learning model centrally by gathering data from users is impractical.

Edge computing is proposed to shift more computation to the network edge, allowing the edge devices to train models locally. However, insufficient data samples and local data shifts

would lead to a worse model. With the landing of federated learning (FL), training deep learning models in parallel with the edge nodes becomes achievable. FL is a distributed computing paradigm that multiple remote edge devices collaboratively train a global model without exchanging their local data. It treats the collaborative edge devices as working clients, training a machine learning model by local data and synchronizing the parameters via the parameter server. Since model parameters instead of raw data are transitted in the training process, the risk of privacy leakage is greatly reduced and the problem of communication overhead is alleviated¹.

As an emerging machine-learning technique, FL is still in the early stages of research [2]. Compared with the traditional distributed machine learning, FL faces the challenge of data heterogeneity caused by the limitation of data transferring. Specifically, for the traditional distributed machine learning, where the training data of clients is sampled in IID (independently and identically distributed) way, the stochastic gradients of local models are unbiased estimates of the full gradients [3], [4]. In this case, all the clients have roughly the same optimal target. The local models move in the same direction. Therefore, the performance is almost identical with the centralized methods even if the local models are synchronized after multiple local iterations [5]–[7]. On the contrary, accuracy significantly degrades under more widespread non-IID data distributions [8], [9]. The heterogeneity of client data has been deemed as a pivotal factor suppressing the performance [10] since the standard algorithm FedAvg [11] is proposed. In this situation, the condition of unbiased estimation is no longer met. Although it can be alleviated by reducing the number of local iterations (e.g., as an extreme case, all the local models are iterated only once in each communication round), intolerable communication overhead would be introduced. Therefore, *how to enhance the learning performance with limited communication resources is a foundational goal of FL*.

Many previous works has been done to try to improve the performance of FL in a variety of aspects. Among them, plenty of studies focus on local training [12]–[14]. For instance, FedProx [12] adds a proximal term to restrict the local model not far from the current global one. FedMMD [13] has the same goal but by making their output distribution similar.

¹Generally, in IoT scenarios, the number of model parameters is relatively small compared to the massive amount of raw data continuously generated.

But they may slow down model updates since they enforce the local models close to the stale model. Besides, from the experimental study [15], FedAvg still performs best in many kinds of FL scenes. Other approaches attempt to improve aggregation scheme [16]–[19], which either requires additional public data for model distillation [17] or introduces expensive training costs for obtaining sufficient model samples [16], which is not suitable for edge devices. Some works [20], [21] attempt to speed up training to reduce the communication rounds, but they have no effect on improving performance or even hurt performance [15]. Thus, in practice FedAvg still is the widely accepted one.

In the non-IID data scenarios, the local model of each client updates iteratively to its local optimum based on the data itself. During the local iterations, local models move in diverse directions, which causes the two following problems. Firstly, the gains from local training would be offset by the aggregation of local models, which slows down or even stalls the global updates to *lead to inefficient communication*. Secondly, local models are far away from each other, which causes the aggregated model distant from all local models, *leading to poor model performance* on all local data. What’s worse, to reduce the communication cost, the clients usually perform multiple SGD steps before aggregation, which enlarges the distances further. From our investigation, these shortcomings can be addressed by reducing the angle between the directions that local models move in. Based on this fact, we propose a new **F**ederated enhancement with **C**osine-similarity penalty (FedCos). We introduce an auxiliary global direction that all clients refer to in the local training phase to reduce the directional inconsistency of local models. This constraint accelerates the global updates and diminishes the distances between local models so that the aggregated model is not far away from all the local optima. By analyzing the execution process, we observe FedCos explores more points around the convergence point of FedAvg in the parameter space, which facilitates more points closer to the global optimum to be found.² Meanwhile, FedCos is auto-adapt to the FL settings. The effect of penalty is tuned automatically according to the degrees of heterogeneity of data. Elaborate hyperparameter selection for different scenes is no longer required. In conclusion, the main contributions are summarized below:

- We investigate how data heterogeneity leads to inefficient communication and performance degradation in detail, and explore the angle between the directions that local models move in is the critical issue.
- We propose FedCos, which reduces the directional inconsistency of local models by introducing a cosine-similarity penalty. FedCos can obtain better models than FedAvg and is auto-adapt to the settings without elaborate hyperparameters selection.
- We construct a wide variety of FL scenes comprising

²In fact, due to the nonconvexity of neural network, SGD and other optimization methods aim to find minima. Here we do not distinguish minima and optima, which is not affect the analysis of this paper.

different degrees of data heterogeneity, varying amounts of participants, under cross-silo and cross-device settings. FedCos outperforms other well-known FL methods (FedAvg, FedProx, FedOpt and FedAvgM) regardless of FL scenes, and can enhance them. Furthermore, FedCos improves the efficiency of communication. It greatly reduces the number of communication rounds to obtain global models with the same performance. To our best knowledge, FedCos is the first enhancement that can persistently outperform and enhance FedAvg and other FL methods in a variety of scenes.

II. RELATED WORK

As an extension of distributed training, McMahan et al. [11] first propose the concept of federated learning and related training paradigm FedAvg. After that, a lot of attention has been attracted to explore its potential and applicability. It has been shown that the data heterogeneity of clients induce performance decline [8], [22], [23]. To address this issue, existing methods primarily adjust the local training. Among them, FedProx [12] adds a regularization to the local loss function, which enables the local parameters not far from the global parameters. FedMMD [13] tries to constrain the distribution of the local model close to the global one. SCAFFOLD [20] manually modifies the drift of local training. However, most of them fail to outperform FedAvg in multiple FL scenarios. Another approach improves the performance by modifying the aggregation scheme. For example, FedBE [16] introduces a Bayesian scheme, where the global model is regarded as the expectation of model distribution, but too much training cost is introduced. FedDF [17] introduces data distillation technology to distill the global model from multiple local models, but additional public data is needed. FedPA [18] constructs a posterior model replacing the weighted average model. The prior assumption of uniform distribution is suspicious. DRFA [19] aims to solve a more general problem with arbitrary weights for local clients, which makes the problem harder.

Additionally, some studies enhance the performance under specific context. GKT [24] focuses on the scenario where the edge has limited resources. FedRobust [25] considers the data distribution has a common drift on one client. CFL [26] and IFCA [27] divide the clients into several classes and learn a model for each class, where the data distribution is heterogeneous. These studies for specific problems are unsuitable for general FL scenes. Besides accuracy, FL also faces other challenges [2] such as communication, privacy, security, personality, which also arouse wide attentions [28]–[32] but are out of the scope of this paper. In addition, many traditional studies such as the Ascent Method [33] and ADMM [34] focus on similar problems with theoretical guarantees, which are not generalizable to the deep learning-based FL, where strong duality is no longer satisfied.

III. PROBLEM STATEMENT

A. Background: FedAvg structure

In this paper, we discuss the distributed network as shown in Fig. 1. It involves N edge devices as clients and one parameter

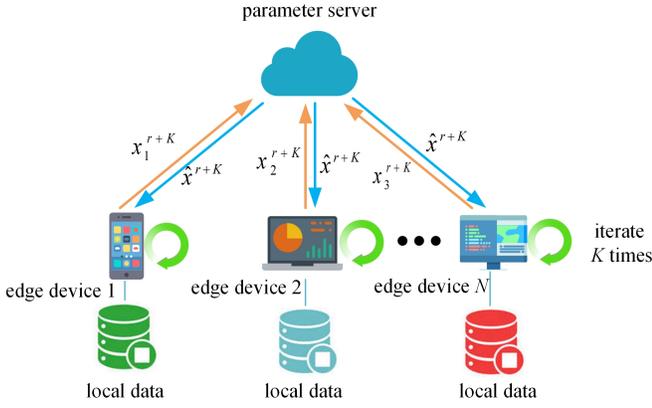


Fig. 1. The federated learning structure for distributed edge data.

server to jointly learn a global model $x \in \mathbb{R}^m$ without data sharing by federated learning (FL). Each device $i \in [1, N]$ owns one local dataset \mathcal{D}_i . The aim can be formulated as the following problem

$$\arg \min_x f(x) = \sum_{i=1}^N \lambda_i f_i(x), \quad (1)$$

where f_i represents the loss function on client i , and λ_i is the weight satisfying $\sum_{i=1}^N \lambda_i = 1$. To facilitate this, Federated Averaging algorithm (FedAvg) is proposed in the federated setting, which tackles the problem iteratively. More specifically, each client performs SGD up to T steps and synchronizes the model every K steps, i.e., synchronization happens at steps $\mathcal{I} = \{nK | n = 1, 2, \dots\}$. There are two phases in each federated learning round:

Local training (on client). For client $i \in \mathcal{S}$, where \mathcal{S} denotes the set of clients participating in the training. The local model iterates from x_i^r ($r \in \mathcal{I}$), which is initialized to the latest received global model \hat{x}^r , and updates as

$$x_i^{r+t+1} = x_i^{r+t} - \eta_i \nabla f_i(x_i^{r+t}, B_i), \quad (2)$$

where $t = \{0, \dots, K-1\}$, B_i is the mini-batch sampled from client i 's local data \mathcal{D}_i , and η_i is the learning rate. After K steps, the client sends the local model x_i^{r+K} to the server.

Global aggregation (on server). The global model is updated as follows:

$$\hat{x}^{r+K} = \sum_{i \in \mathcal{S}} \lambda_i x_i^{r+K}, \quad (3)$$

where $\lambda_i = \frac{|\mathcal{D}_i|}{\sum_{j \in \mathcal{S}} |\mathcal{D}_j|}$, and then the global model parameters are distributed to all or selected clients for the next ‘‘local training and global aggregation’’ round.

In each round, some or all clients participate in the training, which corresponds to distinct FL scenarios called ‘‘cross-device’’ or ‘‘cross-silo’’ respectively.

B. Directional inconsistency

Data heterogeneity of clients, i.e., non-IID data setting, is a common phenomenon in FL, which seriously affects FedAvg's performance. Although numerous methods are proposed, few

approaches exceed FedAvg in all scenes. Li et al [15] exhaustively compare FedAvg and other three typical improved methods (FedProx [12], SCAFFOLD [20] and FedNova [14]), finding out FedAvg still exceeds others in a variety of non-IID scenarios. As the standard solution of FL, FedAvg is still the most appropriate baseline that most studies concern.

In non-IID scenarios, the local model of each client updates iteratively to its local optimum. During the local iterations, local models move in different directions, which causes two problems as follows:

Problem 1 (Communication inefficiency). *Heterogeneity slows down, even stalls the global updates, making the training inefficient, even aborting the training halfway in practice.*

In each round, we define the displacement vector of client i after the local training phase as $d_i = x_i^{r+K} - \hat{x}^r$. Based on equation 3, we obtain the displacement vector of the global model $\hat{d} = \sum_{i=1}^N \lambda_i d_i$ whose norm representing the moving distance of the global model can be denoted as

$$\|\hat{d}\| = \left\| \sum_{i=1}^N \lambda_i d_i \right\| = \sqrt{\sum_{i=1}^N \lambda_i^2 \|d_i\|^2 + \sum_{i \neq j} \lambda_i \lambda_j \|d_i\| \|d_j\| \cos \theta_{ij}}, \quad (4)$$

where θ_{ij} is the angle between two vectors d_i and d_j . When we fix $\|d_i\|$ and λ_i ($1 \leq i \leq N$) (which are dominated by the learning rate and data amount respectively), $\|\hat{d}\|$ is determined by θ_{ij} .

If the data distribution among the parties is homogeneous (i.e., IID), for any x , i and j , $\nabla f_i(x) \approx \nabla f_j(x)$. As a result, all the local models move roughly in the same destination (Fig. 2(a)), namely $\theta_{ij} \approx 0$. If $\|d_i\| = \|d\|$ ($1 \leq i \leq N$), we have $\|\hat{d}\| = \|d\|$. Thus in IID setting, the moving distance of the global model is as large as the local one.

If the data distribution among the clients is heterogeneous (i.e., non-IID), $\nabla f_i(x)$ and $\nabla f_j(x)$ are quite different. In this case, all the local models move in different direction (Fig. 2(b)). Since in high-dimensional space, any two vectors are almost perpendicular, i.e., $\theta_{ij} \approx 90^\circ$. If $\|d_i\| = \|d\|$ and $\lambda_i = \lambda$ ($1 \leq i \leq N$), we have $\|\hat{d}\| = \frac{\|d\|}{\sqrt{N}}$, which means in the non-IID data settings, the global model moves far less distance than local one. When N is large, the global model updates are very slow, even stalled ($N \rightarrow \infty$). To mitigate this problem, from (4), a smaller θ_{ij} leads to a larger $\|\hat{d}\|$, which means the global model moves further. It also be verified in Fig. 3(a) and Fig. 12(c) in the experiment section.

Problem 2 (Poor performance). *The convergence point of FedAvg is far from the true global optimum in non-IID settings.*

Based on Taylor's theorem, if the distance between a point x and the optimum of one client (e.g., x_1^*) is small, x is a suitable model for its local data. Therefore, we aim to find a point close to the optima of all clients. Under the non-IID condition, the distances between the local optima of different clients are large. If x is too close to the optimum of one client, it usually has bad performance for other clients. The global optimum achieves the best balance between different clients,

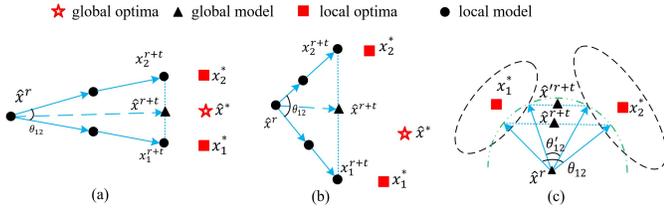


Fig. 2. Examples of FL update with two clients.(a) IID setting; (b) Non-IID setting;(c) adjusting the angle between the directions that local models move in.

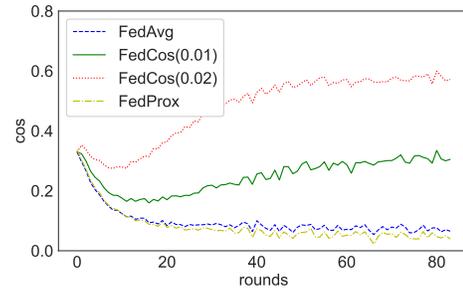
which is not too far away from all the local optima. Actually, Xie et al. [35] prove that SGD favors the flat minima. Izmailov et al. [36] also point out that the optimum is located in a flat region, where the performance is almost the same (referred to as “optimal region” for short). SGD usually stays at a point at the edge of the region. Hence it is feasible to find a representative point on each optimal region, which is closer to the optimal regions of other clients. As shown in Fig. 2(c), if we adjust the angles between the directions that local models move in (e.g., $\theta_{12} \rightarrow \theta'_{12}$), and keep the local models in or at the edges of the optimal regions, the new aggregated model (e.g., \hat{x}^{r+t}) is better than the previous one (e.g., \hat{x}^r) for being closer to the local optima. Suppose there are two clients satisfying $\|d_1\| = \|d_2\| = \|d\|$ and $\lambda_1 = \lambda_2$, and the angle between $\|d_1\|$ and $\|d_2\|$ is θ_{12} . The distance between global model and local model is $\|d\| \cos \frac{\theta_{12}}{2}$. A smaller angle between d_1 and d_2 may lead to a better aggregated model. It also be verified in Fig. 3(a) and Fig. 12(b) in the experiment section.

From the above investigation, reducing the directional inconsistency of local models would be a critical issue to mitigate the influence of data heterogeneity, which motivates us to introduce a cosine-similarity penalty to reduce the angle between the directions that local models move in.

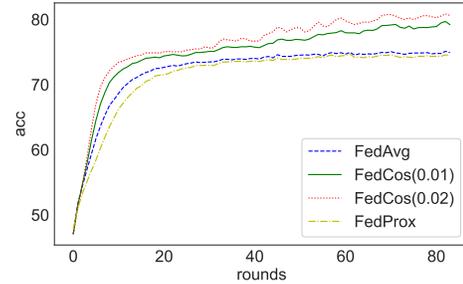
Algorithm 1 FedCos: Federated learning with Cosine-similarity penalty

Input: $K, T, \hat{x}^0, \hat{d}^0 = \vec{0}$

- 1: Server sends \hat{x}^0 and \hat{d}^0 to all clients.
 - 2: **for** $t = 1$ to T **do**
 - 3: (Local training:)
 - 4: Each client in \mathcal{S} updates the local model x_i^t by minimizing (5).
 - 5: **if** $t\%K == 0$ **then**
 - 6: Each client in \mathcal{S} sends x_i^t to the server.
 - 7: (Global aggregation:)
 - 8: The server updates \hat{x}^t and \hat{d}^t by (3) and (7).
 - 9: The server randomly selects a subset of clients \mathcal{S} , and sends \hat{x}^t and \hat{d}^t to the clients in \mathcal{S} .
 - 10: **end if**
 - 11: **end for**
-



(a) The cosine of the angle between two local models in non-IID scenario



(b) The accuracy on test dataset of global models

Fig. 3. Reducing the inconsistency of directions of clients’ models helps improve model performance. The experiment settings are the same as those in Fig. 12.

IV. METHOD

A. FedCos

Based on the above analysis, we try to decrease the angle between the directions that any two local models move to in the local training phase. However, it is impractical to get others’ direction vectors. Therefore, we construct a global direction to constrain the local model updated along this direction. By adding this auxiliary direction, all the local models update towards the same direction to mitigate the direction inconsistency. Meanwhile, we use the auxiliary direction in which the global model moves, the performance degradation of global model is avoided. Based on this design, the loss function of the local model is

$$L_i(x_i; \hat{x}^r, \hat{d}^r) = f(x_i) + \mu_i(1 - \cos \theta_i), \quad (5)$$

where μ_i is the weight of directional consistency, \hat{x}^r, \hat{d}^r are the initial global model and the global direction in this round, and θ_i is the angle between $(x_i - \hat{x}^r)$ and \hat{d}^r , so

$$\cos \theta_i = \frac{\langle x_i - \hat{x}^r, \hat{d}^r \rangle}{\|x_i - \hat{x}^r\| \cdot \|\hat{d}^r\|}. \quad (6)$$

We define $\theta_i = 0$ (i.e., $\cos \theta_i = 1$) if $\|x_i - \hat{x}^r\|$ or $\|\hat{d}^r\|$ is equal to 0. We employ the direction in which the latest global model is moving as the reasonable global direction, i.e.,

$$\hat{d}^r = \hat{x}^r - \hat{x}^{r-K}, \quad (7)$$

where \hat{x}^{r-K} is the global model of the last round. Base on the above formation, we propose our algorithm FedCos (Federated learning with Cosine-similarity penalty) and summarize it in Algorithm 1. FedCos is similar to FedAvg with a few differences. In the local training phase, instead of just minimizing

the loss $f_i(x)$ based on the local data, each client adds a direction penalty. In the global aggregation phase, the server also calculates the new global direction referenced in the next round besides aggregation as FedAvg.

Moreover, we do not limit the form of $f_i(x_i)$ in (5). Therefore, FedCos also can be regarded as an incrementation on other methods of federated learning. For instance, if $f_i(x_i)$ is the cross entropy function, the basic method is FedAvg. If $f_i(x_i)$ is the cross entropy function adding a term $\|x_i - \hat{x}^r\|^2$, the basic method is FedProx.

B. Property analysis

We discuss some advantages of FedCos.

Property 1. FedCos can mitigate the inconsistency of directions of clients' models.

This property is the motivation of FedCos, which is consistent with the practical experiments. In Fig. 3, we illustrate the directional inconsistency by randomly selecting two local models (indexed by i and j). In Fig. 3(a) shows the change of angle of the two local models (i.e., the angle between $x_i^{r+K} - \hat{x}^r$ and $x_j^{r+K} - \hat{x}^r$) in non-IID scenario at each aggregation. For FedAvg, as the training progresses, the direction of local clients become different. By contrast, FedCos reduces the inconsistency with smaller angles. Fig. 3(b) illustrates the performances on test dataset for relevant models. It shows that FedCos with less inconsistency outperforms other two methods. Therefore, reducing the inconsistency of directions of clients' models helps improve model performance.

Property 2. FedCos is auto-adapt to data heterogeneity without selecting penalty weight elaborately.

For FedCos, the update of local model is controlled by $f_i(x)$ and $h(x) = \mu_i(1 - \cos \theta_i)$ (5), where $0 \leq h(x) \leq 2\mu_i$ and usually $f_i(x) \geq 0$. The influence of penalty varies for different data heterogeneity. For IID scenarios, all the local models of clients move in the same direction. The given global direction is consistent with the local optimal ones ($\cos \theta_i \approx 1$). In this situation, FedCos automatically degrades to FedAvg when the penalty term approaches 0. For non-IID scenarios, more serious data heterogeneity leads to more direction inconsistency of local model updates, which enhances the effect of penalty even for a fixed μ_i (i.e., a bigger $h(x)$), and vice versa.

In the local training phase, $f_i(x)$ is dominant and $h(x)$ has little influence in the first few iterations (usually we choose a small μ_i). FedCos has little difference from FedAvg. When $f_i(x)$ approaches 0, i.e., the local model is close to the local minimum, $h(x)$ becomes more effective. This penalty only constrains the orientation of local model updates. Since there is a flattened optimal region [35], [36], direction variation has little influence on accuracy once the model reaches the area, which guarantees $f_i(x)$ and $h(x)$ get small simultaneously. Therefore, FedCos cannot has little influence on the performance of local model. Furthermore, due to reducing the model inconsistencies by the same auxiliary direction, a better aggregated model can be gained.

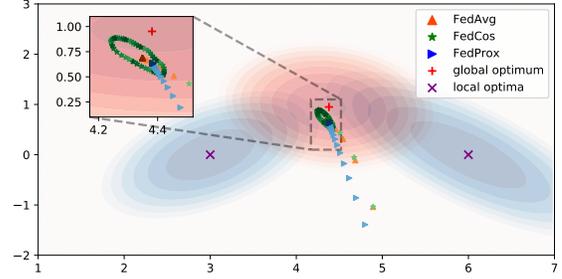


Fig. 4. Comparing the global model trajectories among FedAvg, FedProx and FedCos with 80 rounds. The points with darker colors denote the models from later rounds.

Therefore, FedCos provides a universal performance improvement scheme, which is auto-adapt to the changes of data heterogeneity. We do not need to elaborately select the penalty weight regardless of the non-IID degrees of data. A fixed one would perform well for many kinds of FL settings, which is verified in the following experiments.

In contrast, FedProx reduces the differences between the local models and the aggregated global model by adding the penalty $\|x - \hat{x}\|^2$ with weight μ , which plays an invariant role no matter how heterogeneous the data is. Moreover, for IID scenarios, there is little difference between the local models, but the penalty takes an opposite effect that prevents the model from approaching the optimum. Thus FedProx has to develop a complex heuristic method to tune μ carefully (section 5.3 in [12]).

V. PERFORMANCE ANALYSIS

To illustrate the properties of FedCos, we construct a simple two-dimensional federated example Fig. 4. There are two clients with local convex optimization problems denoted by $f_1(x)$ and $f_2(x)$.³ The global optimization problem is $f(x) = f_1(x) + f_2(x)$. Since $f_1(x)$ and $f_2(x)$ are quite different (which simulates the non-IID scenario), the optimum of $f(x)$ is different with both local ones. With the same parameter settings (such as the learning rate, the number of local steps, and the number of rounds), we run FedCos, FedAvg and FedProx respectively. For FedAvg, it cannot achieve the global optimum. Specifically, the local models (almost) reach the local optima after iterating enough steps at the end of one round. In this case, the global model is close to the middle point of two local optima, rather than the true global optimum. FedProx, which attempts to deal with this problem by decreasing the distances of local models, forces the local models updating near the global model of the last round. Unfortunately, FedProx fails to outperform FedAvg. Furthermore, since it restricts the update of local models, we discover FedProx converges slower than FedAvg from the

³In the example, the explicit functions are $f_1((x_{(0)}, x_{(1)})) = 0.5(x_{(0)} - 6)^2 + 0.75x_{(0)}x_{(1)} + 0.5x_{(1)}^2$, $f_2((x_{(0)}, x_{(1)})) = 0.5(x_{(0)} - 3)^2 - 0.5x_{(0)}x_{(1)} + 0.5x_{(1)}^2$. The initial point of model is $(5.1, -3.1)$.

trajectory of the global model. For FedCos, the global model walks around the convergence point of FedAvg so that it can achieve some points closer to the global optimum. For the scenes with more clients, they can be regarded as the superposition of multiple scenes with two participants, where the same phenomenon exists (more details in Appendix A.1).

Although Fig. 4 only shows a toy example, it helps us to comprehend the advantages of FedCos as follows: *FedCos can achieve the point closer to the global optimum than FedAvg.*

At first we analyze that in each round the aggregated model of FedCos moves further in the given global direction than FedAvg's. Suppose $g(x) = 1 - \frac{\langle x - \hat{x}, \hat{d} \rangle}{\|x - \hat{x}\| \|\hat{d}\|}$. The local optimization for FedCos is $f_i(x) + \mu_i g(x)$, while it for FedAvg is $f_i(x)$. Since rotating the coordinate system does not affect the distance and the relationship between the points or vectors, for convenience we use a matrix $X^{-1} \in \mathbb{R}^{d \times d}$ rotating the coordinate system to a new one so that $\hat{d}X = [\|\hat{d}\|, 0, \dots, 0]^T$, i.e., the global direction is along the first dimension. For simplicity, all the analysis is in the new coordinate system, and we do not times X explicitly, i.e., x means xX . For any component $x_{(i)}$ of x , we have

$$\begin{aligned} \frac{\partial g(x)}{\partial x_{(i)}} &= -\frac{1}{\|\hat{d}\|} \cdot \frac{d_{(i)} \|x - \hat{x}\|^2 - \langle x - \hat{x}, \hat{d} \rangle (x_{(i)} - \hat{x}_{(i)})}{\|x - \hat{x}\|^3} \\ &= -\frac{1}{\|\hat{d}\|} \cdot \frac{d_{(i)} \|x - \hat{x}\|^2 - \|\hat{d}\| (x_{(0)} - \hat{x}_{(0)}) (x_{(i)} - \hat{x}_{(i)})}{\|x - \hat{x}\|^3} \\ &= \begin{cases} \frac{(x_{(0)} - \hat{x}_{(0)})^2 - \|x - \hat{x}\|^2}{\|x - \hat{x}\|^3} \leq 0, & i = 0, \\ \frac{(x_{(0)} - \hat{x}_{(0)}) (x_{(i)} - \hat{x}_{(i)})}{\|x - \hat{x}\|^3}, & i \neq 0. \end{cases} \end{aligned} \quad (8)$$

From (8) the partial derivative is nonnegative in the direction of the auxiliary global direction.

In one round that \hat{d}^r is fixed, if client i starts local training of FedAvg and FedCos with the same model, the local model is updated by (2) for FedAvg and

$$x_i^{r+t+1} = x_i^{r+t} - \eta_i \nabla f_i(x_i^{r+t}, B_i) - \eta_i \mu_i \nabla g(x_i^{r+t}) \quad (9)$$

for FedCos, where $\nabla g(x_i^{r+t}) = [\frac{\partial g(x)}{\partial x_{(0)}}, \dots, \frac{\partial g(x)}{\partial x_{(d-1)}}]^T$. Compared to (2), equation (9) adds $-\eta_i \mu_i \nabla g(x_i^{r+t})$ at each step, which has a nonnegative component along the direction of \hat{d}^r (8). Intuitively speaking, all the local models of FedCos move further than FedAvg's in the direction of the global direction. The parameter μ_i controls the distance of the two points. Thus, the aggregated global model also move further than FedAvg's in this direction.

Then we investigate the relationship of global models between FedAvg and FedCos. When FedAvg converges, i.e., the local models and the aggregated model are fixed. As shown in Fig. 5(a), The local model of FedCos is close to the corresponding one of FedAvg, but has an offset along the global direction. The aggregated model is above the FedAvg's (named reference model). Hence in the next round, the global direction is still pointing up. Due to the constraint of $\nabla f_i(x_i^{r+t}, B_i)$, the local model of FedCos cannot move up endlessly. At some round, the local model of FedCos does not move up, and the aggregated global model does not move up too (Fig. 5(b)).

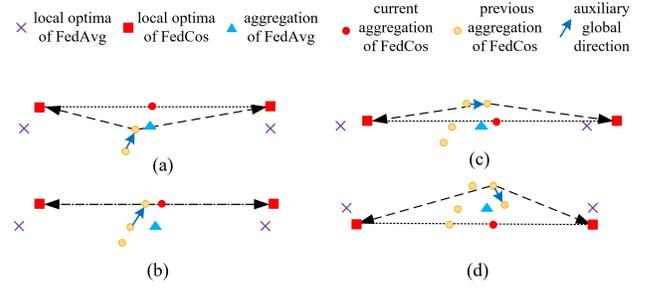


Fig. 5. The description of running FedCos.

In the next round, the global direction is no longer pointing up (e.g., it points horizontally to the right in Fig. 5(c)). In this case, the local model does not achieve the position as high as the previous one. The aggregated model is below the previous one (Fig. 5(c)). Thus in the next round, the global direction is changed to point down, and the aggregated model goes on too (Fig. 5(d)). In this way, the aggregated model of FedCos walks around the reference model round by round, showing the phenomenon in Fig. 4. The "radius" is controlled by μ_i , since it is related to the distance the aggregated model moving.

From the above analysis, the global models of FedCos are around the convergence point of FedAvg. Therefore, with a proper μ_i , which makes the "radius" smaller than the distance between the convergence point of FedAvg and the true global optimum, there are some points of FedCos closer to the true global optimum than the convergence point of FedAvg. In high-dimensional space, FedCos cannot iterate through all the points on the hypersphere determined by μ_i , but it can explore sufficient points around the convergence point of FedAvg to explore points closer to the true global optimum with large probability.

VI. EXPERIMENTS

A. Experimental settings

1) *Baselines, Datasets and Models:* We choose FedAvg [11], FedProx [12], FedAvgM [37], [38] and FedOpt [39] as baselines. Since we need hundreds of edge devices in our experiments, we simulate them on a server as most of previous literatures for federated learning. All methods are implemented with PyTorch1.8 and trained on GeForce RTX 3090 GPUs. We conduct experiments on 4 datasets and corresponding models as:

- **CIFAR10** [40]. CIFAR10 consists of 60000 32x32 color images in 10 categories, with 6000 images per category. The training set contains 50000 training images and the test set contains 10000 images. For CIFAR10, we use a CNN network with 2 convolutional-pooling layers plus 2 fully connected layers. In Section VI-D1 Lenet is used, which also shows similar results.
- **FMNIST** [41]. Fashion-MNIST is an updated version of MNIST. It has 28x28 grayscale images of 70000 fashion products from 10 categories. Each category has 7000 images. The training set contains 60000 images and the

TABLE I
PERFORMANCE (TOP-1 VALIDATION ACCURACY) COMPARISON ON CIFAR10 UNDER DIFFERENT DATA SETTINGS, WHERE $N = 5$.

	Totally non-IID	90% non-IID	70% non-IID	IID
FedAvg	49.56	64.30	66.82	68.78
FedProx(0.1)	44.75	61.75	65.30	66.27
FedAvgM(0.5)	48.35	64.06	66.38	67.95
FedOpt(1.5)	51.21	66.11	68.94	70.57
FedCos(0.01)	54.54 (4.98 \uparrow)	68.33 (4.03 \uparrow)	70.50 (3.68 \uparrow)	71.75 (2.97 \uparrow)
FedCos(0.02)	54.80 (5.24 \uparrow)	69.83 (5.53 \uparrow)	71.46 (4.64 \uparrow)	73.30 (4.52 \uparrow)
FedProx(0.1)+FedCos(0.02)	50.85 (6.10 \uparrow)	67.68 (5.93 \uparrow)	70.65 (5.35 \uparrow)	71.43 (5.16 \uparrow)
FedAvgM(0.5)+FedCos(0.02)	54.30 (5.95 \uparrow)	69.27 (5.21 \uparrow)	71.01 (4.63 \uparrow)	72.39 (4.44 \uparrow)
FedOpt(1.5)+FedCos(0.02)	54.46 (3.25 \uparrow)	69.97 (3.86 \uparrow)	72.64 (3.70 \uparrow)	73.37 (2.80 \uparrow)

TABLE II
PERFORMANCE (TOP-1 VALIDATION ACCURACY) COMPARISON ON FMNIST UNDER DIFFERENT DATA SETTINGS, WHERE $N = 7$.

	Totally non-IID	90% non-IID	70% non-IID	IID
FedAvg	74.91	85.15	86.83	88.04
FedProx(0.1)	74.47	84.35	86.12	87.51
FedAvgM(0.5)	74.76	85.18	86.65	87.98
FedOpt(1.5)	75.47	85.97	87.60	88.43
FedCos(0.01)	79.21 (4.30 \uparrow)	87.37 (2.22 \uparrow)	88.69 (1.86 \uparrow)	89.46 (1.42 \uparrow)
FedCos(0.02)	80.63 (5.72 \uparrow)	87.61 (2.46 \uparrow)	89.08 (2.25 \uparrow)	89.52 (1.48 \uparrow)
FedProx(0.1)+FedCos(0.02)	78.23 (3.76 \uparrow)	87.15 (2.80 \uparrow)	88.56 (2.44 \uparrow)	89.22 (1.71 \uparrow)
FedAvgM(0.5)+FedCos(0.02)	78.72 (3.96 \uparrow)	87.49 (2.31 \uparrow)	88.67 (2.02 \uparrow)	89.31 (1.33 \uparrow)
FedOpt(1.5)+FedCos(0.02)	80.76 (5.29 \uparrow)	87.51 (1.54 \uparrow)	89.03 (1.43 \uparrow)	89.45 (1.02 \uparrow)

test set contains 10000 images. Compared to MNIST, the complexity of the task is increased. For FMNIST, we use a multi-layered perceptron network with 2 fully connected layers.

- **CIFAR100** [40]. CIFAR100 consists of 100 classes, each of which has 500 images for training and 100 images for testing. For CIFAR100, we use ResNet18 with group normalization.
- **Shakespeare** [11] is built from The Complete Works of William Shakespeare. It consists of 143 characters with 517,106 samples. Each client takes all the samples of one character as local data. We use a two-layer LSTM classifier containing 100 hidden units with an 8D embedding layer. The task is a next-character prediction with a sequence of 80 characters as input, and there are 80 classes of characters in total.

2) *Hyperparameter setting*: For all experiments, SGD is used as the optimizer in the local training phase, and the learning rate is 0.01 as default.

For Table I, 5 clients use the same hyperparameters. The batch size is 128. The number of local training steps in each round is 400. The result is after 500 epochs for each client (about 100 rounds of aggregation). For Table II, 7 clients are with the same settings as above. For table III, 5 clients are involved. The batch size is 64. The number of local training steps in each round is 200. The result is after 300 epochs for each client.

For Fig. 6, all the data are sorted by labels and divided into $2N$ equal subsets. Each client is randomly allocated 2 subsets as its local dataset. We set the batch size 64 and the number of local training steps in each round is 200. It shows the result after 500 epochs for each client.

TABLE III
PERFORMANCE (TOP-1 VALIDATION ACCURACY) COMPARISON ON CIFAR100 UNDER DIFFERENT DATA SETTINGS WITH $N = 5$.

	Totally non-IID	90% non-IID	IID
FedAvg	50.51	52.84	64.93
FedProx(0.01)	50.34	50.54	64.58
FedCos(0.01)	52.62	54.49	65.13
FedCos(0.02)	53.37	55.50	65.35

For Table IV the data distribution method is the same as Section VI-B2. We set the batch size 64 and the number of local training steps in each round is 100. All the methods run 250 rounds. For Table V, each client takes all the samples of one character as local data. The batch size is 10. The learning rate is 0.5. The number of local training steps in each round is 50.

B. Performance comparison

To compare the performance more comprehensively, we investigate the performance of the candidate methods for both FL scenarios named “cross-silo” and “cross-device” [42]. In cross-silo scenario, all the clients take part in the local training in each round, while in cross-device scenario only a part of clients is active in each round. We construct experimental scenes with various non-IID data settings and different number of participants to verify the effectiveness of FedCos exhaustively.

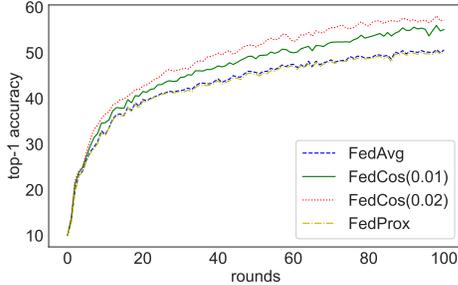
To compare the performance of FedCos, FedAvg and FedProx in different degrees of data heterogeneity, we construct 4 diverse scenes. (1) Totally non-IID setting: M classes of data are sorted by labels and divided into N equal subsets. Each client contains one as its local dataset. (2) 90% non-IID

TABLE IV
PERFORMANCE COMPARISON OF CNN (THE BEST PERFORMANCE OF THE GLOBAL MODEL IN 250 ROUNDS AND THE PERFORMANCE AT THE END OF THE TRAINING).

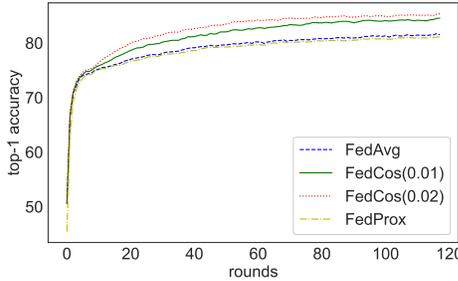
	CIFAR10				FMNIST			
	10%		20%		10%		20%	
	Best	Last	Best	Last	Best	Last	Best	Last
FedAvg	49.91	48.69	51.44	50.61	82.69	79.68	83.29	79.96
FdeProx(0.1)	49.38	49.63	50.66	50.38	82.44	79.73	83.26	79.68
FedCos(0.02)	54.90	52.54	57.97	56.15	84.49	82.21	86.07	84.35
FedCos(0.05)	59.36	57.50	61.80	59.07	85.69	85.22	86.74	85.87

TABLE V
PERFORMANCE COMPARISON OF TWO-LAYER LSTM AFTER 100 ROUNDS IN CROSS-DEVICE FL CIRCUMSTANCES.

	Shakespeare	
	10%	20%
FedAvg	37.46	38.08
FdeProx(0.01)	35.60	35.64
FedCos(0.005)	40.25	40.91
FedCos(0.01)	41.80	42.96



(a) Performance on CIFAR10



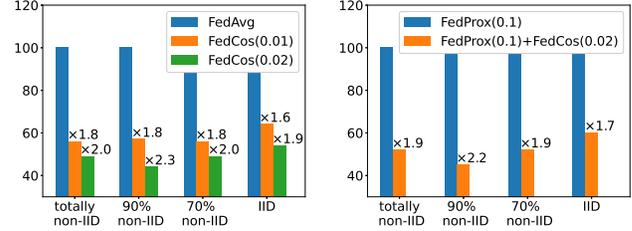
(b) Performance on FMNIST

Fig. 6. Performance comparison under non-IID settings where $N = 20$.

setting: Based on totally non-IID setting, 10% of data of each client are extracted, shuffled and then divided into N parts. Each client has one. Thus each client has 10% of homogeneous data. (3) 70% non-IID setting: Like 90% non-IID setting, each client has 30% of homogeneous data. (4) IID setting: The data are randomly shuffled and divided into N parts. Each client contains one.

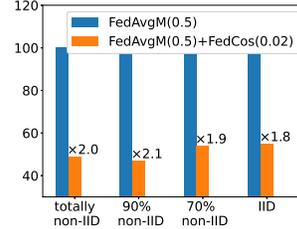
1) *Cross-silo FL with a few number of clients*: Table I shows the performance comparison of FedAvg, FedProx, and FedCos with CNN on CIFAR10, where 5 clients are involved. The penalty weight for FedProx is 0.1.⁴ For FedCos, we investigate two weights, i.e., $\mu_i = 0.01$ and 0.02. For different degrees of data heterogeneity, FedCos leads to a gain of at least 3% over FedAvg and other 3 baselines in all cases. Then we enhance each baseline method with FedCos. It shows that FedCos can significantly improve the performance under various data heterogeneity. Interestingly, in the IID setting, FedCos still improves the baselines. This is because the notion

⁴In these experiments, the performance of FedProx with too small penalty weights (e.g., 0.01) would degenerate to FedAvg.

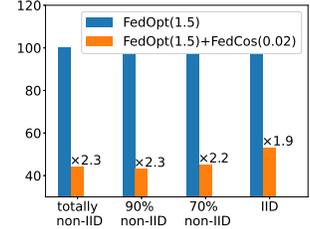


(a) Reduce communication round for FedAvg

(b) Reduce communication round for FedProx



(c) Reduce communication round for FedAvgM



(d) Reduce communication round for FedOpt

Fig. 7. FedCos enhances the communication efficiency of baselines on CIFAR10.

of “IID” is based on statistics, while heterogeneity is pervasive among clients in practice. Table II illustrates the performance of MLP on FMNIST. We harvest the similar conclusions.

We also investigate more complex models and datasets. Table III⁵ illustrates the performance of ResNet18 (with group normalization) on CIFAR100 with 5 clients. FedCos still significantly outperforms others.

2) *Cross-silo FL with multiple clients*: When more clients join in, FedCos still leads the way. We increase the number of clients to 20, and the results are illustrated in Fig. 6. On CIFAR10, after 500 epochs (about 100 rounds), the accuracy of FedCos with $\mu_i = 0.02(0.01)$ reaches 57.14 (55.12), while the accuracy of FedAvg is only 50.63. On FMNIST, FedCos significantly improves the accuracy compared with FedAvg from 81.52 to our best result of 85.36 (84.55) with $\mu_i = 0.02(0.01)$. FedProx has a similar performance to FedAvg but worse than FedCos.

3) *Cross-device FL*: Table IV illustrates specific results, where 100 edge devices join in the FL training, and 10%/20% of clients are randomly chosen in each round. Since different

⁵Since the performance of FedProx with penalty weight 0.1 is too poor, we use 0.01 instead.

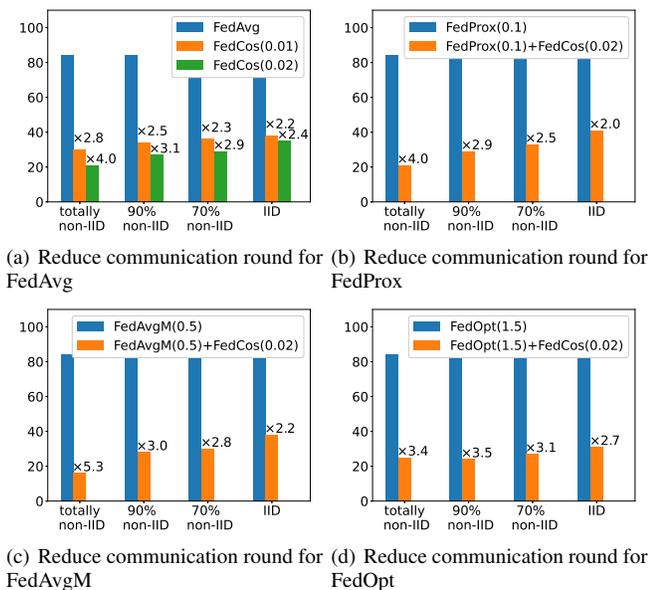


Fig. 8. FedCos enhances the communication efficiency of baselines on FMNIST.

clients contribute the updates in each round, the global model is not as stable as it is in the cross-silo scenario. As a result, we employ a larger penalty to impose the global constraints (e.g., we set $\mu_i = 0.05$). Here, we list the best results in 250 rounds and the results at the end of the training. With different proportions of participants, FedCos is superior to baselines in both measures. Table V illustrates the performance of LSTM on Shakespeare dataset with 143 clients. FedCos is also better than FedAvg and FedProx.

C. Comparison of communication efficiency

FedCos can improve the communication efficiency of each rounds. With the same settings as that in Table I, we firstly perform the 4 baselines for 100 rounds. Then we investigate how many communication rounds are required for the enhanced approaches to train the same performance models. Fig. 7 shows the results on CIFAR10, where FedCos can reduce communication rounds by 50%. It means by the enhancement of FedCos, the communication efficiency of baselines is improved. The same results are in Fig. 8, where the settings on FMNIST are the same as that in Table II. The communication efficiency is improved by 2 to 5 times. Meanwhile, it seems that the improvement is more obvious for larger degrees of heterogeneity of data. It fits our insight that data heterogeneity leads to more serious directional inconsistency of local models, which slows down the global updates.

D. Additional experiments

1) *Detailed Comparison with FedOpt*: From Table I and Table II, FedOpt performs best in 4 baselines, and it can also increase the moving distance of the global model with carefully selected hyperparameters. FedOpt adjusts the learning rate of global updating η_g , where FedAvg is the case with $\eta_g = 1$. Fig. 9 shows the performance comparison

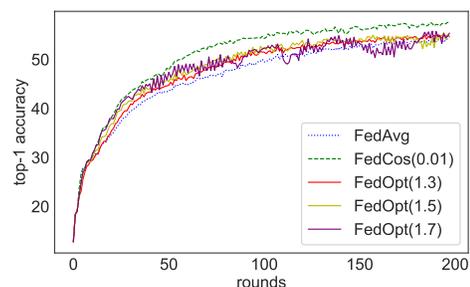


Fig. 9. Comparison of FedCos and FedOpt.

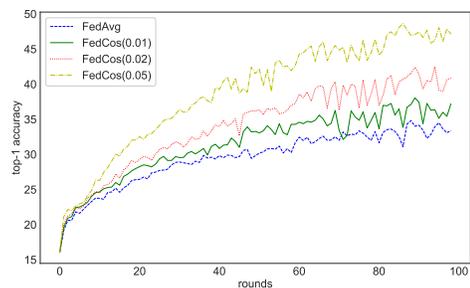


Fig. 10. Performance comparison for ResNet18.

between FedCos and FedOpt with various η_g . For the first several rounds, FedOpt speeds up the training like FedCos. However, no matter what the learning rate is, FedOpt performs as same as FedAvg after 200 rounds. In contrast, FedCos keeps ahead of FedAvg and FedOpt all the time. Furthermore, for FedOpt a larger η_g would make the training instability. In our experiment, when $\eta_g = 1.7$, it has apparent jitter on the performance curve. As η_g continues to increase, the performance would seriously decrease.

2) *Performance comparison under deeper neural networks*: To demonstrate the applicability of FedCos, the performance for more deeper neural networks is investigated. Fig.10 shows the results under total non-IID setting on CIFAR10 for ResNet18, which is typical deep neural network much more complex than previously evaluated models. Same as before, FedCos with various of penalty weights outperforms FedAvg much more. It should be noted that, since batch normalization does not suit the Non-IID setting, we replace it with alter-

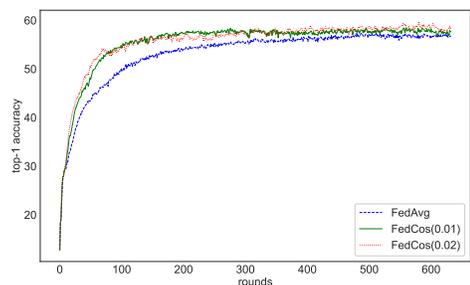


Fig. 11. Performance comparison with 3200 epochs.

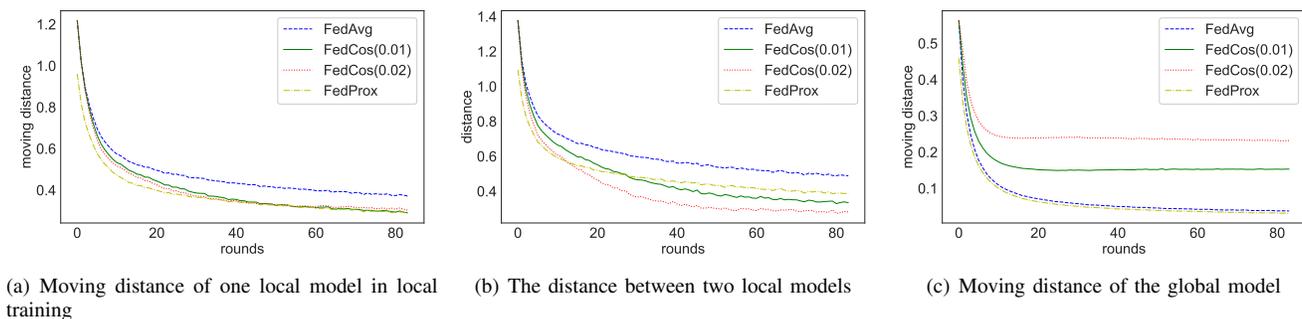


Fig. 12. Mechanism illustration of FedCos.

native group normalization mechanism [43]. Due to the high computational overhead, most of previous works explore the mechanism of FL by simple networks.

3) *Performance comparison under long time training:* In order to investigate the influence of round number, we greatly increase the training time. Fig. 11 shows the comparison under totally non-IID setting on CIFAR10, where each client iterates 3200 epochs (more than 600 rounds). At the end of the training, FedCos with $\mu_i = 0.02(0.01)$ achieves 58.58(57.66), while FedAvg is only 56.83. For the best result on the test dataset, FedCos with $\mu_i = 0.02(0.01)$ achieves 59.73(58.48), while FedAvg is only 57.45. Although the gap between FedCos and FedAvg is narrowed, FedCos still outperforms FedAvg at the end of training. In practice, it is no need to spend so much computational overhead and communication cost, which can be saved by fewer local iteration steps and aggregation round. So it does not affect the results of the previous comparison.

E. Mechanism visualization

As noted in the analysis of previous section, the directional inconsistency of local models is the crucial factor degrading the performance of FL methods in non-IID scenarios, and FedCos eases the inconsistency by the given global direction vector. Fig. 3 and Fig. 12 further verify the results under the totally non-IID scenario in Table II. Specifically, Fig. 12(a) shows the moving distance of one local model. For FedCos, the moving distance of the local model is smaller than it for FedAvg, and the angle between the moving directions of two local models is smaller (Fig. 3(a)). Therefore, the distance of any two local models for FedCos is smaller than FedAvg's (Fig. 12(b)), which validates the analysis in Fig. 2(c). FedCos obtains the local optima (approximately) closer to each other than FedAvg. Thus, the aggregated model is closer to all the local optima and performs better performance.

The global model only has little changes for FedAvg in each round, while it moves farther for FedCos, as shown in Fig. 12(c). Thus the global model updates of FedCos are more efficient, which confirms the view in Fig. 2(b). Moreover, although the moving distance of the local model for FedProx is small (Fig. 12(a)), the angle is large and comparable with FedAvg's (Fig. 3(a)). Accordingly, the distance between the two local models is not as small as expected (Fig. 12(b)).

Meanwhile, because of the strong constraint, the local model for FedProx degrades the efficiency of local training. FedProx performs worse than FedCos/FedAvg, despite with smaller local models' distances than FedAvg's.

VII. CONCLUSION

In this work, we propose FedCos to improve model accuracy in FL. FedCos introduces an auxiliary global direction to guide the direction of local model in training. This scheme improves the performance of the aggregated model by reducing the directional inconsistency of local models. We analyze the properties of FedCos, and point out that FedCos can obtain better models than the standard FL method FedAvg. From the experiments, FedCos vastly outperforms FedAvg in a variety of FL scenes. Due to the complexity of analysis, in this paper more rigorous theoretical proofs are not given, which would be explored in future work.

REFERENCES

- [1] P. Voigt and A. Von dem Bussche, "The eu general data protection regulation (gdpr)," *A Practical Guide, 1st Ed., Cham: Springer International Publishing*, vol. 10, p. 3152676, 2017.
- [2] O. A. Wahab, A. Mourad, H. Otrok, and T. Taleb, "Federated machine learning: Survey, multi-level classification, desirable criteria and future directions in communication and networking systems," *IEEE Communications Surveys & Tutorials*, 2021.
- [3] A. Rakhlin, O. Shamir, and K. Sridharan, "Making gradient descent optimal for strongly convex stochastic optimization," in *Proceedings of the 29th International Conference on International Conference on Machine Learning*, 2012, pp. 1571–1578.
- [4] J. Dean, G. S. Corrado, R. Monga, K. Chen, M. Devin, Q. V. Le, M. Z. Mao, M. Ranzato, A. Senior, P. Tucker *et al.*, "Large scale distributed deep networks," in *Proceedings of the 25th International Conference on Neural Information Processing Systems-Volume 1*, 2012, pp. 1223–1231.
- [5] T. Lin, S. U. Stich, K. K. Patel, and M. Jaggi, "Don't use large mini-batches, use local sgd," in *International Conference on Learning Representations*, 2019.
- [6] H. Yu, S. Yang, and S. Zhu, "Parallel restarted sgd with faster convergence and less communication: Demystifying why model averaging works for deep learning," in *Proceedings of the AAAI Conference on Artificial Intelligence*, 2019, pp. 5693–5700.
- [7] S. Shen, L. Xu, J. Liu, X. Liang, and Y. Cheng, "Faster distributed deep net training: Computation and communication decoupled stochastic gradient descent," in *Proceedings of the Twenty-Eighth International Joint Conference on Artificial Intelligence*, 7 2019, pp. 4582–4589.
- [8] Y. Zhao, M. Li, L. Lai, N. Suda, D. Civin, and V. Chandra, "Federated learning with non-iid data," *arXiv preprint arXiv:1806.00582*, 2018.
- [9] P. Oza and V. M. Patel, "Federated learning-based active authentication on mobile devices," *arXiv preprint arXiv:2104.07158*, 2021.

- [10] Z. Chai, H. Fayyaz, Z. Fayyaz, A. Anwar, Y. Zhou, N. Baracaldo, H. Ludwig, and Y. Cheng, "Towards taming the resource and data heterogeneity in federated learning," in *2019 {USENIX} Conference on Operational Machine Learning (OpML 19)*, 2019, pp. 19–21.
- [11] B. McMahan, E. Moore, D. Ramage, S. Hampson, and B. A. y Arcas, "Communication-efficient learning of deep networks from decentralized data," in *Artificial Intelligence and Statistics*. PMLR, 2017, pp. 1273–1282.
- [12] T. Li, A. K. Sahu, M. Zaheer, M. Sanjabi, A. Talwalkar, and V. Smith, "Federated optimization in heterogeneous networks," *arXiv preprint arXiv:1812.06127*, 2018.
- [13] X. Yao, T. Huang, C. Wu, R.-X. Zhang, and L. Sun, "Federated learning with additional mechanisms on clients to reduce communication costs," *arXiv preprint arXiv:1908.05891*, 2019.
- [14] J. Wang, Q. Liu, H. Liang, G. Joshi, and H. V. Poor, "Tackling the objective inconsistency problem in heterogeneous federated optimization," *Advances in Neural Information Processing Systems*, vol. 33, 2020.
- [15] Q. Li, Y. Diao, Q. Chen, and B. He, "Federated learning on non-iid data silos: An experimental study," *arXiv preprint arXiv:2102.02079*, 2021.
- [16] H.-Y. Chen and W.-L. Chao, "Fed{be}: Making bayesian model ensemble applicable to federated learning," in *International Conference on Learning Representations*, 2021.
- [17] T. Lin, L. Kong, S. U. Stich, and M. Jaggi, "Ensemble distillation for robust model fusion in federated learning," in *Advances in Neural Information Processing Systems*, H. Larochelle, M. Ranzato, R. Hadsell, M. F. Balcan, and H. Lin, Eds. Curran Associates, Inc., 2020, pp. 2351–2363.
- [18] M. Al-Shedivat, J. Gillenwater, E. Xing, and A. Rostamizadeh, "Federated learning via posterior averaging: A new perspective and practical algorithms," in *International Conference on Learning Representations*, 2021.
- [19] Y. Deng, M. M. Kamani, and M. Mahdavi, "Distributionally robust federated averaging," in *Advances in Neural Information Processing Systems*, H. Larochelle, M. Ranzato, R. Hadsell, M. F. Balcan, and H. Lin, Eds. Curran Associates, Inc., 2020, pp. 15 111–15 122.
- [20] S. P. Karimireddy, S. Kale, M. Mohri, S. Reddi, S. Stich, and A. T. Suresh, "Scaffold: Stochastic controlled averaging for federated learning," in *International Conference on Machine Learning*. PMLR, 2020, pp. 5132–5143.
- [21] D. A. E. Acar, Y. Zhao, R. Matas, M. Mattina, P. Whatmough, and V. Saligrama, "Federated learning based on dynamic regularization," in *International Conference on Learning Representations*, 2021.
- [22] F. Haddadpour and M. Mahdavi, "On the convergence of local descent methods in federated learning," *arXiv preprint arXiv:1910.14425*, 2019.
- [23] X. Li, K. Huang, W. Yang, S. Wang, and Z. Zhang, "On the convergence of fedavg on non-iid data," in *International Conference on Learning Representations*, 2019.
- [24] C. He, M. Annavaram, and S. Avestimehr, "Group knowledge transfer: Federated learning of large cnns at the edge," in *Advances in Neural Information Processing Systems*, H. Larochelle, M. Ranzato, R. Hadsell, M. F. Balcan, and H. Lin, Eds. Curran Associates, Inc., 2020, pp. 14 068–14 080.
- [25] A. Reiszadeh, F. Farnia, R. Pedarsani, and A. Jadbabaie, "Robust federated learning: The case of affine distribution shifts," in *Advances in Neural Information Processing Systems*, H. Larochelle, M. Ranzato, R. Hadsell, M. F. Balcan, and H. Lin, Eds. Curran Associates, Inc., 2020, pp. 21 554–21 565.
- [26] F. Sattler, K.-R. Müller, and W. Samek, "Clustered federated learning: Model-agnostic distributed multitask optimization under privacy constraints," *IEEE Transactions on Neural Networks and Learning Systems*, pp. 1–13, 2020.
- [27] A. Ghosh, J. Chung, D. Yin, and K. Ramchandran, "An efficient framework for clustered federated learning," in *Advances in Neural Information Processing Systems*, H. Larochelle, M. Ranzato, R. Hadsell, M. F. Balcan, and H. Lin, Eds. Curran Associates, Inc., 2020, pp. 19 586–19 597.
- [28] P. P. Liang, T. Liu, L. Ziyin, N. B. Allen, R. P. Auerbach, D. Brent, R. Salakhutdinov, and L.-P. Morency, "Think locally, act globally: Federated learning with local and global representations," *arXiv preprint arXiv:2001.01523*, 2020.
- [29] T. Li, S. Hu, A. Beirami, and V. Smith, "Ditto: Fair and robust federated learning through personalization," 2021.
- [30] P. Yu, L. Wynter, and S. H. Lim, "Fed+: A family of fusion algorithms for federated learning," *arXiv preprint arXiv:2009.06303*, 2020.
- [31] E. Diao, J. Ding, and V. Tarokh, "Hetero{fl}: Computation and communication efficient federated learning for heterogeneous clients," in *International Conference on Learning Representations*, 2021.
- [32] M. S. Ozdayi, M. Kantarcioglu, and Y. R. Gel, "Defending against backdoors in federated learning with robust learning rate," *arXiv preprint arXiv:2007.03767*, 2020.
- [33] S. Shalev-Shwartz and T. Zhang, "Stochastic dual coordinate ascent methods for regularized loss minimization," *Journal of Machine Learning Research*, vol. 14, no. 2, 2013.
- [34] S. Boyd, N. Parikh, and E. Chu, *Distributed optimization and statistical learning via the alternating direction method of multipliers*. Now Publishers Inc, 2011.
- [35] Z. Xie, I. Sato, and M. Sugiyama, "A diffusion theory for deep learning dynamics: Stochastic gradient descent exponentially favors flat minima," in *International Conference on Learning Representations*, 2021.
- [36] P. Izmailov, D. Podoprikin, T. Garipov, D. Vetrov, and A. G. Wilson, "Averaging weights leads to wider optima and better generalization," *arXiv preprint arXiv:1803.05407*, 2018.
- [37] T.-M. H. Hsu, H. Qi, and M. Brown, "Measuring the effects of non-identical data distribution for federated visual classification," *arXiv preprint arXiv:1909.06335*, 2019.
- [38] Z. Huo, Q. Yang, B. Gu, L. C. Huang *et al.*, "Faster on-device training using new federated momentum algorithm," *arXiv preprint arXiv:2002.02090*, 2020.
- [39] S. J. Reddi, Z. Charles, M. Zaheer, Z. Garrett, K. Rush, J. Konečný, S. Kumar, and H. B. McMahan, "Adaptive federated optimization," in *International Conference on Learning Representations*, 2021.
- [40] A. Krizhevsky, G. Hinton *et al.*, "Learning multiple layers of features from tiny images," 2009.
- [41] H. Xiao, K. Rasul, and R. Vollgraf, "Fashion-mnist: a novel image dataset for benchmarking machine learning algorithms," *arXiv preprint arXiv:1708.07747*, 2017.
- [42] P. Kairouz, H. B. McMahan, B. Avent, A. Bellet, M. Bennis, A. N. Bhagoji, K. Bonawitz, Z. Charles, G. Cormode, R. Cummings *et al.*, "Advances and open problems in federated learning," *arXiv preprint arXiv:1912.04977*, 2019.
- [43] K. Hsieh, A. Phanishayee, O. Mutlu, and P. Gibbons, "The non-iid data quagmire of decentralized machine learning," in *International Conference on Machine Learning*. PMLR, 2020, pp. 4387–4398.

VIII. APPENDIX

A. FL example with multiple clients

For the simple federated example with more participants, the same phenomenon exists. Fig. 13 illustrates 3 clients. The initial point is (4.53, 0.38). As in the example in Fig. 4, the trajectory of the global model in FedCos is around the stationary point of FedAvg. FedCos can obtain better models closer to the global optimum than FedAvg.

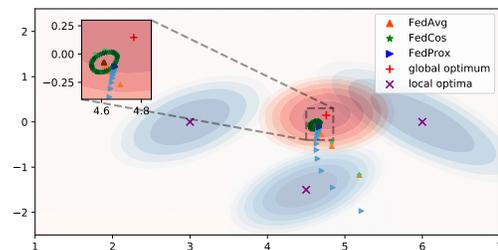


Fig. 13. Comparison of the trajectories of global models by FedAvg, FedProx and FedCos on a simple example with 3 participants. All the methods perform 80 rounds. The points with darker colors denote the models from later rounds.