

Mutual Information Driven Federated Learning

Md Palash Uddin, Yong Xiang, *Senior Member, IEEE*, Xuequan Lu, *Member, IEEE*,
John Yearwood, *Senior Member, IEEE*, Longxiang Gao, *Senior Member, IEEE*

Abstract—Federated Learning (FL) is an emerging research field that yields a global trained model from different local clients without violating data privacy. Existing FL techniques often ignore the effective distinction between local models and the aggregated global model when doing the client-side weight update, as well as the distinction of local models for the server-side aggregation. In this paper, we propose a novel FL approach with resorting to mutual information (MI). Specifically, in client-side, the weight update is reformulated through minimizing the MI between local and aggregated models and employing Negative Correlation Learning (NCL) strategy. In server-side, we select top effective models for aggregation based on the MI between an individual local model and its previous aggregated model. We also theoretically prove the convergence of our algorithm. Experiments conducted on MNIST, CIFAR-10, ImageNet, and the clinical MIMIC-III datasets manifest that our method outperforms the state-of-the-art techniques in terms of both communication and testing performance.

Index Terms—Distributed Learning, Federated Learning, Parallel Optimization, Data Parallelism, Information Theory, Mutual Information, Communication Bottleneck, Data Heterogeneity.

1 INTRODUCTION

In this age of digital transformation, there are around 3 billion smartphone users and 7 billion connected Internet of Things (IoT) devices along with a huge number of other computing and embedded devices [1]. These end-edge devices even including standalone smart sensors in distributed networks generate immense amount of user private data everyday. Federated Learning (FL), as an emerging privacy-preserving Machine Learning (ML) technique, can train these private data on-device to enhance intelligence for many real-world applications, such as image analysis, sentiment analysis, clinical care activities (in-hospital mortality prediction or heart disease prediction), burglary analysis, activity analysis, and semantic location [2], [3].

FL is an iterative scheme of training a shared neural network model in parallel over clients' own private data and aggregating the trained models into a global trained model in the server. Specifically, the server-initialized global model is downloaded by the clients for training with their private data. After training, the clients upload their models to the server. The server aggregates the local models to generate an updated global model, which is subsequently downloaded by the clients for retraining again. This process of training and aggregation is repeated until reaching a learning convergence, depicted in Fig. 1.

FL is challenging due to the following two main issues: (1) communication overhead (communication between clients and server), and (2) data heterogeneity (e.g., non-independent and identically distributed). In specific, the well-established classical data center distributed ML strategies [4], [5], [6] would be infeasible, as they just consider

Independent and Identically Distributed (IID) data over a small number of clients. [2], [7] and [8] confirmed these issues and provided detailed discussion on the differences between classical distributed ML and FL. The Federated Averaging (FA) method [7] uses classical Stochastic Gradient Descent (SGD) for model training at the client-side, and further averages the local trained models in the server. Federated Stochastic Variance Reduced Gradient (FSVRG) [8] adopts Stochastic Variance Reduced Gradient (SVRG) [9] for federated optimization with the expensive full gradient calculation. Nevertheless, these two methods are seminal FL works and consider little on the communication overhead and data heterogeneity (These approaches are provided in supplementary document A.). Recent research such as [10] proposes Momentum FA (MFA) which utilizes the momentum Gradient Descent (GD) instead of the conventional GD for accelerating the training speed of models in FA. The authors of [11] further improved their previous work [12], by proposing an updated layer-wise weight averaging for model aggregation, called Federated Matched Averaging (FMA) to replace the coordinate-wise weight averaging of the baseline FA, which requires more additional internal communications between the server and clients for transferring each layer's weights individually. However, these methods are not essentially designed to mitigate the above two issues, because they attempt to either update client-side model training (e.g., [8], [10]) or improve server-side model aggregation (e.g., [11]) at a time, rather than conducting both at the same time.

The existing FL methods often ignore the effective distinction between current local model and the server-sent global model for client-side weight update, i.e., the client replaces its current model by the downloaded global model for training on its private dataset at every communication round. Along with this, all local trained models uploaded to the server are considered indiscriminately for server-side model aggregation at every communication round. As data

• The authors are with the School of Information Technology, Deakin University, Geelong, VIC 3220, Australia (corresponding author: Yong Xiang).

E-mail: {mpuddin, yong.xiang, xuequan.lu, john.yearwood, longxiang.gao}@deakin.edu.au

Manuscript received May 2020; revised September 2020; accepted October 2020

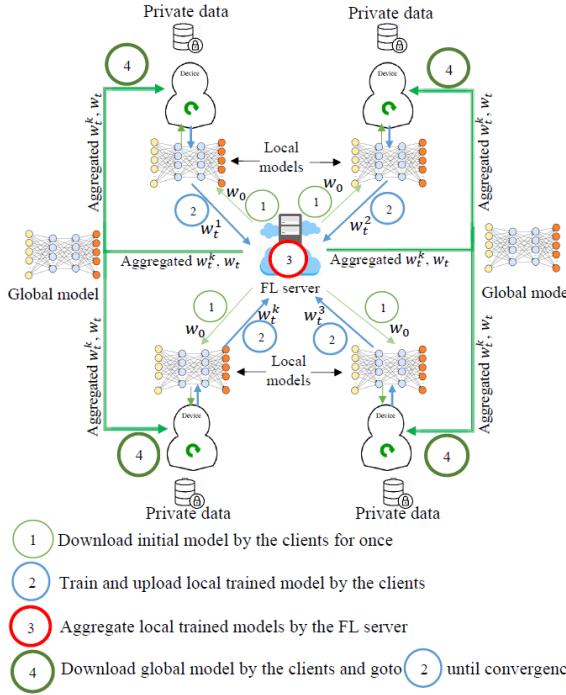


Fig. 1: Interaction among clients and server in FL paradigm. The clients train the downloaded global model ($w_t, t \geq 0$) using their own private data in parallel and upload the trained models to the FL server. The server aggregates these models ($w_t^k, t \geq 1$) and sends the aggregated model to the clients. This process repeats until reaching convergence.

are often not evenly distributed among the clients and are often of non-IID type, the client's current trained model can moderately be utilized to minimize its training loss when training the server-sent global model. Also, the training performance of all local models might not be effective, that is, some local models can be outliers in respect to the overall training objective.

To address the aforementioned issues, we, therefore, propose a novel FL approach in this paper. In particular, we propose to update the model weights at client-side by exploiting the Mutual Information (MI) between the current local model and the shared global model. To enable effective selection of local trained models for aggregation, we further attempt to measure the MI between local models and the shared global model, and choose the local models based on the computed MI values. We also show the convergence analysis for our algorithm, and visualize the learning progresses through MI. Extensive experiments on IID and non-IID data have been conducted on different datasets (i.e., MNIST, CIFAR-10, ImageNet, and MIMIC-III) to validate the proposed method. Results show that it outperforms state-of-the-art techniques, in terms of communication speedup and testing performance. The main contributions of this paper are

- an MI-driven federated optimization approach,
- a novel model update algorithm for local models,
- a novel model aggregation algorithm, and
- a theoretical convergence analysis for our method

The remainder of this paper is formulated as follows. We

provide the related works in Section 2. Section 3 describes the proposed FL approach systematically. Subsequently, we explain the theoretical convergence analysis of the proposed optimization strategy in Section 4. Section 5 analyzes and compares the experimental outcomes whereas Section 6 summarizes the observations and accomplishes the paper.

2 RELATED WORK

In this section, we only review the most relevant researches to this work. Federated Learning (FL) dates back to several years ago, and has attracted increasing attention since then. The FA, i.e., vanilla FL algorithm, performed the classical SGD parallelism for client-side model training and averaged the local trained models for the server-side aggregation [7]. FSVRG was presented in [8] that aims to utilize a constant learning rate and keep tracking of the global model for federated optimization where all clients' participation is required to calculate the expensive full gradients using SVRG [9] at each communication round. Recently, [13] focused on optimizing the structure of Deep Neural Network (DNN) models used in FA training through the adaptation of genetic algorithms. [14] introduced FedProx which modifies FA substantially through adding a proximal term to client's local objective function. MFA was proposed in [10] which suggests to use the momentum GD instead of the classical GD and thus achieves faster training of the local models in comparison to the FA. A self-balancing FL approach named Astraea was proposed in [15] which relieves global imbalance through data augmentation and downsampling operations and local imbalance through creating mediators. It uses Kullback–Leibler divergence of clients' data for rescheduling the training process of the clients. The work in [12] was further improved by [11] which presents a layer-wise weight averaging algorithm called FMA for model aggregation instead of the coordinate-wise weight averaging of FA. As a result, FMA outperforms FA and FedProx, but demands more internal communications between the server and clients for transferring each layer's weights.

3 PROPOSED APPROACH

3.1 Approach Overview

Our approach consists of two steps: model weights update and aggregation. For the first step (called "MIFL-Update"), we propose to update model weights by exploiting the MI for current local model and the aggregated shared model. This is to discriminate their contributions. Regarding the second step (called "MIFL-Aggregation"), we propose to rank the local models based on the MI measurement and pick up the top effective models for aggregation. The overall approach termed as MIFL is shown in Fig. 2.

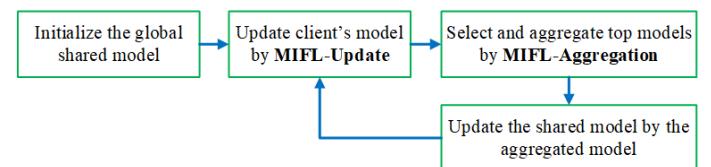


Fig. 2: Overview of the proposed MIFL.

3.2 Model Weights Update

Unlike vanilla FA, both the shared global model and the existing local model at each client are emphasized for next weight update in MIFL-Update approach. The goal is to optimally utilize individual local device's current trained model for minimizing overall training loss by the aggregated global model using the private dataset of each client. This can be defined as

$$\min_{w \in R^d} f(w) = \sum_{k=1}^K \frac{N_k}{N} L_k(w), \quad (1)$$

where $L_k(w) = \frac{1}{N_k} \sum_{i \in M_k} l(w, i)$ represents the k^{th} client loss calculated using the current model w over N_k samples. K is the total number of participating clients, and N is the total number of data samples of the clients. $N_k = |M_k|$ i.e., M_k is the set of data indices whose size is N_k , and $l(w, i)$ is the loss for the i^{th} sample.

We reformulate the existing loss calculation strategy for conventional SGD algorithm [16] by exploiting the error measure between current local model and the aggregated global model. Specifically, we assume that the output F_k of current local network, k and the output F_g for the aggregated global network, g follow Gaussian functions with respective variances σ_k^2 and σ_g^2 . Then, the MI between F_g and F_k , denoted by $MI(F_g; F_k)$, is defined as [17], [18], [19]

$$MI(F_g; F_k) = H(F_g) + H(F_k) - H(F_g, F_k), \quad (2)$$

where $H(F_g)$ is the entropy of F_g , $H(F_k)$ is the entropy of F_k , and $H(F_g, F_k)$ is the joint entropy of F_g and F_k , which are defined as

$$\begin{aligned} H(F_g) &= \frac{1}{2} [1 + \log(2\pi\sigma_g^2)] \\ H(F_k) &= \frac{1}{2} [1 + \log(2\pi\sigma_k^2)] \\ H(F_g, F_k) &= 1 + \log(2\pi) + \frac{1}{2} \log[\sigma_g^2 \sigma_k^2 (1 - \rho_{gk}^2)] \end{aligned} \quad (3)$$

In Eq. (3), the correlation coefficient, ρ_{gk} between F_g and F_k can be computed as

$$\rho_{gk} = \frac{E[(F_g - E[F_g])(F_k - E[F_k])]}{\sigma_g^2 \sigma_k^2}, \quad (4)$$

where $E[.]$ denotes the mathematical expectation. Substituting Eq. (3) into Eq. (2), we get

$$MI(F_g; F_k) = -\frac{1}{2} \log(1 - \rho_{gk}^2). \quad (5)$$

From Eq. (5), it can be observed that, i) when the outputs F_g and F_k are uncorrelated, the correlation ρ_{gk} is zero and their MI becomes very small; ii) when the outputs F_g and F_k become highly positively correlated, the correlation ρ_{gk} is close to 1 and their MI tends to be noticeably large. To enable the update using client's current local and aggregated global networks, the correlation coefficient ρ_{gk} between the current local network and the aggregated global network should be minimized, which is equivalent to minimizing the MI between them, to extract distinct information about the inputs [17], [18], [19]. Supposing that the aggregated global model represents local models of other clients, we formulate this minimization problem via introducing a correlation penalty term in the traditional loss function in Negative Correlation

Learning (NCL) [17], [18], [19]. The final classical loss function E_g , a typical notion of $l(w)$, is therefore

$$\begin{aligned} E_g &= \frac{1}{N_k} \sum_{i=1}^{N_k} E_g(x_i) \\ &= \frac{1}{N_k} \sum_{i=1}^{N_k} [\frac{1}{2} (F_g(x_i) - y_i)^2 + \lambda \alpha_g(x_i)], \end{aligned} \quad (6)$$

where $E_g(x_i)$ denotes the loss function for the i^{th} sample which is represented by sample-label pair as (x_i, y_i) and λ is a threshold ($0 \leq \lambda \leq 1$). $\alpha_g(x_i)$ is the correlation penalty for the x_i^{th} sample that can be defined as

$$\alpha_g(x_i) = -\frac{1}{2} (F_g(x_i) - F_k(x_i))^2. \quad (7)$$

To enable effective training of FL local model at client-side through minimizing α_g , we take the partial derivative of Eq. (6) with respect to F_g , and obtain

$$\begin{aligned} \frac{\partial E_g(x_i)}{\partial F_g(x_i)} &= F_g(x_i) - y_i - \lambda (F_g(x_i) - F_k(x_i)) \\ &= (1 - \lambda)(F_g(x_i) - y_i) + \lambda(F_k(x_i) - y_i), \end{aligned} \quad (8)$$

where λ is a regularization parameter, which is used to dispose the strength of the penalty. Eq. (8) is then added to the standard SGD algorithm for minimization [17], [18], [19]. Note that we also provide the derivation for the cross-entropy based loss in our supplementary document B.

Notice that we need to make a bit adjustment for the first training round of the clients. After downloading the initial or intermediate aggregated model, a client takes it as its local model for training with conventional loss calculation function of SGD, following the vanilla FA [7]. Our MIFL-Update approach is summarized in Algorithm 1.

3.3 Model Aggregation

Due to data heterogeneity or data discrepancy among clients, some local trained models can be similar to others, and some other models can be greatly different from others. Over-working models with excessively high performance and primitive models with excessively low performance are collectively called outlier models. It is obvious that local models with similar performance can increase redundancy in models' aggregation and also increase communication overhead significantly [2]. As such, we propose to quantitatively estimate the performance of all local trained models and rank their priorities before determining the models to be uploaded to the server for aggregation. Inspired by the input features ranking for different vision and visual tasks based on MI [20], [21], [22], [23], we consider the local models as features to be selected for the aggregation task. This way finds an effective subset of local trained models for aggregation, and the communication cost among the clients and server can be thus reduced.

Given the trained local models via Algorithm 1, the MI between the outputs of the local trained network and the previous global network is calculated via Eq. (5), where the correlation coefficient ρ_{gk} is

$$\rho_{gk} = \frac{\sum_{i=1}^{N_k} (F_g(x_i) - E[F_g(x_i)])(F_k(x_i) - E[F_k(x_i)])}{\sqrt{\sum_{i=1}^{N_k} (F_g(x_i) - E[F_g(x_i)])^2 (F_k(x_i) - E[F_k(x_i)])^2}}. \quad (9)$$

Algorithm 1 MIFL-Update.

Input: Server-sent shared model, w_t at round t
Output: Updated model, w_{t+1}^k and MI, MI_{t+1}^k

- 1: **procedure** *MIFL-Update* (k, t, w_t) : ▷ Run on FL client k
- 2: **if** k^{th} 's first update, **then**
- 3: Set local model, $F_k := w_t$
- 4: Set client's model, $F_g := w_{t+1}^k := w_t$ to be trained
- 5: $\mathbb{B} :=$ (Split training data into batches of size B)
- 6: **for each** local epoch i from 1 to E **do**
- 7: **for each** batch $b \in \mathbb{B}$ **do**
- 8: $x(b) :=$ Load data points of b^{th} batch
- 9: $y(b) :=$ Load true labels for $x(b)$
- 10: $\Omega := (F_g(x(b)) - y(b))$ ▷ Classical loss calculation
- 11: $w_{t+1}^k := w_{t+1}^k - \eta\Omega$ ▷ Local model's weight update
 ▷ η is the learning rate
- 12: **end for**
- 13: **end for**
- 14: **else**
- 15: Set local model, $F_k := w_t^k$ ▷ w_t^k is client's previous model
- 16: Set client's model, $F_g := w_{t+1}^k := w_t$ to be trained
- 17: $\mathbb{B} :=$ (Split training data into batches of size B)
- 18: **for each** local epoch i from 1 to E **do**
- 19: **for each** batch $b \in \mathbb{B}$ **do**
- 20: $x(b) :=$ Load data points of b^{th} batch
- 21: $y(b) :=$ Load true labels for $x(b)$
- 22: $l_g := F_g(x(b)) - y(b))$
- 23: $l_k := F_k(x(b)) - y(b))$
- 24: $c_v := \frac{\sigma(l_g, l_k)}{\sum(l_g, l_k)}$
- 25: **if** $l_g = l_k$ **then**
- 26: $\lambda := \frac{1}{2}$
- 27: **else**
- 28: **if** $c_v <= \frac{1}{2}$ **then**
- 29: $\lambda := c_v$
- 30: **else**
- 31: $\lambda := 1 - c_v$
- 32: **end if**
- 33: **end if**
- 34: $\Omega := (1 - \lambda)(F_g(x(b)) - y(b)) + \lambda(F_k(x(b)) - y(b))$ ▷ Modified loss calculation
- 35: $w_{t+1}^k := w_{t+1}^k - \eta\Omega$ ▷ Local model's weight update
 ▷ η is the learning rate
- 36: **end for**
- 37: **end for**
- 38: **end if**
- 39: $F_g := w_{t+1}^k$
- 40: $MI_{t+1}^k := MI(F_g; F_k)$
- 41: **return** w_{t+1}^k, MI_{t+1}^k
- 42: **end procedure**

The MI between the outputs of the local model and global model is uploaded to the server. The server ranks the local trained models' MI values in an increasing order to identify the outlier models (over-working and primitive).

Algorithm 2 MIFL-Aggregation.

Input: Clients' models, w_{t+1}^k and MI, MI_{t+1}^k , and targeted accuracy, $T_{accuracy}$
Output: Aggregated model, w_{t+1}

- 1: **procedure** *MIFL-Aggregation()* : ▷ Run on FL server
- 2: Initialize w_0
- 3: **for each** round $t := 0, 1, 2, \dots, T - 1$ **do**
- 4: **if** $T_{accuracy}$ not achieved **then**
- 5: Choose $U := \lceil (C \times K) \rceil$ clients, $0 < C \leq 1$
- 6: **for each** client $k \in U$ in parallel **do**
- 7: $(w_{t+1}^k, MI_{t+1}^k) := MIFL-Update(k, t, w_t)$
- 8: **end for**
- 9: $K_{Top} :=$ Compute the set of effective clients
 pruning outliers based on MI_{t+1}^k
- 10: $w_{t+1} = \sum_{k=1}^{K_{Top}} \frac{N_k}{N_{Top}} w_{t+1}^k$ ▷ Aggregation of effective models
- 11: Send w_{t+1} to outlier clients
- 12: **else**
- 13: **break**
- 14: **end if**
- 15: **end for**
- 16: **end procedure**

tive models). A lower MI means a better learning as it forces the network to have different weights with different inputs. The models with excessively low MI values are over-working, while the models with too high MI are thus primitive. The server picks up K_{Top} top effective models (say, $w_{t+1}^k, k \in [1, K_{Top}]$) by eliminating the over-working and the primitive models for aggregation. The aggregation manner is then similar to vanilla FA, shown in Eq. (10).

$$w_{t+1} = \sum_{k=1}^{K_{Top}} \frac{N_k}{N_{Top}} w_{t+1}^k, \quad (10)$$

where w_{t+1} is the aggregated global model and N_{Top} represents the total number of samples of K_{Top} clients. Notice that the selection of K_{Top} local models does not affect the number of total parameters of the global model, and it influences the values of the parameters of the global model. This is because the outlier models are not used for calculating the global model, to alleviate the misleading of the overall training process and the delay of achieving the convergence.

After aggregation, the global model is sent to the newly selected clients to train this aggregated shared model along with their respective local models using the MIFL-Update algorithm (Algorithm 1). As with FA [7], a small fraction of total clients is randomly selected as the new clients. However, this updated global model is also delivered to those clients which are over-working and primitive in this round, for replacing their previous global model of the previous round. Similarly, the clients' updated models are ranked and selected for the next aggregation. This process of training and aggregation is repeated until a targeted learning convergence is reached. Algorithm 2 shows the model aggregation procedure.

4 CONVERGENCE ANALYSIS

To analyse the convergence of our proposed federated optimization algorithm (MIFL), we focus on the IID linear regression task where each sample is represented by (x_i, y_i) and the associated classical loss is calculated as $l(w, i) = (y_i - x_i^\top w)^2$. Notice that non-IID data appear to be diverse, such as class imbalance, size imbalance and feature imbalance. As a result, it involves various factors and still requires researchers' significant efforts on modelling the non-IID data setting in FL. That is why we do not carry out the convergence analysis on the non-IID data, and we believe that the convergence analysis on non-IID data will be a promising research direction in FL context in future.

4.1 Setting

Suppose that the conventional loss function, $l : \mathbb{R}^d \rightarrow \mathbb{R}$ in linear regression is defined as follows.

$$l(w) = \frac{1}{2} \mathbb{E}_{(x,y) \sim \mathbb{P}} (y - w^\top x)^2, \quad (11)$$

where \mathbb{P} denotes a fixed distribution on $\mathbb{R}^d \times \mathbb{R}$, \mathbb{P}_X represents the marginal distribution on \mathbb{R}^d and $\mathbb{P}(\cdot|x)$ represents a conditional distribution of \mathbb{R} given $x \in \mathbb{R}^d$. In the FL setup on IID data distribution, each client is assumed to have access to independent samples from \mathbb{P} and the goal is to minimize the overall loss l . In our proposed MIFL, the MIFL-Update algorithm calculates $l(w) = (1 - \lambda)(F_g(x - y)) + \lambda(F_k(x - y))$ through an adaptive weighted-combination (by λ) of losses of a client's current local model (F_k) and the aggregated global model (F_g). While *Client_Update()* of vanilla FA calculates $l(w) = (F_g(x - y))$ using F_g only. Besides, the MIFL-Aggregation algorithm of MIFL selects priority models to be effectively aggregated in each round while *Server_Compuation()* of vanilla FA takes all models equally for aggregation.

4.2 Assumptions

The first assumption is about the uniformity of clients' upload time [24], the second is about marginal \mathbb{P}_X and the last one is about conditional $\mathbb{P}(\cdot|x)$ [25].

4.2.1 Assumption 1

Each selected client uploads its trained model within a predefined allocated time to the server for aggregation.

4.2.2 Assumption 2

(i) For few $R > 0$, $\mathbb{E}_{x \sim \mathbb{P}_X} [|x|^2 x x^\top] \preceq R^2 \mathbb{E}[x x^\top]$ exists, where the representation $A \preceq B$ signifies that $B - A$ is positive semidefinite for symmetric matrices $A, B \in \mathbb{R}^{d \times d}$. Particularly, this relationship exists when $\|x\| \leq R$ in the marginal distribution \mathbb{P}_X . (ii) The lowest eigenvalue $\mu = \Lambda_{\min} \mathbb{E}_{x \sim \mathbb{P}_X} [x x^\top]$ of the hessian $\mathbb{E}_{x \sim \mathbb{P}_X} [x x^\top]$ of l is strictly positive.

In linear regression, (i) is a common assumption [25], while (ii) implies that l is μ -strongly convex, where a condition number κ is defined as $\kappa = \frac{R^2}{\mu}$.

4.2.3 Assumption 3

There holds a vector $w_* \in \mathbb{R}^d$ for which for all $x \in \mathbb{R}^d$, $y(x)$ distributed through $\mathbb{P}(\cdot|x)$ complies with $y(x) = x^\top w_* + \zeta$. ζ is a random variable complying: (i) ζ is not dependant on x , (ii) $\mathbb{E}\zeta = 0$, and (iii) $\mathbb{E}\zeta^2 \leq \sigma^2$ for some $\sigma^2 < \infty$.

If all three assumptions hold, then w_* implies that it is an unique minimizer of the overall loss function l . Moreover, it is clear from Assumption 3 that as ζ is independent of x , it represents that a well-specified model is produced.

4.3 Convergence Result

We now attempt to express the convergence of MIFL with tail-averaged SGD as *MIFL-Update()*, while the aggregation approach *MIFL-Aggregation()* is expected to produce a result in $\text{conv}\{w^1, \dots, w^K\}$ (K is the total number of clients).

4.3.1 Proposition

Considering l from Eq. (11), we suppose that the above assumptions exist. Also, we suppose the aggregation procedure *MIFL-Aggregation()* $\in \text{conv}\{w^1, \dots, w^K\}$ for all $w^1, \dots, w^K \in \mathbb{R}^d$, and $E > (8 \log 2)\kappa \log \kappa$. Then, the aggregation output w_T using K_{Top} clients' models, $\eta = (1/2R^2)$ and the amount of SGD iterations $E_t = E$ for all t , satisfy the following [25], [26]

$$\mathbb{E}l(w_T) - l(w_*) \leq 2^{-T}(l(w_0) - l(w_*)) + \frac{16d\sigma^2}{E}. \quad (12)$$

Alternatively, if we consider local SGD epochs $E_t = 2^t E$, then we have,

$$\mathbb{E}l(w_T) - l(w_*) \leq 2^{-T}(l(w_0) - l(w_*)) + \frac{8d\sigma^2 T}{2^{T-1} E}. \quad (13)$$

4.3.2 Proof

Suppose that ℓ_t is sigma algebra produced by w_t and w_t^k is the updated model of each selected client. The output model of *MIFL-Update()* having batch size $B = 1$ complies with the following guarantee [25].

$$\mathbb{E}[l(w_t^k) | \ell_t] - l(w_*) \leq 2e^{-\frac{E_t}{8\kappa \log \kappa}} (l(w_t) - l(w_*)) + \frac{8d\sigma^2}{E_t}. \quad (14)$$

As $E_t > (8 \log 2)\kappa \log \kappa$ for each t , we get,

$$\mathbb{E}[l(w_t^k) | \ell_t] - l(w_*) \leq \frac{1}{2}(l(w_t) - l(w_*)) + \frac{8d\sigma^2}{E_t}. \quad (15)$$

Since our aggregation strategy *MIFL-Aggregation()* picks a point in the convex hull of $\{w_t^1, \dots, w_t^k\}$, it satisfies the convexity of l as follows.

$$l(w_{t+1}) \leq l(w_t^k), \quad (16)$$

for each client $k, k \in K_{Top}$. If we take the expectation of (15) over ℓ_t and use (16), we then get,

$$\mathbb{E}l(w_{t+1}) - l(w_*) \leq \frac{1}{2}(\mathbb{E}l(w_t) - l(w_*)) + \frac{8d\sigma^2}{E_t}. \quad (17)$$

To this end, we unfold this sum over t global epochs and obtain

$$\mathbb{E}l(w_T) - l(w_*) \leq 2^{-T}(l(w_0) - l(w_*)) + 8d\sigma^2\left(\frac{1}{E_{T-1}} + \frac{1}{2E_{T-2}} + \dots + \frac{1}{2^{T-1}E_0}\right). \quad (18)$$

When $E_t = E$ for each t , then the series in the second term of right hand side of Eq. (18) is upperly bounded by $2/E$, which finishes the proof of Eq. (12). On the other hand, when $E_t = 2^t E$, then the second term sums to $8d\sigma^2 T/(2^{T-1}E)$, which thus proves Eq. (13).

4.4 Discussion

- It can be seen that the bounds in the proposition do not depend on the number of fraction of clients (C) participated in the optimization, which entails that a large number of clients may be engaged for training convergence. Particularly, the bounds convey that the model aggregation in convex hull cannot be worse than executing same optimization method with a single client having $\frac{1}{U^t}$ of the entire data.
- When Assumption 3 does not exist i.e., if the model is mis-specified, similar convergence can be achieved with a tiny learning rate η and a great number of local iterations as $E_t = 2^t E$. This is practical because the computation cost of local epochs is relatively inexpensive in comparison to the communication [25].
- It is evident from [27] that combining losses is faster at reaching convergence. The proposed MIFL-Update also benefits from this since it combines losses from both client's current local and global models to calculate $l(w)$ (Eq. (11)). Since MIFL-Aggregation picks the most effective local models using MI values for the aggregation, Eq. (16) must hold for the proposed method, and we can conclude that our method converges.

5 EXPERIMENTAL RESULTS

5.1 Experimental settings

To better interpret the effectiveness of our proposed MIFL, we inherit and modify a bit the experimental settings in [7]. In particular, we use similar model architectures i.e., Multi-Layer Perceptron (MLP) and Convolutional Neural Network (CNN) with same specifications for similar image datasets. We first consider to extend parallelism using different values of client fraction C and then increase computations in each client. In contrast to [7], we evaluate a real-world clinical critical care dataset-MIMIC-III [28] and a large image classification dataset-ImageNet [29]. Besides, we visualize the training and testing progresses for analysing the learning behaviour from the MI perspective and consider noise free outputs for calculating the MI between the local model and global model.

Datasets. The widely used image classification datasets MNIST, CIFAR-10, and ImageNet along with MIT's health care dataset MIMIC-III have been utilized to validate the superiority of the proposed MIFL over vanilla FA. The parametric information of the datasets is provided in

Table 1. There are 26 tables in MIMIC-III such as patients' admissions, diagnosis codes. Among them, we extract three raw tables (ADMISSIONS, PATIENTS and PRESCRIPTIONS) for our experiments. We first form two new tables from the three ones. The first table named PERSONAL_INFORMATION is made through an inner-join of ADMISSIONS and PATIENTS, and contains AGE_GROUP, GENDER and MORTALITY (survival status) of all patients. The second table called SUBJECT_DRUG_TABLE is extracted from PRESCRIPTIONS, and contains the usage of drugs during first 48 hours of stay (STARTDATE-ENDDATE=48 hours) of each patient. The two tables are further joined on SUBJECT_ID to prepare a dataset of 30,000 samples, from which we randomly select 22,500 samples for the training set and the rest is the test set. In this dataset, DRUGS involving 1,473 different drug features (e.g., D5W, Heparin Sodium, Nitro-glycerine, Insulin, Atropine Sulphate etc.) prescribed to the patients are used as data points while MORTALITY is the class label to be predicted.

TABLE 1: Datasets.

Dataset	Training set	Testing set	Labels	Channels	Resolution or features
MNIST	60K	10K	10	1	28×28
CIFAR-10	50K	10K	10	3	32×32
ImageNet	1.2M	50K	1000	3	224×224
MIMIC-III	22.5K	7.5K	2	1	1473

Model structure. Similar to the baseline FA [7], we employ two commonly used DNN model categories i.e., MLP and CNN with similar architecture. For the MNIST digit recognition task, we use 2-hidden layers of MLP (termed as 2NN MNIST), with 200 units in each hidden layer and ReLu activation (network parameters: 199,210), as well as CNN with two 5×5 convolutional layers, a fully connected layer with 128 units having ReLu activation (network parameters: 454,922) and a final 10-class softmax output layer for both models. The numbers of channels in the convolutional layers are respectively 32 and 64, and each convolutional layer is followed by a 2×2 max pooling. For MIMIC-III, 3-hidden layers of MLP (denoted as 3NN MIMIC) with respectively 200, 100 and 50 neurons are implemented with the ReLu activation, and finally a single-neuron sigmoid output layer is employed (network parameters: 320,051). For the CIFAR-10 dataset, we use CNN with three 3×3 convolutional layers, a fully connected layer with 64 units involving ReLu activation and a final 10-class softmax output layer (network parameters: 122,570). The numbers of channels in the convolutional layers are 32, 64 and 64, respectively, while each convolutional layer is followed by a 2×2 max pooling. Finally for ImageNet dataset, we employ the mostly used network architecture called ResNet-50 (denoted as ResNet ImageNet; network parameters: 25,557,032) [30].

Data distribution. To validate the FL algorithms, it is important to discuss the data distributions (IID and non-IID) over the participating clients. Aligned with [7], the MNIST dataset is assumed to be both IID and non-IID, CIFAR-10 and ImageNet IID, and MIMIC-III non-IID. For the IID MNIST, the entire 60,000 training samples are shuffled and then partitioned into 100 clients, each of which has 600 samples of all classes. For the non-IID MNIST, the training samples are sorted by digit labels and divided into 200

TABLE 2: Impact of fraction of clients C . (a) 2NN MNIST with $E=1$. (b) 2NN MNIST and CNN MNIST, with $E=5$. Each cell specifies the communication round numbers required to produce the targeted test-set accuracy. Speedups of our MIFL in parentheses are calculated based on the corresponding cell results of FA.

Data distribution	B	Method	2NN MNIST, $E=1$: targeted test-set accuracy: 97%				
			C				
			0.0	0.1	0.2	0.5	1.0
IID	∞	FA	1298	1092	1172	1579	1931
		MIFL	1081 (1.20 \times)	792 (1.38 \times)	881 (1.33 \times)	1137 (1.39 \times)	1138 (1.70 \times)
	10	FA	82	58	53	51	50
		MIFL	59 (1.39 \times)	46 (1.26 \times)	39 (1.36 \times)	40 (1.28 \times)	39 (1.28 \times)
Non-IID	∞	FA	4387	1652	1594	1854	- (86.54% at 1929)
		MIFL	3142 (1.40 \times)	1329 (1.24 \times)	1145 (1.39 \times)	1349 (1.37 \times)	1535 (Immeasurable \times)
	10	FA	3500	261	197	182	183
		MIFL	3492 (1.00 \times)	150 (1.74 \times)	129 (1.53 \times)	115 (1.58 \times)	141 (1.30 \times)

(a)

Data distribution	B	Method	2NN MNIST, $E=5$					CNN MNIST, $E=5$				
			Targeted test-set accuracy: 97%					Targeted test-set accuracy: 99%				
			C					C				
IID	∞	FA	1109	871	812	843	1224	258	248	288	183	181
		MIFL	610 (1.82 \times)	538 (1.62 \times)	540 (1.50 \times)	591 (1.43 \times)	1092 (1.12 \times)	78 (3.31 \times)	81 (3.06 \times)	96 (3.00 \times)	63 (2.90 \times)	66 (2.74 \times)
	10	FA	397	47	39	34	34	21	18	17	20	24
		MIFL	376 (1.06 \times)	33 (1.42 \times)	26 (1.50 \times)	23 (1.48 \times)	23 (1.48 \times)	11 (2.00 \times)	7 (2.57 \times)	9 (1.89 \times)	8 (2.50 \times)	12 (2.00 \times)
Non-IID	∞	FA	1196	1096	1147	1396	1466	1172	1007	825	1038	980
		MIFL	965 (1.24 \times)	896 (1.22 \times)	854 (1.34 \times)	993 (1.41 \times)	1132 (1.30 \times)	686 (1.71 \times)	320 (3.15 \times)	350 (2.36 \times)	417 (2.49 \times)	402 (2.44 \times)
	10	FA	3497	134	117	103	96	1091	239	201	295	180
		MIFL	3496 (1.00 \times)	90 (1.49 \times)	73 (1.60 \times)	61 (1.69 \times)	58 (1.66 \times)	552 (1.98 \times)	157 (1.52 \times)	120 (1.68 \times)	129 (2.29 \times)	136 (1.32 \times)

(b)

fragments of size 300. After that, 2 fragments are assigned to each of 100 clients that signify skewed non-IID manner as each client has only samples of two classes. An instance of IID and non-IID distribution for a set of 10 clients is depicted in Fig. 3. For the IID CIFAR-10 and ImageNet, the entire training datasets are divided among 100 clients so that each client has 500 and 12,000 training examples, respectively. However, for the non-IID MIMIC, the training set is sorted according to AGE_GROUP and GENDER, and is equally split and assigned to 100 clients.

are tuned consistently for both existing FL algorithms and ours, to allow a better understanding of different methods. λ is empirically tuned based on the models' performance (observing the error values of global and local models), which would induce a decent contribution of a client's local model in training the global model.

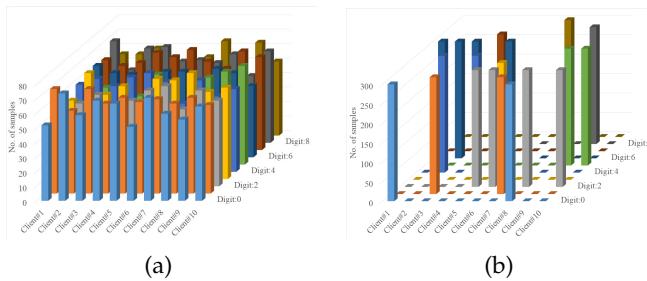


Fig. 3: Data distribution on MNIST for 10 typical clients used in the experiments. (a) IID case, and (b) Non-IID case.

Hyperparameters tuning. There are some model-oriented hyperparameters, e.g., learning rate (η), local batch size (B) and local epochs (E), some existing federated optimization hyperparameters like fraction of clients (C), and our hyperparameter (λ in Eq. (8)). For η , we simply set it to 0.1 and 0.001 for image datasets and MIMIC dataset respectively, based on the naive-approach recommendation and some recent works such as [13] and [3]. B , E and C

Evaluation metric. The mitigation of communication overhead is evidenced by the reduction of communication rounds to reach a targeted test-set accuracy for any data distribution setups; meanwhile, the mitigation of data heterogeneity issue is achieved implicitly. Particularly, the non-IID data signifies an instance of data heterogeneity difficulty, where the training data have been assigned among the clients in an imbalanced way. In general, the non-IID data create more difficulties for the training algorithms and thus demand comparatively more training rounds to reach the convergence than IID data. A potential way to prove that the non-IID data problem (heterogeneity) has been mitigated is to provide faster training convergence. Therefore, the reduction of communication rounds also serves a proof for mitigating heterogeneity. As the proposed MIFL-Update and MIFL-Aggregation algorithms respectively follow the local model training and global model aggregation procedures of the baseline FA, our MIFL has no additional overall overhead in comparison to the vanilla-FA. Note that the modified loss function in MIFL involves the additional output calculation of client's current model to train the server-sent global model, which may increase negligible computation according to our empirical observations.

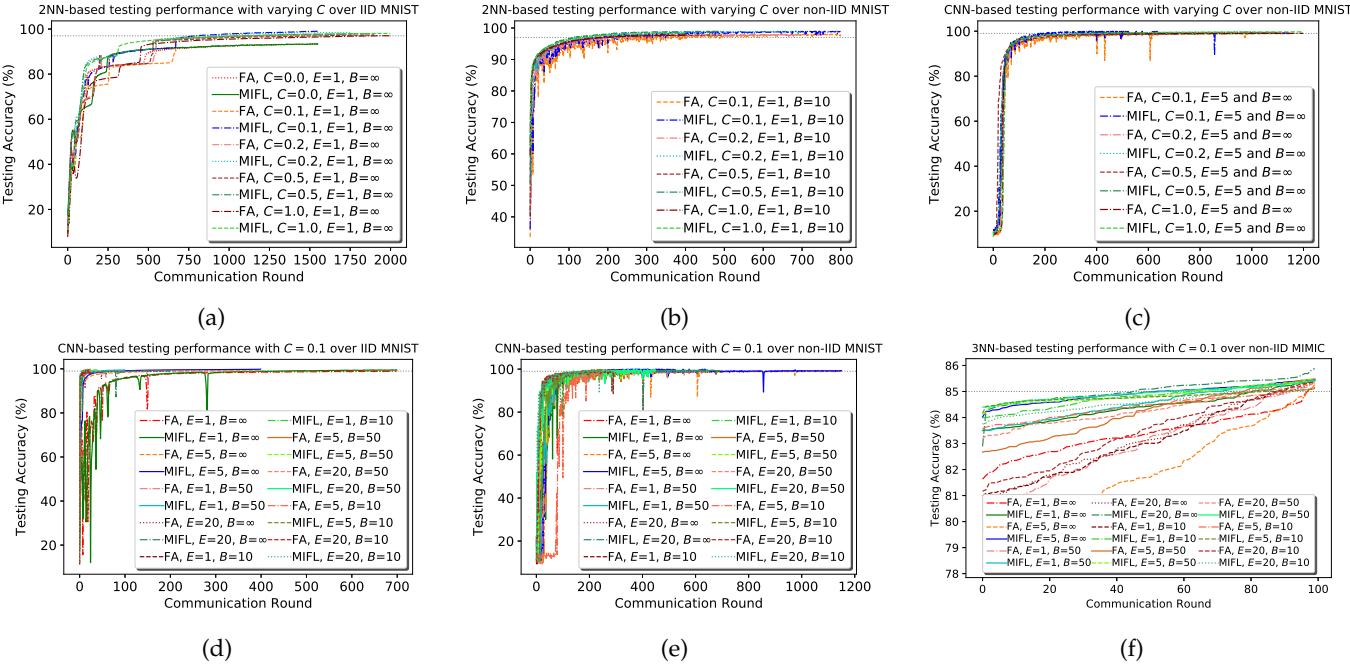


Fig. 4: Test-set accuracies versus communication rounds. (a) Varying C , $E = 1$ and $B = \infty$ over IID MNIST using 2NN. (b) Varying C , $E = 1$ and $B = 10$ over non-IID MNIST using 2NN. (c) Varying C , $E = 5$ and $B = \infty$ using CNN over non-IID MNIST. (d) $C = 0.1$, and varying E and B using CNN over IID MNIST. (e) $C = 0.1$, and varying E and B using CNN over non-IID MNIST. (f) $C = 0.1$, and varying E and B using 3NN over non-IID MIMIC. The targeted accuracies are marked by grey lines.

TABLE 3: Numbers of communication rounds required to reach the targeted test-set accuracy for both MIFL and FA. The value of $u = EN/(KB)$ signifies the expected number of local updates per communication round.

CNN MNIST, targeted test-set accuracy: 99%				3NN MIMIC-III, targeted test-set accuracy: 85%			
E, B, u	FA		MIFL	E, B, u	FA		MIFL
	IID	Non-IID	IID		Non-IID	Non-IID	Non-IID
1, ∞ , 1	700	417	431 (1.62 \times)	318 (1.31 \times)	1, ∞ , 1	98	82 (1.20 \times)
5, ∞ , 5	248	1007	81 (3.06 \times)	320 (3.15 \times)	5, ∞ , 5	98	53 (1.85 \times)
1, 50, 12	57	530	38 (1.50 \times)	311 (1.70 \times)	1, 50, 4.5	93	69 (1.35 \times)
20, ∞ , 20	187	633	108 (1.73 \times)	229 (2.76 \times)	20, ∞ , 20	93	46 (2.02 \times)
1, 10, 60	25	269	16 (1.56 \times)	193 (1.39 \times)	1, 10, 22.5	92	64 (1.44 \times)
5, 50, 60	28	214	17 (1.65 \times)	181 (1.18 \times)	5, 50, 22.5	79	66 (1.20 \times)
20, 50, 240	31	393	16 (1.94 \times)	190 (2.07 \times)	20, 50, 90	81	56 (1.45 \times)
5, 10, 300	18	239	7 (2.57 \times)	157 (1.52 \times)	5, 10, 112.5	92	85 (1.08 \times)
20, 10, 1200	12	275	6 (2.00 \times)	139 (1.98 \times)	20, 10, 450	89	75 (1.19 \times)

5.2 Parallelism

We first consider to increase multi-client parallelism through experimenting on different values of client fraction (C) for the proposed MIFL and existing FA optimization approaches. The results for both 2NN and CNN models over MNIST are illustrated in Table 2. By pre-setting a targeted accuracy for the test set, we record the number of communication rounds. We generate a testing-accuracy curve for each combination of parameters, enabling the curve almost monotonically updating through the best testing result produced over all prior rounds. The number of communication rounds is calculated where the curve crosses the pre-set accuracy. This can be better perceived through the reference of Fig. 4 in which the targeted accuracies are marked by grey lines. $B=\infty$ represents that all samples (full local batch size i.e., 600, 500 and 225 samples per client for MNIST, CIFAR-10 and MIMIC datasets respectively) for a client is used as a mini-batch for training in each round.

It is noticed that there is still a speedup by the proposed MIFL over the baseline FA when $B=\infty$ [7]. MIFL endows a noticeable speedup for a small batch size ($B=10$), especially in the non-IID case when $C \geq 0.1$. Based on the experimental observations, we fix $C = 0.1$ in the rest of experiments, which usually obtains a tradeoff between the convergence rate and the computational efficiency. For the IID MNIST, our optimization reduces the amount of communication rounds and achieves speedups of (1.20-1.70 \times) and (1.05-1.82 \times) using 2NN with $E = 1$ and $E = 5$, respectively, and (2.00-3.34 \times) using CNN with $E = 5$. Similarly, the speedups are (1.00-1.74 \times) and (1.00-1.70 \times) using 2NN with $E = 1$ and $E = 5$, respectively, and (1.33-3.15 \times) using CNN with $E = 5$ for the non-IID MNIST. In general, it can be observed that our approach can enable noticeable speedups for communication. This is because our both weight update and aggregation reduce the computation.

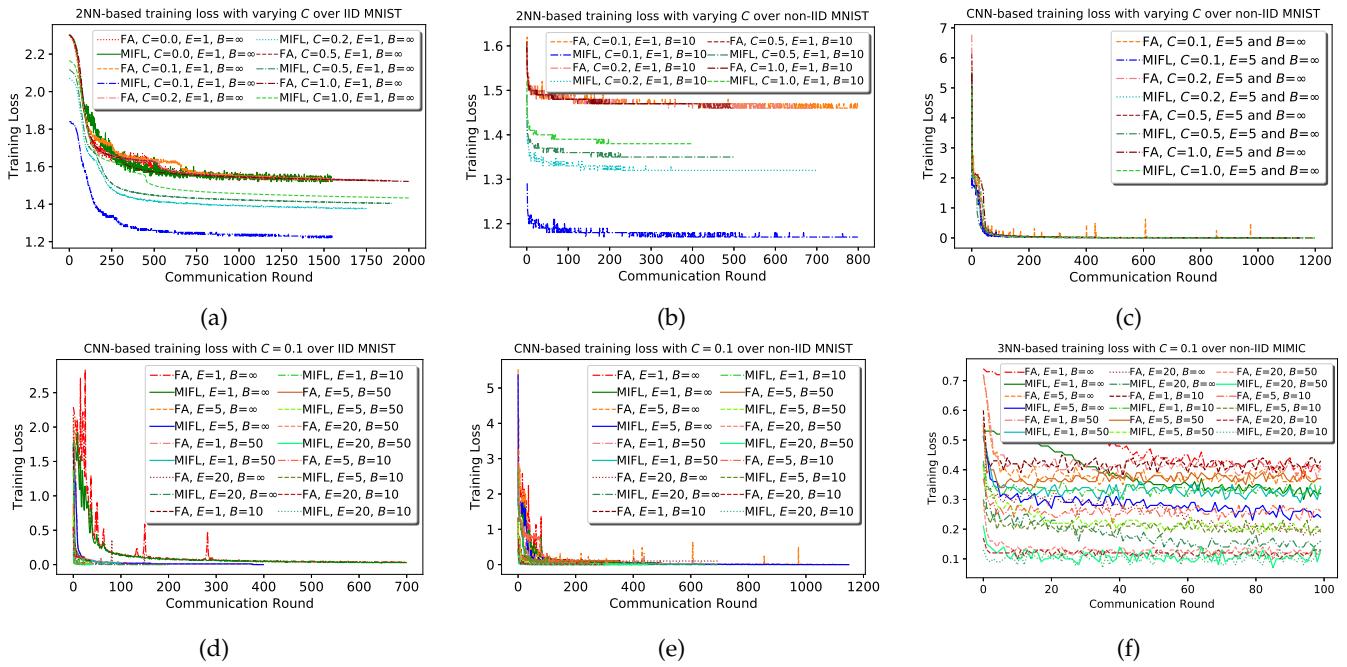


Fig. 5: Training losses versus communication rounds. (a) Varying C , $E = 1$ and $B = \infty$ using 2NN over IID MNIST. (b) Varying C , $E = 1$ and $B = 10$ using 2NN over non-IID MNIST. (c) Varying C , $E = 5$ and $B = \infty$ using CNN over non-IID MNIST. (d) $C = 0.1$, and varying E and B using CNN over IID MNIST. (e) $C = 0.1$, and varying E and B using CNN over non-IID MNIST. (f) $C = 0.1$, and varying E and B using 3NN over non-IID MIMIC.

5.3 Increasing Computation in Each Client

We now increase computation in each client by fixing $C=0.1$ in every communication round. This is accomplished by varying either B or E , or both. Some core results are illustrated in Fig. 4 (Some other results are provided in supplementary document C.1.). Table 3 shows the results of communication rounds. We calculate u , the expected number of updates for each client in each round, as $u = (\mathbb{E}[N_k]/B)E = NE/(KB)$, and arrange the table's rows with this value. It can be seen that varying B and E for increasing u is effective, as the communication rounds are generally reduced while reaching convergence (i.e., reaching the pre-set accuracy). For CNN MNIST, the speedup ranges from $(1.51\text{-}3.09)\times$ for the IID and $(1.18\text{-}3.15)\times$ for the non-IID setup, while 3NN MIMIC-III for the non-IID case produces a speedup of $(1.08\text{-}2.02)\times$. These speedups confirm that our MIFL has a great potential to be deployed in real-world applications. Although we primarily concern about the generalization performance of the proposed MIFL with the targeted test-set accuracy, it also shows the effectiveness at largely decreasing training loss, illustrated in Fig. 5 (Some other plots are provided in supplementary document C.2.).

5.4 Visualization of Training and Testing Progresses

We also calculate and visualize the MI, $MI(Y; \hat{Y})$, between the ground truth information (Y) and the logits (outputs, \hat{Y}) produced by every global model at each communication round using both training and testing sets. Some results are depicted in Fig. 6 (Some other results are provided in Appendix C.3.), which illustrates that the proposed MIFL with presented local models' update and aggregation strategies produces better MI values than those by the vanilla FA.

The training MI values in the non-IID setup, are fluctuating, because the data is more challenging to be converged. The way of calculating MI could be an potential attempt towards interpreting the black-box of DNN in the FL setup.

5.5 MI-based Loss Calculation for Model Training and Model Pruning for Model Aggregation

In the proposed federated optimization procedure, client's current local model is considered along with current aggregated global model to calculate the final loss to be adjusted by the global model on client's own private data. Doing this for all selected clients, the overall training loss can be significantly reduced by MIFL. This impact is shown in Fig. 5 and in supplementary document (C.2). It can be observed that the combination of losses in MIFL is converged faster than the loss of the standalone global model in vanilla FA. Besides this, the proposed federated optimization introduces the pruning of primitive and over-working outlier models before aggregation in each round. We pick 2.5% of the selected clients' models per round as primitive and over-working separately, based on the MI values between the updated local model and the previous aggregated model. The primitive models are identified by the highest MI values and the over-working models by the lowest MI values. The proposed MIFL is remarkably improved by this model pruning strategy, as illustrated in Fig. 7. From the curves of both IID and non-IID cases, it can be seen that the MI values of clients' local models first increase or decrease from low or high initial values, respectively, and then attempt to convergence through tiny fluctuations.

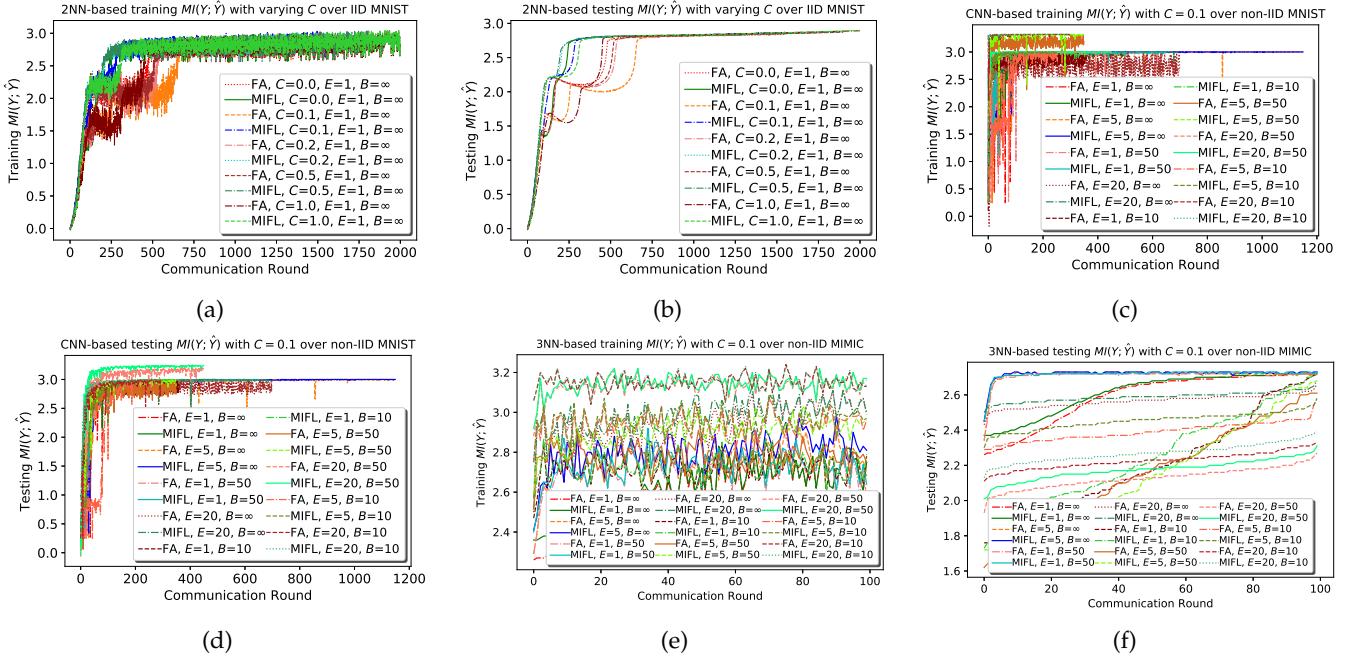


Fig. 6: Training and testing $MI(Y; \hat{Y})$ with increasing communication rounds. (a) and (b): varying C , $E = 1$ and $B = \infty$ using 2NN over IID MNIST (training and testing). (c) and (d): $C = 0.1$, and varying E and B using CNN over non-IID MNIST (training and testing). (e) and (f): $C = 0.1$, and varying E and B using 3NN over non-IID MIMIC (training and testing).

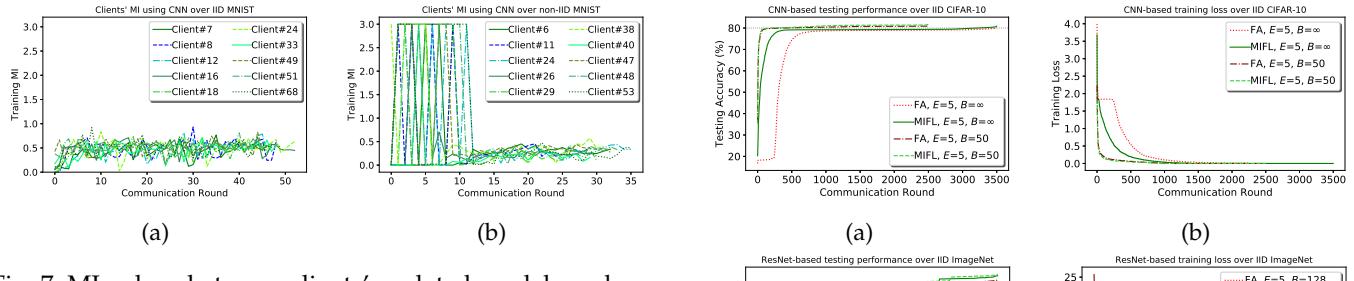


Fig. 7: MI values between clients' updated models and corresponding previous global model using CNN over MNIST training rounds of the clients ($C = 0.1$ and $E=5$). (a) IID case with $B=600$. (b) Non-IID case with $B=10$.

5.6 Experiments on CIFAR-10 and ImageNet

Based on the experiments and observations on MNIST, we also consider complex image classification task on the CIFAR-10 and ImageNet datasets. Unlike FA using testing accuracy for comparisons, we set the targeted test-set accuracies to 70-80% for CIFAR-10 and 60-70% (Top-1) for ImageNet, since we validate our proposed optimization algorithm using the standard CNN and ResNet models, respectively. Table 4 shows the required numbers of communication rounds to achieve the targeted test-set accuracies, along with corresponding speedups (1.09 – 3.32 \times) for CIFAR-10 and (1.14 – 1.66 \times) for ImageNet. Fig. 8 shows the testing accuracies and training losses with the increase of communication rounds. (MI-based learning visualization plots are provided in supplementary document C.4).

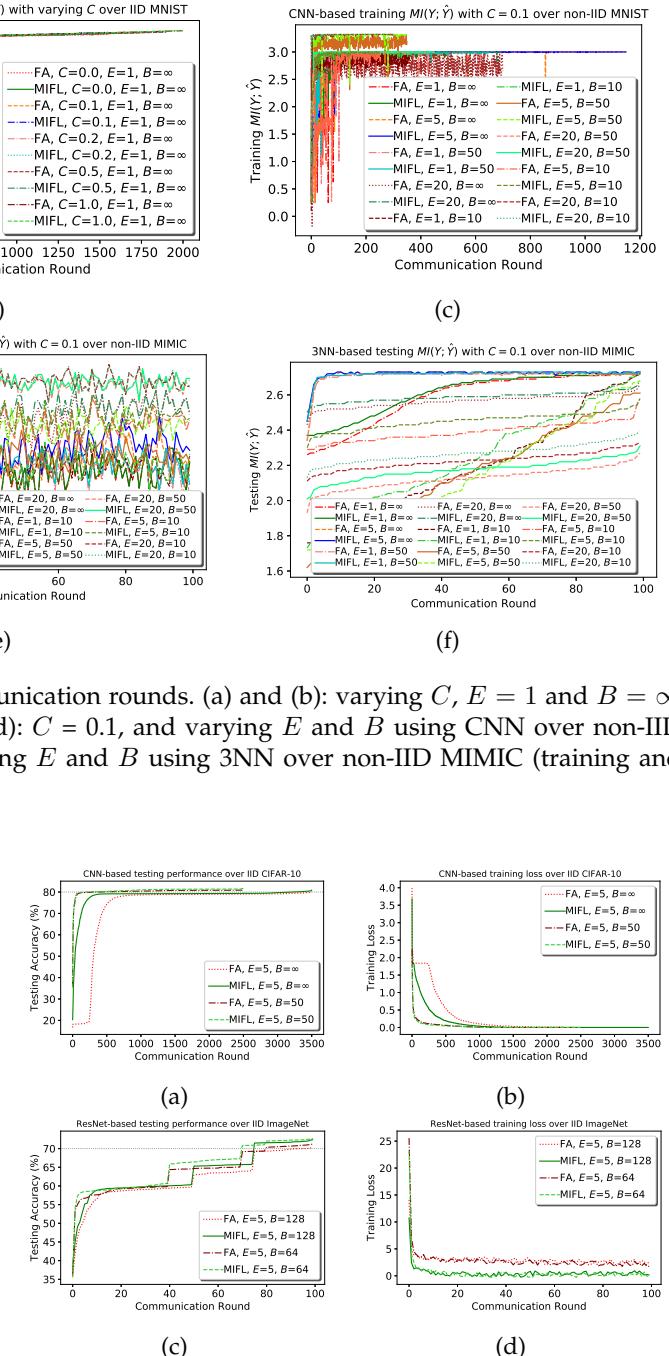


Fig. 8: CNN and ResNet-50 over CIFAR-10 and ImageNet, respectively. (a) and (b): testing accuracies and training losses for CIFAR-10. (c) and (d): testing accuracies and training losses for ImageNet.

5.7 Extended Comparison

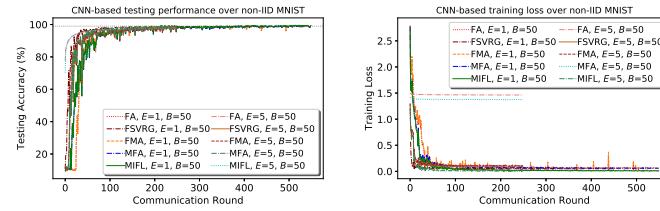
In addition to the above results, we also compare the FSVRG [8], MFA [10], and FMA [11] with our MIFL. Similar to above experiments, we set a targeted test-set accuracy, and run the aforementioned FL methods and MIFL on both IID and non-IID MNIST with the same CNN architecture and on non-IID MIMIC-III with same MLP architecture. We set $C=0.1$ for all methods except for FSVRG, as it requires all clients for training and full gradient calculation at each

TABLE 4: Numbers of communication rounds required to reach the targeted test-set accuracies with CNN over CIFAR-10 and ResNet-50 over ImageNet using $C = 0.1$ and $E=5$.

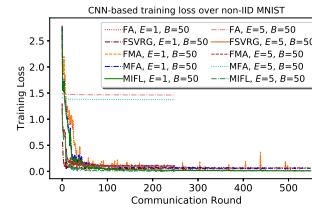
CNN CIFAR-10						ResNet-50 ImageNet					
B	Method	Targeted test-set accuracy			B	Method	Targeted test-set accuracy (Top-1)			B	Method
		70%	75%	80%			60%	65%	70%		
∞	FA	432	517	3497	128	FA	51	76	96	64	MIFL
	MIFL	130 (3.32 \times)	187 (2.76 \times)	3215 (1.09 \times)		MIFL	38 (1.34 \times)	51 (1.49 \times)	76 (1.26 \times)		FA
50	FA	17	32	337	64	FA	41	68	81		MIFL
	MIFL	13 (1.31 \times)	27 (1.19 \times)	207 (1.63 \times)		MIFL	32 (1.28 \times)	41 (1.66 \times)	71 (1.14 \times)		FA

TABLE 5: Numbers of communication rounds required to reach a targeted test-set accuracy for our MIFL and the compared FL techniques.

CNN MNIST, targeted test-set accuracy: 99%										3NN MIMIC-III, targeted test-set accuracy: 85%					
E, B	FA		FSVRG		MFA		FMA		MIFL		FA	FSVRG	MFA	FMA	MIFL
	IID	Non-IID	IID	Non-IID	IID	Non-IID	IID	Non-IID	IID	Non-IID	Non-IID	Non-IID	Non-IID	Non-IID	Non-IID
1, 50	57	530	49	392	50	406	51	415	38	311	93	89	89	90	69
5, 50	28	214	26	199	23	193	24	194	17	181	79	70	73	74	66



(a)



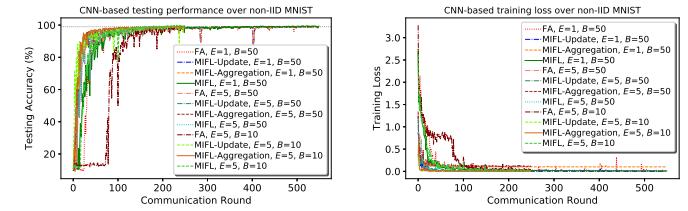
(b)

Fig. 9: CNN over non-IID MNIST for the comparison of the proposed MIFL with other FL methods. (a) Testing accuracies. (b) Training losses.

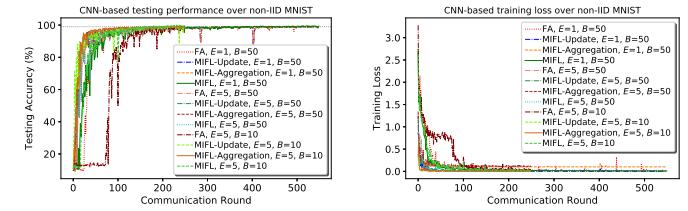
round, and set $E=\{1, 5\}$ and $B=50$. We provide the results in Table 5, and in Fig. 9 (other results are provided in supplementary document C.5.). It manifests that the MIFL still requires relatively less communications to achieve the targeted accuracy than the compared approaches.

5.8 Ablation Studies

We perform algorithm-level and data distribution-level ablation studies to investigate the influence of each component of our proposed MIFL. First, we remove the second component (MIFL-Aggregation) of MIFL and run the first component (MIFL-Update) on both IID and non-IID data independently, and run the aggregation scheme of vanilla FL. Then, we remove the first component (MIFL-Update) of MIFL, use the model update scheme of vanilla FL and run the second component (MIFL-Aggregation) on both IID and non-IID data individually. Additionally, we provide the results produced by our entire method (MIFL) and the vanilla FL. The results are provided in Table 6. Also, some accuracy and training loss plots are provided in Fig. 10, and other results plots are given in supplementary document C.6. From the results, it can be clearly seen that MIFL-Update and MIFL-Aggregation can both outperform the vanilla FL in reducing communication rounds, and together they boost the communication reduction.



(a)



(b)

Fig. 10: CNN over non-IID MNIST for the ablation study. (a) Testing accuracies. (b) Training losses.

6 CONCLUSION

We have introduced a novel federated optimization approach in this work. We first present an MI-driven client-side loss calculation for model training with considering both client's local model and the aggregated global model at each communication round. The final loss from these two models leads to faster training convergence. Secondly, we propose a server-side outlier model pruning strategy for model aggregation, based on the MI values between clients' trained models and the previous aggregated model. Top effective local models are thus selected for better aggregation and faster learning convergence. We also present the convergence proof for our proposed optimization algorithm MIFL. We have conducted extensive experiments to validate our method, and demonstrated that it outperforms state-of-the-art techniques, in terms of communication rounds for reaching a targeted test-set accuracy.

It is thus obvious from the observations that the adoption of MI in context of client-side model training and server-side outlier model pruning could be an interesting potential research dimension. Moreover, the presented MI-based shared global model's behaviour analysis mechanism signifies that the information bottleneck principle [31] might be utilized to unfold the 'blackbox' of shared DNN model learning. In addition, as in baseline FL, there is still a direction to employ the use of additional privacy preserving schemes, such as differential privacy [32], [33], encryption strategy [34] and secure multi-party computation [35].

TABLE 6: Impact of each component (MIFL-Update and MIFL-Aggregation) of MIFL with $E=\{1, 5\}$ and $B=\{10, 50\}$. (a) 2NN MNIST and 3NN MIMIC-III. (b) CNN MNIST and CNN CIFAR.

2NN MNIST, targeted test-set accuracy: 97%										3NN MIMIC-III, targeted test-set accuracy: 85%			
E, B	FA		MIFL-Update		MIFL-Aggregation		MIFL		FA	MIFL-Update	MIFL-Aggregation	MIFL	
	IID	Non-IID	IID	Non-IID	IID	Non-IID	IID	Non-IID					
1, 50	89	213	73 (1.22 \times)	178 (1.20 \times)	80 (1.11 \times)	188 (1.13 \times)	60 (1.48 \times)	157 (1.36 \times)	93	75 (1.24 \times)	83 (1.12 \times)	69 (1.35 \times)	
5, 50	68	173	54 (1.26 \times)	151 (1.15 \times)	59 (1.15 \times)	162 (1.07 \times)	46 (1.48 \times)	133 (1.30 \times)	79	72 (1.10 \times)	73 (1.08 \times)	66 (1.20 \times)	
5, 10	47	134	39 (1.21 \times)	107 (1.25 \times)	41 (1.15 \times)	118 (1.14 \times)	33 (1.42 \times)	90 (1.49 \times)	92	89 (1.03 \times)	90 (1.02 \times)	85 (1.08 \times)	

(a)

CNN MNIST, targeted test-set accuracy: 99%								CNN CIFAR, targeted test-set accuracy: 80%				
E, B	FA		MIFL-Update		MIFL-Aggregation		MIFL		FA	MIFL-Update	MIFL-Aggregation	MIFL
	IID	Non-IID	IID	Non-IID	IID	Non-IID	IID	Non-IID				
1, 50	57	530	47 (1.21 \times)	346 (1.53 \times)	48 (1.19 \times)	458 (1.16 \times)	38 (1.50 \times)	311 (1.70 \times)	546	350 (1.56 \times)	386 (1.41 \times)	278 (1.96 \times)
5, 50	28	214	23 (1.22 \times)	193 (1.11 \times)	25 (1.12 \times)	195 (1.10 \times)	17 (1.65 \times)	181 (1.18 \times)	337	248 (1.36 \times)	281 (1.20 \times)	207 (1.63 \times)
5, 10	18	239	12 (1.50 \times)	188 (1.27 \times)	15 (1.20 \times)	213 (1.12 \times)	7 (2.57 \times)	157 (1.52 \times)	292	223 (1.31 \times)	254 (1.15 \times)	179 (1.63 \times)

(b)

REFERENCES

- [1] K. L. Lueth, "State of the IoT 2018: Number of IoT Devices Now at 7B – Market Accelerating," *IOT Analytics*, 2018.
- [2] J. Park, S. Samarakoon, M. Bennis, and M. Debbah, "Wireless Network Intelligence at the Edge," *Proceedings of the IEEE*, vol. 107, pp. 2204–2239, Nov. 2019.
- [3] L. Huang, Y. Yin, Z. Fu, S. Zhang, H. Deng, and D. Liu, "LoAdaBoost: Loss-based AdaBoost federated machine learning with reduced computational complexity on IID and non-IID intensive care data," *PLOS ONE*, vol. 15, no. 4, pp. 1–16, 2020.
- [4] J. Dean, G. S. Corrado, R. Monga, K. Chen, M. Devin, Q. V. Le, M. Z. Mao, M. Ranzato, A. Senior, P. Tucker, K. Yang, and A. Y. Ng, "Large Scale Distributed Deep Networks," in *International Conference on Neural Information Processing Systems (NIPS 2012)*, (Lake Tahoe, Nevada), 2012.
- [5] D. Povey, X. Zhang, and S. Khudanpur, "Parallel training of Deep Neural Networks with Natural Gradient and Parameter Averaging," in *International Conference on Artificial Intelligence and Statistics (AISTATS)*, (Reykjavik, Iceland), 2014.
- [6] S. Zhang, A. Choromanska, and Y. LeCun, "Deep learning with Elastic Averaging SGD," in *International Conference on Neural Information Processing Systems (NIPS 2015)*, (Montreal, Canada), 2015.
- [7] H. B. McMahan, E. Moore, D. Ramage, S. Hampson, and B. A. y. Arcas, "Communication-Efficient Learning of Deep Networks from Decentralized Data," in *20th International Conference on Artificial Intelligence and Statistics (AISTATS)*, (Florida, USA), 2017.
- [8] J. Konečný, H. B. McMahan, and P. Richtarik, "Federated Optimization: Distributed Machine Learning for On-Device Intelligence," *Computing Research Repository (CoRR)*, vol. Computing Research Repository (CoRR), p. 38, 2016.
- [9] R. Johnson and T. Zhang, "Accelerating Stochastic Gradient Descent using Predictive Variance Reduction," in *International Conference on Neural Information Processing Systems (NIPS 2013)*, (Lake Tahoe, Nevada), 2013.
- [10] W. Liu, L. Chen, Y. Chen, and W. Zhang, "Accelerating Federated Learning via Momentum Gradient Descent," *IEEE Transactions on Parallel and Distributed Systems*, vol. 31, no. 8, pp. 1754–1766, 2020.
- [11] H. Wang, M. Yurochkin, Y. Sun, D. Papailiopoulos, and Y. Khazaeni, "Federated Learning with Matched Averaging," in *International Conference on Learning Representations*, (Addis Ababa, Ethiopia), 2020.
- [12] M. Yurochkin, M. Agarwal, S. Ghosh, K. Greenewald, T. N. Hoang, and Y. Khazaeni, "Bayesian Nonparametric Federated Learning of Neural Networks," in *36th International Conference on Machine Learning*, (California, USA), 2019.
- [13] H. Zhu and Y. Jin, "Multi-Objective Evolutionary Federated Learning," *IEEE Transactions on Neural Networks and Learning Systems*, vol. 31, no. 4, pp. 1310–1322, 2020.
- [14] T. Li, A. K. Sahu, M. Zaheer, M. Sanjabi, A. Talwalkar, and V. Smith, "Federated Optimization in Heterogeneous Networks," in *MLSys Conference*, (TX, USA), 2020.
- [15] M. Duan, D. Liu, X. Chen, R. Liu, Y. Tan, and L. Liang, "Self-Balancing Federated Learning With Global Imbalanced Data in Mobile Systems," *IEEE Transactions on Parallel and Distributed Systems*, vol. 32, no. 1, pp. 59–71, 2020.
- [16] M. Li, T. Zhang, Y. Chen, and A. J. Smola, "Efficient mini-batch training for stochastic optimization," in *20th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining - KDD '14*, (New York, USA), pp. 661–670, 2014.
- [17] Yong Liu and Xin Yao, "Simultaneous training of negatively correlated neural networks in an ensemble," *IEEE Transactions on Systems, Man and Cybernetics, Part B (Cybernetics)*, vol. 29, no. 6, pp. 716–725, 1999.
- [18] Y. Liu and X. Yao, "Ensemble learning via negative correlation," *Neural Networks*, vol. 12, pp. 1399–1404, 1999.
- [19] X. Yao and Y. Liu, "Evolving Neural Network Ensembles by Minimization of Mutual Information," *International Journal of Hybrid Intelligent Systems*, vol. 1, pp. 12–21, Sept. 2004.
- [20] R. Battiti, "Using mutual information for selecting features in supervised neural net learning," *IEEE Transactions on Neural Networks*, vol. 5, pp. 537–550, July 1994.
- [21] Hanchuan Peng, Fuhui Long, and C. Ding, "Feature selection based on mutual information criteria of max-dependency, max-relevance, and min-redundancy," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 27, pp. 1226–1238, Aug. 2005.
- [22] P. Estevez, M. Tesmer, C. Perez, and J. Zurada, "Normalized Mutual Information Feature Selection," *IEEE Transactions on Neural Networks*, vol. 20, pp. 189–201, Feb. 2009.
- [23] S. Cang and H. Yu, "Mutual information based input feature selection for classification problems," *Decision Support Systems*, vol. 54, pp. 691–698, Dec. 2012.
- [24] R. Zhang and J. T. Kwok, "Asynchronous Distributed ADMM for Consensus Optimization," in *International Conference on Machine Learning*, (Beijing, China), 2014.
- [25] P. Jain, P. Netrapalli, S. M. Kakade, R. Kidambi, and A. Sidford, "Parallelizing Stochastic Gradient Descent for Least Squares Regression: Mini-batching, Averaging, and Model Misspecification," *Journal of Machine Learning Research*, vol. 18, no. 223, pp. 1–42, 2018.
- [26] P. Jain, S. M. Kakade, R. Kidambi, P. Netrapalli, V. K. Pillutla, and A. Sidford, "A Markov Chain Theory Approach to Characterizing the Minimax Optimality of Stochastic Gradient Descent (for Least

- Squares)," in *IARCS Annual Conference on Foundations of Software Technology and Theoretical Computer Science*, (Kanpur, India), 2017.
- [27] H. Hajibabadi, D. Molla-Aliod, R. Monsefi, and H. S. Yazdi, "Combination of loss functions for deep text classification," *International Journal of Machine Learning and Cybernetics*, 2019.
- [28] A. E. Johnson, T. J. Pollard, L. Shen, L.-w. H. Lehman, M. Feng, M. Ghassemi, B. Moody, P. Szolovits, L. Anthony Celi, and R. G. Mark, "MIMIC-III, a freely accessible critical care database," *Scientific Data*, vol. 3, 2016.
- [29] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and L. Fei-Fei, "ImageNet: A Large-Scale Hierarchical Image Database," in *CVPR09*, 2009.
- [30] K. He, X. Zhang, S. Ren, and J. Sun, "Deep Residual Learning for Image Recognition," in *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, (Las Vegas, NV, USA), 2016.
- [31] N. Tishby and N. Zaslavsky, "Deep learning and the information bottleneck principle," in *IEEE Information Theory Workshop (ITW)*, (Jerusalem, Israel), pp. 1–5, 2015.
- [32] N. Agarwal, A. T. Suresh, F. Yu, S. Kumar, and H. B. McMahan, "cpSGD: Communication-efficient and differentially-private distributed SGD," in *International Conference on Neural Information Processing Systems (NIPS 2018)*, (Montreal, Canada), 2018.
- [33] Y. Lu, X. Huang, Y. Dai, S. Maharjan, and Y. Zhang, "Differentially Private Asynchronous Federated Learning for Mobile Edge Computing in Urban Informatics," *IEEE Transactions on Industrial Informatics*, vol. 16, no. 3, pp. 2134–2143, 2020.
- [34] L. Lyu, J. Yu, K. Nandakumar, Y. Li, X. Ma, J. Jin, H. Yu, and K. S. Ng, "Towards Fair and Privacy-Preserving Federated Deep Models," *IEEE Transactions on Parallel and Distributed Systems*, vol. 31, no. 11, pp. 2524–2541, 2020.
- [35] K. Bonawitz, V. Ivanov, B. Kreuter, A. Marcedone, H. B. McMahan, S. Patel, D. Ramage, A. Segal, and K. Seth, "Practical Secure Aggregation for Privacy-Preserving Machine Learning," in *Conference on Computer and Communications Security*, (Texas, USA), 2017.



Md Palash Uddin received the B. Sc. degree in Computer Science and Engineering from Hajee Mohammad Danesh Science and Technology University (HSTU), Bangladesh and the M. Sc. degree in Computer Science and Engineering from Rajshahi University of Engineering & Technology, Bangladesh. He is currently pursuing the Ph.D. degree with the School of Information Technology, Deakin University, Australia. He is also an academic faculty member in HSTU, Bangladesh. His research interests include machine learning, blockchain and remote sensing image analysis.



Yong Xiang (SM'12) received the Ph.D. degree in Electrical and Electronic Engineering from The University of Melbourne, Australia. He is a Professor at the School of Information Technology, Deakin University, Australia. His research interests include information security and privacy, signal and image processing, data analytics and machine intelligence, Internet of Things, and blockchain. He has published 5 monographs, over 150 refereed journal articles, and numerous conference papers in these areas.

Professor Xiang is the Senior Area Editor of IEEE Signal Processing Letters, Associate Editor of IEEE Communications Surveys and Tutorials and IEEE Access, and Guest Editor of IEEE Transactions on Industrial Informatics. He was the Associate Editor of IEEE Signal Processing Letters and the Guest Editor of IEEE Multimedia. He has served as Honorary Chair, General Chair, Program Chair, TPC Chair, Symposium Chair and Track Chair for many conferences, and was invited to give keynotes at a number of international conferences.



Xuequan Lu is a Lecturer (Assistant Professor) at Deakin University, Australia. He spent more than two years as a Research Fellow in Singapore. Prior to that, he earned his Ph.D at Zhejiang University (China) in June 2016. His research interests mainly fall into the category of visual computing, for example, geometry modeling, processing and analysis, animation/simulation, 2D data processing and analysis. More information can be found at <http://www.xuequanlu.com>.



Professor John Yearwood is Head of the School of Information Technology at Deakin University and was previously Executive Dean, Faculty of Science and Technology, and Director, Centre of Informatics and Applied Optimization at Federation University, Australia. He was instrumental in setting up the Internet Commerce Security Laboratory with Westpac, IBM and the Victorian State Government as a joint industry-focused and data-driven laboratory on cybersecurity in the financial sector. He has held a number of ARC grants and was a QEII Fellow working on computational narrative and argumentation in decision science. Professor Yearwood's work in data mining and computational intelligence has led to the development of new machine learning and hybrid learning algorithms for artificial neural networks, as well as new data and text mining and pattern recognition approaches. His work in decision science has developed the use of argumentation structures for the modelling of knowledge and collaborative decision making in complex domains. He has published over 200 journal and refereed conference papers including 2 books. Professor Yearwood is currently a CI on the ARC funded Discovery Project 'Enhancing and supporting deliberation in multi-disciplinary team decision-making'. He is Editor-in-Chief of the Journal of Research Practice in Information Technology and a reviewer for a large number of journals and competitive research grant programs including the Australian Research Council grant program, the NHMRC grant program, and for the Dutch Government in the assessment of their NWO/ToKeN2000.



Longxiang Gao (SM'17) received his PhD in Computer Science from Deakin University, Australia. He is currently a Senior Lecturer and Co-founder of Deakin Blockchain Innovation Lab at School of Information Technology, Deakin University. Before joined Deakin University, he was a post-doctoral research fellow at IBM Research Development, Australia. His research interests include fog/edge computing, blockchain, data analysis and privacy protection. Dr. Gao has over 80 publications, including patent, monograph, book chapter, journal and conference papers. Some of his publications have been published in the top venue, such as IEEE TMC, IEEE TPDS, IEEE IoTJ, IEEE TDSC, IEEE TVT, IEEE TCSS, IEEE TII and IEEE TNSE. He has being Chief Investigator (CI) for more than 20 research projects (the total awarded amount is over \$5million), from pure research project to contracted industry research. Dr. Gao is a Senior Member of IEEE and active in IEEE Communication Society. He has served as the TPC co-chair, publicity co-chair, organization chair and TPC member for many international conferences. He is the Associate Editor of IEEE Access and the Assessor of Australian Research Council (ARC) projects.