

DOI:10.18686/ahe.v7i29.10721

Guiding Federated Learning with Low Rank Metrics

Toby Yang Ahmed Elkoushy

University of Toronto Toronto Ontario Canada M5S 1A1

Abstract: Federated learning (FL) is a distributed machine learning method that protects the privacy of local data, and thus is used in scenarios like medical applications. However, FL algorithms are known to demonstrate deteriorated performance when the local data of each client is not independently and identically distributed (Non-IID). In this work, we explore the possibility of guiding the aggregation process inFL using low-rank based layer-wise metrics, and propose an aggregation framework that allows the use of any layer-wise metric to guide the aggregation process. Specifically, a layer-wise metric is used as the weight for the weighted averaging process in the aggregation, and the metric is optionally smoothed by exponential averaging.

In this work, we propose two low-rank layer-wise metrics. However, experiments do not show significant improvement of the proposed methods over the baseline method (FedAvg), thus this paper mainly serves as a record of our exploration and the experimental results.

Keywords: Federated learning (FL); Low rank; Non-IID data; Explainable metrics

1. Introduction

Federated Learning (FL) trains distributed machine learning models, preserving data privacy in cross-silo settings like healthcare. Non-IID data often leads to client drift, harming knowledge preservation ^[4, 6]. We focus on guiding aggregation with explainable low rank factorized metrics in FL, enabling layer-wise metric usage.

2. Related Works

2.1 Federated Learning Problem Formulation

Traditionally, most machine learning optimization problems (involving one client) are of the following form:

$$\min_{w \in \mathbb{R}^d} \quad f(w) \quad \text{where} \quad f(w) \equiv \frac{1}{n} \sum_{i=1}^n f_i(w) \tag{1}$$

Where f(w) is defined to be the empirical loss and is simply the average of the loss on each data point l(xi, yi) or l(yi, yi) where y'i refers to the output of the model.

In the federated case, we would like to construct a problem where we could leverage already well defined optimization methods such as SGD. [5] showed that one can form local subproblems of similar structure on each client and define the following local loss:

$$F_k(w) \equiv \frac{1}{n_k} \sum_{i \in P_k} f_i(w) \tag{2}$$

Where k refers to the kth client from K total clients, P_k is the data partition for client k where $nk = |P_k|$ and $f_i(w)$ is the same as before but only specific to a data point within P_k . The goal is to use local optimization algorithms and define the global loss in terms of all of the local subproblems in order to come up with some kind of aggregation technique to merge the results of the local optimization. Traditionally, this is done through the Federated Averaging algorithm. If $n \equiv \sum n_k$, the global empirical loss can be expressed as the following according to [5]:

$$f(w) = \sum_{k=1}^{K} \frac{n_k}{n} F_k(w) = \sum_{k=1}^{K} \frac{n_k}{n} \cdot \frac{1}{n_k} \sum_{i \in P_k} f_i(w)$$
(3)

Finally, the new federated optimization problem becomes

$$\min_{w \in \mathbb{R}^d} \quad f(w) \quad \text{where} \quad f(w) \equiv \sum_{k=1}^K \frac{n_k}{n} F_k(w) \tag{4}$$

2.2 Federated Averaging

As previously discussed, FedAvg is the traditional aggregation algorithm in which each client is weighted according to its corresponding fraction $\frac{nk}{2}$ as seen in (4) ^[8]. In practice, this is done by having each client run a local optimizer such (typically SGD or a variant) and every E epochs, a communication round is done and the aggregation step is performed. This translates to two different weight update steps:

 $w_{t+1}^k \leftarrow w_t^k + \eta \cdot \nabla F_k(w_t)$ or $w_{t+1}^k \leftarrow ClientUpdate(w_t^k)$ Update Step for Client K (5)

 $w_{t+1} \leftarrow \sum_{k=1}^{K} \frac{n_k}{n} w_{t+1}^k$ Update Step for the Server after E epochs (6)

3. Low Rank Metrics

3.1 Low Rank Factorization

Within a Convolutional Neural Network (CNN), the weights of a convolution layer take the form of 4-D tensors of dimension $k \times k \times n_i \times n_o$ where ni is the number of input channels and no is the number of output channels. The weight W_1 for convolution layer lacts as a learned encoder for feature extraction and so ^[3] was able to unfold the tensor W on the input channel or output channel to obtain a 2D matrix $W_1 \in {}^{whno} \times {}^{ni}$ or $W_1 \in {}^{whni} \times {}^{no}$. This matrix is then separated from noise using Variational Bayesian Matrix Factorization to obtain a low-rank singular value decomposition $\hat{W}_1 = U \Lambda V^T$ alongside some noise matrix:

 W_1 [Weight Matrix] $\rightarrow \hat{W}_1$ [Low-Rank Structure] + E_1 [Noise Perturbation]

3.2 Stable Rank

According to [3], the stable rank is defined as the norm energy of the singular values of a given matrix, however they propose a modified definition of stable rank for the low rank m×n matrix \hat{W}_1 given by the following:

$$s(\hat{W}_l) = \frac{1}{n\sigma_1^2(\hat{W}_l)} \sum_{k=1}^{n'} \sigma_k^2(\hat{W}_l)$$
(7)

Where $\sigma_1 \ge \sigma_2 \ge ... \ge \sigma_n'$ are the low rank singular values in descending order and $n \le m$. This measure is normalized by $n\sigma^2(\hat{W}_1)$ and $s(\hat{W}_1)$ [0, 1]. A higher measure is theoretically indicative of a better encoder.

3.3 Temporal Knowledge Gain

We define temporal knowledge gain of client i S for layer l as the difference of stable rank between round t and round t + 1 as evaluated on the low rank weight structure:

$$KG(\hat{W}_{l}^{i}) = s(\hat{W}_{l}^{i,t+1}) - s(\hat{W}_{l}^{i,t})$$
(8)

If this difference is positive, it indicates an increase in stable rank and a relative improvement in themapping done on layer l for client i. The opposite holds true if it is negative and the magnitude of the knowledge gain represents the amount of knowledge lost or gained. Since $s(\hat{M}) \in [0,1]$, KG \in

[-1,1] with 1 being the desirable extreme representing a complete gain of knowledge.

3.4 Uniqueness Score

The goal of this metric is to compare the learning of one client in comparison to how consistent it is with other clients in order to combat client drift. In order to measure this, we use cosine similarity. Consider a set of clients S and two clients j and k S. The low rank weight tensors of layer l are concatenated into vectors in the following manner:

$$\hat{w}_{l}^{j} = \begin{bmatrix} | \\ \hat{W}_{l}^{j} \\ | \end{bmatrix} \text{ and } \hat{w}_{l}^{k} = \begin{bmatrix} | \\ \hat{W}_{l}^{k} \\ | \end{bmatrix}$$

To calculate the similarity between the weights of these two clients we then simply compute

$$Sim(\hat{w}_{l}^{j}, \hat{w}_{l}^{k}) = \frac{\hat{w}_{l}^{j} \cdot \hat{w}_{l}^{k}}{\|\hat{w}_{l}^{j}\| \|\hat{w}_{l}^{k}\|}$$
(9)

We then define the uniqueness score in the following manner:

$$U(\hat{W}_{l}^{i}) = \sum_{k \in S} (1 - Sim(\hat{w}_{l}^{i,t}, \hat{w}_{l}^{k,t}))$$
(10)

Various intuitions about metric's role in aggregation exist due to clients' diverse training on unique samples. There are multiple intuitions at play here in terms of how this metric can be used to guide aggregation. Given that each client is being trained on different samples, each client may develop new patterns that other clients cannot develop. Discovery of these useful patterns leads to rapid divergence by responsible clients. Prioritizing the most unique client for new information may be beneficial. Sometimes, minimizing drift by favoring agreeing clients is useful. The most useful approach may change, but for simplicity, we use the uniqueness score (defined in (10)) for aggregation. Cosine similarity applies to any vectors; low rank factorization isn't strictly needed but aids noise reduction for better results.

4. Proposed Framework

The low rank metrics mentioned are among many possible ones to explore factorized weights and assess learning quality. Our focus is on temporal knowledge gain and similarity score in our experiments. This framework operates per layer, modifying federated averaging to introduce smoothness in weighted averaging and implement aggregation using a specified low rank metric. For general purposes, we will refer to a low rank metric as some measure on the weights $M(W^{\circ} i)$ for layer 1 and client i. That means, in our example, $M(W^{\circ} i)$ can be equal to KG(W^{\circ} i) or U(W^{\circ} i). Note that for a certain round t, $M(W^{\circ} i)$ might include past values for weights from previous rounds such as in the case for knowledge gain and can also of other clients such as is required to calculate uniqueness score. $M(\hat{W}^{i})$ is simply a freehand that shows that the measurement is being done to rank the client itself.

These metrics work layer by layer, converting 4D tensors in convolutional blocks for low rank measurements. Linear layers are weighted similarly to convolutional layers, each requiring individual weighted aggregation, unlike FedAvg and other strategies.

Algorithm 1 Framework for Metric Guided Aggregation 1: Inputs: exponential decay factor β , exponential update factor ζ 2: Initialize: $w^0 = [W_1^0, W_2^0, ..., W_L^0]$ for L layers, $\alpha_l^{i,0} = 1$ for all layers l and clients i w^n : global model parameters at round n. $\alpha_l^{i,n}$: exponential averaging memory variable for layer l, client model i, at round n. 3: for round t = 0, 1, ..., T - 1 do 4: for client $i \in S$ in parallel do $w_i^{t+1} \leftarrow ClientUpdate(w_i^t)$ 5: ▷ Local update step for client $i \in S$ do 6: 7: Initialize $\alpha_p^i = 1$ Keep track of visited conv layer. α_p^i : temporary variable used for keeping track the score of the previous conv layer for layer l = 1, 2, ..., L in $w_i^{t+1} = [W_1^{i,t+1}, W_2^{i,t+1}, ..., W_L^{i,t+1}]$ do 8: if $dim(W_t^{i,t+1}) = 4$ then 9: ▷ If convolutional layer $\alpha_p^i \leftarrow \beta \cdot \alpha_l^{i,t} + \zeta \cdot M(\hat{W}_l^i) \Rightarrow \text{assign the smoothed score to the temp variable}$ 10: $\alpha_l^{i,t+1} \leftarrow \alpha_p^i \qquad \qquad \triangleright \text{ move the temp value into the exponential averaging memory}$ 11: $\bar{\alpha}_l^{i} \leftarrow \frac{N_l}{\sum_{k \in S} N_k} \cdot \alpha_l^{i,t+1} \qquad \triangleright \text{ calculate a normalized score based on number of data}$ 12: samples per client for layer l = 1, 2, ..., L in $w_i^{t+1} = [W_1^{i,t+1}, W_2^{i,t+1}, ..., W_L^{i,t+1}]$ do 13: // Normalize $\bar{\alpha_l}^i$ for each client i and aggregate 14: $W_l^{t+1} \leftarrow \sum_{i \in S} \bar{\alpha}_l^i \cdot W_l^i$ where $\bar{\alpha}_l^i = \bar{\alpha}_l^i / \sum_{i \in S} \bar{\alpha}_l^i$ 15: $w^{t+1} \leftarrow [W_1^{t+1}, W_2^{t+1}, ..., W_L^{t+1}]$ 16: Concatenate updated layers

To prevent discrepancies, the calculated score is normalized for weighted averaging. Multiple metrics can be applied by chaining multiplication (see normalization in Algorithm 1). Scores are often chained with the number of client samples, similar to FedAvg. Note, in the implementation, any zero/negative α values are limited to 1e-3 for stability, ensuring no negative or zero α in aggregation but allowing very small values for specific clients. Furthermore, as seen in line 10 of Algorithm 1, the factor β is introduced in order to increase thesmoothness. This hyperparameter can be tuned in order to obtain the desired stability.

The two algorithms being used in our experimentation will be termed as FedKG where $M(\hat{W}_{1}^{i}) = KG(\hat{W}_{1}^{i})$ where KG is as defined in (8) and FedUnique where $M(\hat{W}_{1}^{i}) = U(\hat{W}_{1}^{i})$ where the Uniqueness Score is as defined in (10).

5. Experiments

We experimented with various algorithms and hyper parameters with our proposed framework and explored how much each algorithm improve the final accuracy in an non-IID dataset setting and how they perform in combination.

5.1 Experimental Setup

This work primarily focuses on image classification tasks, using the widely accepted Resnet-18 model and Adam optimizer. Experiments utilize CIFAR10, CIFAR100, and TinyImageNet200 datasets, partitioned non-IID using Dirichlet distribution. Each subset corresponds to one client's private data. Implementation leverages the Flower framework and Dirichlet partitioning with FedML. FedAvg serves as the baseline for comparison. Hyper-parameter search is constrained to 250 training epochs, with the best combination used in the final experiment. Each experiment averages accuracy from three trials. See Figure 1.



Figure 1: Flowchart of the framework.

5.2 Generating Non-IID Data

Non-IID covers various data distribution types, including horizontal and vertical splitting in sample and label spaces. Our experiment focuses on horizontally split data that is Non-IID in both spaces, resulting in similar content and format among local datasets with varying sample counts per class. We created Non-IID datasets using a customized Dirichlet partitioning algorithm based on FedML. We simulated two Non-IID levels with $\alpha = 0.3$ and $\alpha = 0.0001$, as illustrated in Figure 2.



Figure 2: Non-IID data distribution visualization. Darker color indicates more data per class label per client. This representation reflects the data proportion, not the actual scale of CIFAR100 dataset samples.

5.3 Accuracy Results

We conducted experiments on six datasets: CIFAR10, CIFAR100, TinyImageNet200, with Dirichlet alpha=0.3 and 0.0001. Each algorithm underwent three trials for each setting, and the highest test accuracy per trial was averaged. Test accuracy is the mean of an epoch's accuracy across all four clients. Results are depicted in Figure 3.



Figure 3: Experiment result. The background colour shows the improvement in accuracy compared to the baseline, deeper colour is preferable. The accuracy shown is the average of 3 trials with random initialization.

6. Conclusion

We introduced a federated learning framework enabling layer-wise low-rank metric usage to guide aggregation and explored its effectiveness. However, experiment results didn't reveal a strong pattern. Performance of FedKG and FedUnique compared to FedAvg depends on the dataset and lacks statistical significance (<1%). This paper mainly serves as an exploration record. Additional methods explored but not included can be found in the appendices.

References:

- [1]Daniel J. Beutel, Taner Topal, Akhil Mathur, Xinchi Qiu, Titouan Parcollet, and Nicholas D. Lane. Flower: A friendly federated learning research framework. CoRR, abs/2007.14390, 2020.
- [2]Chaoyang He, Songze Li, Jinhyun So, Mi Zhang, Hongyi Wang, Xiaoyang Wang, Praneeth Vepakomma, Abhishek Singh, Hang Qiu, Li Shen, Peilin Zhao, Yan Kang, Yang Liu, RameshRaskar, Qiang Yang, Murali Annavaram, and Salman Avestimehr. Fedml: A research library and benchmark for federated machine learning. CoRR, abs/2007.13518, 2020.
- [3]Mahdi S Hosseini, Mathieu Tuli, and Konstantinos N Plataniotis. Exploiting explainable metrics for augmented sgd. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pages 10296–10306, 2022.
- [4]Tzu-Ming Harry Hsu, Hang Qi, and Matthew Brown. Measuring the effects of non-identical data distribution for federated visual classification. CoRR, abs/1909.06335, 2019.
- [5]Jakub Konečný, H. Brendan McMahan, Daniel Ramage, and Peter Richtárik. Federated op- timization: Distributed machine learning for on-device intelligence. CoRR, abs/1610.02527, 2016.
- [6]Gihun Lee, Yongjin Shin, Minchan Jeong, and Se-Young Yun. Preservation of the global knowledge by not-true self knowledge distillation in federated learning. CoRR, abs/2106.03097,2021.
- [7]Qinbin Li, Yiqun Diao, Quan Chen, and Bingsheng He. Federated learning on non-iid data silos: An experimental study. CoRR, abs/2102.02079, 2021.
- [8]H. Brendan McMahan, Eider Moore, Daniel Ramage, Seth Hampson, and Blaise Agüera y Arcas. Communication-efficient learning of deep networks from decentralized data. 2016.
- [9]Bjarne Pfitzner, Nico Steckhan, and Bert Arnrich. Federated learning in a medical context: A systematic literature review, 07 2020.