# AUTO-WEIGHTED ROBUST FEDERATED LEARNING WITH CORRUPTED DATA SOURCES

### A PREPRINT

**Shenghui Li**
Uppsala University, Sweden
shenghui.li@it.uu.se

**Edith Ngai**
The University of Hong Kong, China
chngai@eee.hku.hk

**Fanhua Ye**
University College London, UK
fanghua.ye.19@ucl.ac.uk

**Thiemo Voige**
Uppsala University, Sweden
thiemo.voigt@it.uu.se

July 15, 2022

## ABSTRACT

Federated learning provides a communication-efficient and privacy-preserving training process by enabling learning statistical models with massive participants without accessing their local data. Standard federated learning techniques that naively minimize an average loss function are vulnerable to data corruptions from outliers, systematic mislabeling, or even adversaries. In this paper, we address this challenge by proposing Auto-weighted Robust Federated Learning (ARFL), a novel approach that jointly learns the global model and the weights of local updates to provide robustness against corrupted data sources. We prove a learning bound on the expected loss with respect to the predictor and the weights of clients, which guides the definition of the objective for robust federated learning. We present an objective that minimizes the weighted sum of empirical risk of clients with a regularization term, where the weights can be allocated by comparing the empirical risk of each client with the average empirical risk of the best $p$ clients. This method can downweight the clients with significantly higher losses, thereby lowering their contributions to the global model. We show that this approach achieves robustness when the data of corrupted clients is distributed differently from the benign ones. To optimize the objective function, we propose a communication-efficient algorithm based on the blockwise minimization paradigm. We conduct extensive experiments on multiple benchmark datasets, including CIFAR-10, FEMNIST, and Shakespeare, considering different neural network models. The results show that our solution is robust against different scenarios including label shuffling, label flipping, and noisy features, and outperforms the state-of-the-art methods in most scenarios.

## 1 Introduction

Federated learning (McMahan et al., 2017; Li et al., 2021; Konečný et al., 2016a) has recently attracted more and more attention due to the increasing concern of user data privacy. In federated learning, the server trains a shared model based on data originating from remote clients such as smartphones and IoT devices, without the need to store and process the data in a centralized server. In this way, federated learning enables joint model training over privacy-sensitive data in a wide range of applications (Yang et al., 2018, 2019b), including natural language processing (Chen et al., 2019), computer vision (Luo et al., 2019), and speech recognition (Guliani et al., 2020).

However, standard federated learning strategies fail in the presence of data corruption (Rodríguez-Barroso et al., 2020; Mahloujifar et al., 2019). Data collected from different clients, or data sources (in this article we use the two terms interchangeably), may vary greatly in data quality and thus reduce the reliability of the learning task. Some of the clients may be unreliable or even malicious. For instance, distributed sensor networks are vulnerable to cybersecurity attacks, such as false data injection attacks (Liu and Ngai, 2019). In crowdsourcing scenarios, the problem of noisy data

is unignorable due to biased or erroneous workers (Wais et al., 2010). As a consequence, the clients can update their local models using corrupted data and send the parameters to the server. After averaging these harmful updates, the accuracy and the convergence of the shared global model can be compromised. To mitigate this limitation, different robust learning approaches have been proposed in the literature. Some of these techniques rely on robust statistical estimations (e.g. geometric median estimators) to update the aggregated model (Blanchard et al., 2017; Pillutla et al., 2019; Yin et al., 2018). Although these techniques are widely used in traditional distributed learning scenarios with i.i.d. datasets, it is not straightforward to generalize them for non-i.i.d. data settings, i.e., some of the clients have significantly different data distribution than the others. Other approaches require some trusted clients or samples to guide the learning (Li et al., 2020b; Peng et al., 2019; Sui et al., 2020) or detect the updates from corrupted clients (Han and Zhang, 2020; Konstantinov and Lampert, 2019). Unfortunately, the credibility of these trusted clients and samples are usually not guaranteed. Since the data is stored locally in the clients, the server is insensible to the corruption behaviors and unable to measure the quality of data sources due to privacy constraints.

In this article, we aim to improve the robustness of federated learning when some of the clients are malicious, i.e., they actively corrupt the local training set by changing the features or the labels. We propose a novel solution, named Auto-weighted Robust Federated Learning (`ARFL`), to jointly learn the global model and the weights of local updates. More specifically, we first prove a learning bound on the expected risk with respect to the predictor and the weights of clients. Based on this theoretical insight, we present our objective that minimizes a weighted sum of the empirical risk of clients with a regularization term. We then theoretically show that the weights in the objective can be allocated by comparing the empirical risk of each client with the best $p$ clients. When the corrupted clients have significantly higher losses comparing with the benign ones, their contributions to the global model will be downweighted or even zero-weighted, so as to play less important roles in the global model. Therefore, by using `ARFL`, we can exclude potentially corrupted clients and keep optimizing the global model with the benign clients. To solve the problem in federated learning settings, we further propose a communication-efficient optimization method based on the blockwise updating paradigm (Zheng et al., 2019). Through extensive experiments on multiple benchmark datasets (i.e., CIFAR-10, FEMNIST and Shakespeare) with different neural network models, we demonstrate the robustness of our approach compared with the state-of-the-art approaches (McMahan et al., 2017; Pillutla et al., 2019; Blanchard et al., 2017; Sattler et al., 2020), showing up to $30\%$ improvement in model accuracy.

## 2    Related work

The concept of federated learning has been proposed for collaboratively learning a model without collecting users' data (McMahan et al., 2017; Li et al., 2021; Konečnỳ et al., 2016a,b). The research work on federated learning can be divided into three categories, i.e., horizontal federated learning, vertical federated learning, and federated transfer learning, based on the distribution characteristics of the data. Due to space limits, we refer to Yang et al. (Yang et al., 2019a) for detailed explanations. In this paper we focus on horizontal federated learning where datasets in all clients share the same feature space but different samples. In 2017, one of the most famous horizontal federated learning framework, Federated Averaging (`FedAvg`) has been proposed to update global parameters with a weighted average of the model parameters sent by a subset of clients after several local iterations (McMahan et al., 2017).

Due to the nature of data decentralization and the requirement of collaborative training from multiple clients, federated learning is vulnerable to malicious corruption of training data from remote clients (Lyu et al., 2020). Notably, it has been shown that traditional federated learning approaches (e.g., `FedAvg`) are fragile in the presence of malicious or corrupted clients (Fang et al., 2020; So et al., 2020). To mitigate the impact of malicious clients and improve the robustness of federated training, several solutions have been proposed in the literature (So et al., 2020; Pillutla et al., 2019; Kairouz et al., 2021; Li et al., 2021). Among these robust approaches, robust statistical estimations have received much attention in particular. Typical estimation rules include Geometric Median (GM) (Chen et al., 2017; Pillutla et al., 2019), trimmed mean (Yin et al., 2018), and Krum (Blanchard et al., 2017). For instance, in 2017, Chen et al. (2017) proposed to apply GM as a gradient aggregation protocol in robust distributed learning and showed that it can tolerate up to half malicious clients while estimating the underlying true gradients. In 2019, Pillutla et al. (2019) used GM to aggregate parameters in a robust FL solution. However, the robustness of the traditional estimators are limited as they rely on the assumption that data are i.i.d. and balanced among the clients, i,e, they distribute identically and have the same (or a similar) number of training data points. Hence, those approaches can be inefficient when some of the clients have significantly more data than others. In addition, Li et al. (2021) proposed a multi-task learning objective called Ditto, for federated learning that provides robustness via personalization in 2021. The optimization of Ditto, however, requires extra training overhead for personalized models.

Others (Li et al., 2020a; Han and Zhang, 2020; Sattler et al., 2020; Cao et al., 2021), couple the process of teaching and learning based on a few trusted instances to produce a robust prediction model. For example, in 2020, Sattler et al. (2020) proposed to separate the client population into different groups (e.g., benign and corrupted groups) based on

the pairwise cosine similarities between their parameter updates. Also in 2020, Li et al. (2020a) suggested allowing the server to learn to detect and remove the malicious model updates using an encoder-decoder based detection model. These approaches require some trusted clients or samples to guide the learning or detect the updates from corrupted clients. In 2021, Cao et al. (2021) presented a Byzantine-robust method that allows the service provider itself collect a clean small training datase for the learning task and maintain a model based on it to bootstrap trust. Unfortunately, the credibility of these trusted clients and samples are usually not guaranteed since the data is isolated and stored locally. Thus, the server is not aware of these corruption behaviors and does not have the ability to measure the quality of data at the sources due to privacy and communication constraints. Different from previous studies, we propose a robust approach that can learn both the global model and the weights of clients automatically from a mix of reliable and unreliable clients, without the need of any pre-verified trusted instances.

## 3 Preliminaries and Motivation

### 3.1 Federated Learning

In federated learning tasks, a general assumption is that the target distribution for which the centralized model is learned is a weighted mixture of distributions associated with all clients (or data sources), that is, if we denote by $\mathcal{D}_i$ the distribution associated with the $i$-th client, the centralized model is trained to minimize the risk with respect to $\mathcal{D}_{\boldsymbol{\alpha}} = \sum_{i=1}^{N} \alpha_i \mathcal{D}_i$, where $N$ is the total number of clients, $\boldsymbol{\alpha} = (\alpha_1, ..., \alpha_N)^{\top}$ is the vector of source-specific weights. We also have $\boldsymbol{\alpha} \in \mathbb{R}_+^n$ and $\mathbf{1}^{\top} \boldsymbol{\alpha} = 1$ Mohri et al. (2019); Hamer et al. (2020); Li et al. (2021).

Let $\ell_h(\boldsymbol{z})$ be the loss function that captures the error of a predictor $h \in \mathcal{H}$ (where $\mathcal{H}$ is the hypothesis class) on the training data $\boldsymbol{z} = (\boldsymbol{x}, y)$ (where $(\boldsymbol{x}, y)$ is the input and output pair), and $\mathcal{L}_{\mathcal{D}_{\boldsymbol{\alpha}}}(h)$ be the risk of a predictor $h$ on the mixture data distribution $\mathcal{D}_{\boldsymbol{\alpha}}$, we have:

$$\mathcal{L}_{\mathcal{D}_{\boldsymbol{\alpha}}}(h) = \sum_{i=1}^{N} \alpha_i \mathcal{L}_i(h) = \sum_{i=1}^{N} \alpha_i \mathbb{E}_{\boldsymbol{z} \sim \mathcal{D}_i}(\ell_h(\boldsymbol{z})), \tag{1}$$

where $\mathcal{L}_i(h) = \mathbb{E}_{\boldsymbol{z} \sim \mathcal{D}_i}(\ell_h(\boldsymbol{z}))$ is the expected loss of a predictor $h$ on the data distribution $\mathcal{D}_i$ of the $i$-th client.

Most prior work in federated learning has assumed that all samples are uniformly weighted, where the underlying assumption is that the target distribution is $\overline{\mathcal{U}} = \sum_{i=1}^{N} \frac{m_i}{M} \mathcal{D}_i$, where $m_i$ is the number of samples from client $i$ and $M = \sum_{i=1}^{N} m_i$. Thus the risk becomes:

$$\mathcal{L}_{\overline{\mathcal{U}}}(h) = \sum_{i=1}^{N} \frac{m_i}{M} \mathbb{E}_{\boldsymbol{z} \sim \mathcal{D}_i}(\ell_h(\boldsymbol{z})). \tag{2}$$

In practice, the goal is to minimize a traditional empirical risk $\hat{\mathcal{L}}_{\overline{\mathcal{U}}}(h)$ as follows:

$$\hat{\mathcal{L}}_{\overline{\mathcal{U}}}(h) = \sum_{i=1}^{N} \frac{m_i}{M} \frac{1}{m_i} \sum_{j=1}^{m_i} \ell_h(\boldsymbol{z}_{i,j}), \tag{3}$$

which can be minimized by sampling a subset of clients randomly at each round, then running an optimizer such as stochastic gradient descent (SGD) for a variable number of iterations locally on each client. These local updating methods enable flexible and efficient communication compared to traditional mini-batch methods, which would simply calculate a subset of the gradients Wang and Joshi (2019); Stich (2018); Yu et al. (2019). One of the most well-known methods to minimize Eq. (3) in non-convex settings is FedAvg McMahan et al. (2017), which runs simply by letting each selected client apply a fixed number of epochs of SGD locally and then averaging the resulted local models.

### 3.2 Threat Model: Data Corruption

Guerraoui et al. has shown that a few clients with corrupted data can lead to inaccurate models Guerraoui et al. (2018). The problem stems from a mismatch between the target distribution and $\overline{\mathcal{U}}$. That is, in corruption scenarios, the target distribution may in general be quite different from $\overline{\mathcal{U}}$, since $\overline{\mathcal{U}}$ includes some corrupted components. We expect that the data distributions are more similar among benign clients compared with the corrupted ones even when the data is non-i.i.d. More specifically, we model the target distribution with corrupted clients as,

$$\mathcal{D}_{\boldsymbol{\alpha}} = \sum_{i=1}^{N} \alpha_i \mathcal{D}_i = \sum_{i=1}^{N} \alpha_i (\eta_i \mathcal{D}_{i,b} + (1 - \eta_i) \mathcal{D}_{i,c}), \tag{4}$$

3

(a) Clean                     (b) Feature Corruption             (c) Label Corruption
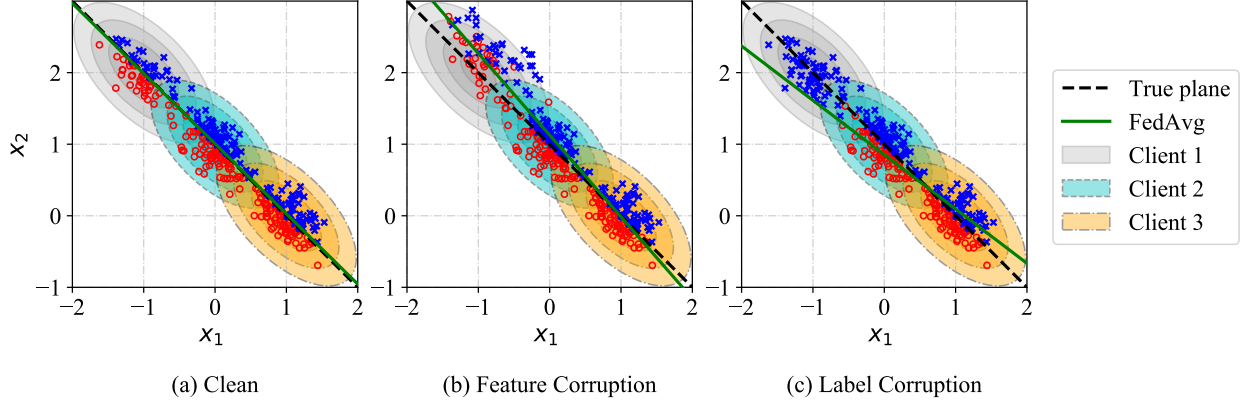
Figure 1: Illustration of federated learning using the standard `FedAvg` with potential corruptions in local datasets, where red circles represent data of class 0 and blue crosses represent class 1. (a) All samples are clean. (b) The features in Client 1 are shifted by random noise. (c) The labels in Client 1 are forced to class 1.

where $\eta_i \in \{0, 1\}$ denotes whether the local data distribution $\mathcal{D}_i$ is a benign distribution $\mathcal{D}_{i,b}$ (when $\eta_i = 1$) or a corrupted distribution $\mathcal{D}_{i,c}$ (when $\eta_i = 0$). Ideally, the corruption can be usually measured by the difference between $\mathcal{D}_{i,b}$ and $\mathcal{D}_{i,c}$ using statistical distance (e.g., the Kullback-Leibler divergence), if both distributions are known Konstantinov and Lampert (2019). In practice, data corruption can be achieved by modifying features or labels of the local training dataset.

When $\sum_{i=1}^{N} \eta_i = N$, all components of the mixture distribution are benign (i.e. $\mathcal{D}_i = \mathcal{D}_{i,b}$). Assuming that the target distribution is uniform $\overline{\mathcal{U}}$, minimizing Eq. (3) can lead to an accurate global model. However, when there are corrupted data sources (i.e. $\sum_{i=1}^{N} \eta_i < N$), the mixture of the distributions will include some corrupted components $\mathcal{D}_{i,c}$. In this case, optimizing the empirical risk with respect to $\overline{\mathcal{U}}$ will not lead to an accurate global model. We make the following assumptions regarding the adversaries:

1. Each adversary controls exactly one non-colluding and corrupted client. The data distributions of any corrupted clients are independent of each other. Since all clients do not collude, the effect of malicious updates from each adversary to the global model is limited Bhagoji et al. (2019).

2. A corrupted client has a higher loss with respect to the best predictor under the uniformly weighted assumption in Eq. (2), i.e., we have

$$\hat{\mathcal{L}}_{D_{i,c}}(h^*) > \hat{\mathcal{L}}_{D_{i,b}}(h^*), \tag{5}$$

where $h^*$ is the optimal global predictor that minimizes the empirical risk. Base on this, we will identify the corrupted clients according to their empirical losses as discussed in the next section.

### 3.3 Introductory Example

As an example, we train a binary classification model from three clients, where the data is originally generated from three distributions with linearly separable classes. A Logistic Regression (LR) model is learned using the standard `FedAvg` approach. As shown in Fig. 1(a), the learned separating plane is close to the true plane when the all datasets are clean. When Client 1 is corrupted on either feature $x$ (Fig. 1(b)) or label $y$ (Fig. 1(c)), the learned plane is driven away from the true one.

One possible solution is to exclude all the corrupted components of the mixture distribution and seek for a new target distribution considering only benign clients, i.e. $\overline{\mathcal{U}}_b = \sum_{i \in \mathcal{B}} \frac{m_i}{M_\mathcal{B}} \mathcal{D}_i$, where $\mathcal{B}$ is the set of all benign clients and $M_\mathcal{B} = \sum_{i \in \mathcal{B}} m_i$. The challenge with this approach is that the centralized server is agnostic to the corruptions on the local clients and hence it is impossible to measure the data quality directly.

## 4   Robust Federated Learning

In this section, we describe our Auto-weighted Robust Federated Learning (`ARFL`) in detail. Specifically, we present a novel objective according to a learning bound with respect to both the predictor $h$ and the weights $\boldsymbol{\alpha}$. We analyze the robustness of the proposed objective against corruptions and optimize it with a federated learning based algorithm.

### 4.1   Learning Bound

In the following theorem, we present a learning bound on the risk on the weighted mixture distribution with respect to the predictor $h$ and the weights $\boldsymbol{\alpha}$. A proof is provided in Appendix A.

**Theorem 1.** *We denote $\hat{\mathcal{L}}_i(h)$ as the corresponding empirical counterparts of $\mathcal{L}_{\mathcal{D}_i}(h)$. Assume that the loss function $\ell_h(\boldsymbol{z})$ is bounded by a constant $\mathcal{M} > 0$. Then, for any $\delta \in (0,1]$, with probability at least $1 - \delta$ over the data, the following inequality holds:*

$$\mathcal{L}_{\mathcal{D}_{\boldsymbol{\alpha}}}(h) \leq \sum_{i=1}^{N} \alpha_i \hat{\mathcal{L}}_i(h) + 2 \sum_{i=1}^{N} \alpha_i \mathcal{R}_i(\mathcal{H}) + 3 \sqrt{\frac{\log\left(\frac{4}{\delta}\right) \mathcal{M}^2}{2}} \sqrt{\sum_{i=1}^{n} \frac{\alpha_i^2}{m_i}},$$

*where, for each client $i = 1, \ldots, N$,*

$$\mathcal{R}_i(\mathcal{H}) = \mathbb{E}_{\sigma}\left(\sup_{f \in \mathcal{H}}\left(\frac{1}{m_i}\sum_{j=1}^{m_i} \sigma_{i,j} \ell_f(z_{i,j})\right)\right)$$

*and $\sigma_{i,j}$s are independent uniform random variables taking values in $\{-1, +1\}$. These variables are called Rademacher random variables (Yin et al., 2019).*

Although the Rademacher complexities in the bound are functions of both the underlying distribution and the hypothesis class (Yin et al., 2019), in practice one usually works with a computable upper bound of $\mathcal{R}_i(\mathcal{H})$ that is distribution-independent (e.g. using VC dimension) (Shalev-Shwartz and Ben-David, 2014; Konstantinov and Lampert, 2019; Bousquet et al., 2003). In our setting the hypothesis space $\mathcal{H}$ is fixed and hence these bounds would be identical for all $i$. Therefore, we expect the $\mathcal{R}_i(\mathcal{H})$ to be of similar order for all clients and the impact of $\boldsymbol{\alpha}$ in the second term to be negligible. We refer to Konstantinov et al. (Konstantinov and Lampert, 2019, Appendix B.1) for detailed explanations and some examples on this point. In general, we can ignore this term and concentrate on optimizing the remaining terms.

The rest of the terms for the learning bound in Theorem 1 thus suggest a trade-off between reducing the weighted sum of the empirical losses and minimizing a weighted norm of the weights $\sqrt{\sum_{i=1}^{n} \frac{\alpha_i^2}{m_i}}$. Note that reducing the weighted sum of the empirical losses will encourage trusting the clients who provide data with the smallest loss but it may lead to large and sparse weights, which will make the model fit only very few local datasets. On the contrary, minimizing the norm will increase the smoothness of the weight distribution, but it may also increase the weighted sum of empirical losses. To control this trade-off, we introduce a tuning hyperparameter $\lambda$ to the weighted norm. For the convenience of derivation, the square root on the sum of weights can be removed by tuning $\lambda$.

### 4.2   Problem Formulation

We present our problem formulation in the following. The learning bound derived above suggests minimizing the following objective with respect to the model parameters $\boldsymbol{w}$ together with the weights $\boldsymbol{\alpha}$:

$$\min_{\boldsymbol{w}, \boldsymbol{\alpha}} \quad \sum_{i=1}^{N} \alpha_i \hat{\mathcal{L}}_i(\boldsymbol{w}) + \frac{\lambda}{2} \sum_{i=1}^{N} \frac{\alpha_i^2}{m_i}, \tag{6}$$
$$s.t. \quad \boldsymbol{\alpha} \in \mathbb{R}_+^n, \mathbf{1}^{\top} \boldsymbol{\alpha} = 1,$$

where $\boldsymbol{w}$ is a vector of parameters defining a predictor $h$. Here and after we use notation $\hat{\mathcal{L}}_i(\boldsymbol{w})$ to replace $\hat{\mathcal{L}}_i(h)$, representing the empirical risk of hypothesis $h$ (corresponding to $\mathbf{w}$) on client $i$. Note that removing the square root of the last term does not affect the optimal solution, since the total sum $\sum_{i=1}^{N} \frac{\alpha_i^2}{m_i}$ is the same in both forms, and the square root is simply a scaling for importance of the total sum. In practice, we can achieve the equivalent purpose by adjusting $\lambda$.

The second term of the objective is small whenever the weights are distributed proportionally to the number of samples of the client. As $\lambda \to \infty$, we have $\alpha_i(\boldsymbol{w}) = \frac{m_i}{M}$, which means that all clients are assigned with weights proportional to their number of training samples and the model minimizes the empirical risk over all the data, regardless of the losses of the clients. Thus, the objective becomes the same as the standard `FedAvg` in Eq. (3). In contrast, as $\lambda \to 0$, the regularization term in Eq. (6) vanishes, so that the client with the lowest empirical risk will dominate the objective by setting its weight to 1. $\lambda$ thus acts as a form of regularization by encouraging the usage of information from more clients.

### 4.3    Robustness of the Objective

We now study the optimal weights $\boldsymbol{\alpha}$ of the problem in Eq. (6) to understand how the objective yields robustness against corrupted data sources. We notice that the objective is a convex quadratic program problem over $\boldsymbol{\alpha}$. Given that $\boldsymbol{\alpha} = N^{-1}\mathbf{1}$ is a strictly feasible point, the problem satisfies Slater's condition, which indicates the strong duality of the problem. Thus, the optimal weights $\boldsymbol{\alpha}$ can be obtained using the Karush-Kuhn-Tucker (KKT) conditions (Boyd et al., 2004). Here we give the closed-form solution in Theorem 2. The detailed proof is provided in Appendix B.

**Theorem 2.** *For any $\boldsymbol{w}$, when $\lambda > 0$ and $\{\hat{\mathcal{L}}_i(\boldsymbol{w})\}_{i=1}^N$ are sorted in increasing order: $\hat{\mathcal{L}}_1(\boldsymbol{w}) \leq \hat{\mathcal{L}}_2(\boldsymbol{w}) \leq ... \leq \hat{\mathcal{L}}_N(\boldsymbol{w})$, by setting:*

$$p = \underset{k}{\arg\max}\{1 + \frac{M_k(\overline{\mathcal{L}}_k(\boldsymbol{w}) - \hat{\mathcal{L}}_k(\boldsymbol{w}))}{\lambda} > 0\}, \tag{7}$$

*where $M_k = \sum_{i=1}^k m_i$,*

$$\overline{\mathcal{L}}_k(h) = \frac{\sum_{i=1}^k m_i \hat{\mathcal{L}}_i(\boldsymbol{w})}{M_k} \tag{8}$$

*is the average loss over the first $k$ clients that have the smallest empirical risks. Then the optimal $\boldsymbol{\alpha}$ to the problem (6) is given by:*

$$\alpha_i(\boldsymbol{w}) = \frac{m_i}{M_p}[1 + \frac{M_p(\overline{\mathcal{L}}_p(\boldsymbol{w}) - \hat{\mathcal{L}}_i(\boldsymbol{w}))}{\lambda}]_+, \tag{9}$$

*where $[\cdot]_+ = max(0, \cdot)$.*

When $\lambda \in (0, \infty)$, plugging $\alpha_i(\boldsymbol{w})$ back into Eq. (6) yields the equivalent concentrated objective

$$\sum_{i=1}^N [\frac{m_i}{M_p} + \frac{m_i(\overline{\mathcal{L}}_p(\boldsymbol{w}) - \hat{\mathcal{L}}_i(\boldsymbol{w}))}{\lambda}]_+(\hat{\mathcal{L}}_i(\boldsymbol{w}) + \frac{1}{2}[\frac{\lambda}{M_p} + \overline{\mathcal{L}}_p(\boldsymbol{w}) - \hat{\mathcal{L}}_i(\boldsymbol{w})]_+), \tag{10}$$

which consists of $N$ components and each is related to the empirical risk of the corresponding client and the average loss over the first $p$ clients with the smallest losses, which is called the $p$-average loss. An intuitive interpretation is that the $p$ clients act as a consensus group to reallocate the weights and encourage trusting clients that provide empirical losses that are smaller than the average while downweighting the clients with higher losses. Given a suitable $\lambda$, more benign clients will be in the consensus group to dominate the model and exclude the outliers.

In a situation where the majority of clients are benign, $\overline{\mathcal{L}}_p(\boldsymbol{w})$ becomes relatively low as the benign clients achieve the minimum empirical risk. The components with corrupted datasets will be downweighted as they have higher losses. Especially, the $i$-th component becomes zero when $\hat{\mathcal{L}}_i(\boldsymbol{w}) \geq \frac{\lambda}{M_p} + \overline{\mathcal{L}}_p(\boldsymbol{w})$, which means that a client is considered to be corrupted and does not contribute to the objective if its empirical risk is significantly larger than the $p$-average loss, where the threshold $\frac{\lambda}{M_p} + \overline{\mathcal{L}}_p(\boldsymbol{w})$ is controlled by $\lambda$. From Eq. (7) we can also conclude that the optimal solution has only $p$ non-zero components and the remaining components will be exactly zero.

On the other hand, if the corrupted clients try to bias the model to fit their corrupted datasets, the $p$-average loss $\overline{\mathcal{L}}_p(\boldsymbol{w})$ becomes higher because the model does not fit the samples from the majority. The threshold $\frac{\lambda}{M_p} + \overline{\mathcal{L}}_p(\boldsymbol{w})$ will also be enlarged, which makes $\alpha_i(\boldsymbol{w})$ fail to downweight the component with high losses. Thus, the optimization problem in Eq. (6) will exceed the minimum.

### 4.4    Blockwise Minimization Algorithm

To solve the robust learning problem in federated learning settings, we propose an optimization method based on the blockwise updating paradigm, which is guaranteed to converge to a critical point when the parameter set is closed and convex (Grippo and Sciandrone, 2000). The key idea is to divide the problem into two parts. One sub-problem

for estimating the model parameters and the other sub-problem for automatically weighting the importance of client updates. Then we minimize the objective iteratively w.r.t. one variable each time while fixing the other one.

The pseudocode of the optimization procedure is given in Algorithm 1. At the beginning of the algorithm, we initialize $\hat{\mathcal{L}} = [\hat{\mathcal{L}}_1, \hat{\mathcal{L}}_2, ..., \hat{\mathcal{L}}_N]^\top$ by broadcasting the initial global model $\boldsymbol{w_0}$ to each client to measure its training loss and return it to the server.

---

**Algorithm 1** Optimization of `ARFL`

---

   **Server executes:**
1: Initialize $\boldsymbol{w}_0, \hat{\mathcal{L}}, \boldsymbol{\alpha}$
2: **for** each round $t = 1, 2, \ldots$ **do**
3:    Select a subset $S_t$ from $N$ clients at random
4:    Broadcast the global model $\boldsymbol{w}_t$ to selected clients $S_t$
5:    **for** each client $i \in S_t$ **in parallel do**
6:       $\boldsymbol{w}_{t+1}^i, \hat{\mathcal{L}}_i \leftarrow \text{ClientUpdate}(i, \boldsymbol{w}_t)$
7:    **end for**
8:    Update $\boldsymbol{w}_{t+1}$ according to Eq. (11)
9:    Update $\boldsymbol{\alpha}$ according to Theorem 2
10: **end for**
11:
   **ClientUpdate**$(i, \boldsymbol{w})$:   *// Run on client $i$*
12: $\mathcal{L}_i \leftarrow$ (evaluate training loss using training set)
13: $\mathcal{B} \leftarrow$ (split local training set into batches of size $B$)
14: **for** each local epoch $i$ from 1 to $E$ **do**
15:    **for** batch $b \in \mathcal{B}$ **do**
16:       $\boldsymbol{w} \leftarrow \boldsymbol{w} - \eta \nabla \ell(\boldsymbol{w}; b)$
17:    **end for**
18: **end for**
19: return $\boldsymbol{w}$ and $\hat{\mathcal{L}}_i$

---

**Updating $\boldsymbol{w}$.** When $\boldsymbol{\alpha}$ is fixed, similar to the standard `FedAvg` approach, at round $t$, the server selects a subset $S_t$ of clients at random (Line 3) and broadcasts the global model $\boldsymbol{w}_t$ to the selected clients (Line 4). For each client $i$, it firstly evaluates its training loss $\mathcal{L}_i$ using its local dataset. Then, the model parameters can be updated by local computation with a few steps of SGD (Line 14-18), after which the client uploads the new model parameters $\boldsymbol{w}$ along with $\mathcal{L}_i$ to the server (Line 6). While at the aggregation step, the server assembles the global model as:

$$\boldsymbol{w}_{t+1} \leftarrow \sum_{i \in S_t} \frac{\alpha_i}{\sum_{i \in S_t} \alpha_i} \boldsymbol{w}_{t+1}^i. \tag{11}$$

**Updating $\boldsymbol{\alpha}$.** When $\boldsymbol{w}$ is fixed, we update $\boldsymbol{\alpha}$ using Eq. (9) in Theorem 2. Intuitively, in order to update $\boldsymbol{\alpha}$, the server should broadcast the updated model parameters to all clients to obtain their training loss before updating $\boldsymbol{\alpha}$. Unfortunately, such behavior might significantly increase the burden of the communication network. To improve communication efficiency, we only update the losses from those selected clients while keeping the others unchanged. In other words, the weights of clients are reallocated according to their latest empirical losses. If a client is not selected in the current round, the last updated loss is used instead.

## 5 Experimental Results

### 5.1 Experimental Setup

We implement `ARFL` and the considered baseline methods in TensorFlow (Abadi et al., 2016) Version 2.3[1], simulating a federated learning system with one server and $N$ clients. We perform our experimental evaluation on three datasets that are commonly used in previous work (Li et al., 2019; Wang et al., 2019; McMahan et al., 2017), namely CIFAR-10 (van Dyk and Meng, 2001), FEMNIST (Cohen et al., 2017; Caldas et al., 2018), and Shakespeare (Caldas et al., 2018; McMahan et al., 2017). Their basic information is listed in Table 1. For the CIFAR-10 dataset, we consider an i.i.d. partition where each local client has approximately the same amount of samples and in proportion to each of the classes.

---

[1]Code is available at `https://github.com/lishenghui/arfl`

Table 1: Dataset description and parameters

| Dataset | #Classes | #Clients | #Samples | i.i.d. | Model used | $l_r$ | $E$ | Batch size | $|S_t|$ | #Rounds |
|---------|----------|----------|----------|--------|------------|-------|-----|------------|---------|---------|
| CIFAR-10 | 10 | 100 | 60,000 | Yes | CNN | 0.01 | 5 | 64 | 20 | 2000 |
| FEMNIST | 62 | 1039 | 236,500 | No | CNN | 0.01 | 20 | 64 | 32 | 2000 |
| Shakespeare | 80 | 71 | 417,469 | No | LSTM | 0.6 | 1 | 10 | 16 | 100 |

We use the original test set in CIFAR-10 as our global test set for comparing the performance of all methods. For the Shakespeare and FEMNIST datasets, we treat each speaking role or writer as a client and randomly sample subsets of all clients. We assume that data distributions vary among clients in the raw data, and hence we regard this sampling process as non-i.i.d.. The raw data in FEMNIST and Shakespeare is preprocessed using the popular benchmark LEAF (Caldas et al., 2018), where the data on each local client is partitioned into an 80% training set and a 20% testing set. The details of the network models we use in the experiments are as follows:

- The model for CIFAR-10 is a Convolutional Neural Network (CNN) chosen from Tensorflow's website[2], which consists of three 3x3 convolution layers (the first with 32 channels, the second and third with 64, the first two followed by 2x2 max pooling), a fully connected layer with 64 units and ReLu activation, and a final softmax output layer. To improve the performance, data augmentation (random shift and flips) is used in this dataset (Wang et al., 2019).

- For the FEMNIST dataset, we train a CNN with two 5x5 convolution layers (the first with 32 channels, the second with 64, each followed by 2x2 max pooling), a fully connected layer with 126 units and ReLu activation, and a final softmax output layer.

- For the Shakespeare dataset, we learn a character-level language model to predict the next character over *the Complete Works of Shakespeare* (Shakespeare, [n.d.]). The model takes a series of characters as input and embeds each of these into an 8-dimensional space. The embedded features are then processed through two stacked Long Short-Term Memory (LSTM) layers, each with 256 nodes and a dropout value of 0.2. Finally, the output of the second LSTM layer is sent to a softmax output layer with one node per character.

For each dataset we consider four different scenarios: 1) **Normal operation** (*clean*): we use the original datasets without any corruption. 2) **Label shuffling** (*shuffling*): the labels of all samples are shuffled randomly in each corrupted client. 3) **Label flipping** (*flipping*): the labels of all samples are switched to a random one in each corrupted client, which means that all labels of training samples are flipped as the same one for each corrupted client. 4) **Noisy clients** (*noisy*): for CIFAR-10 and FEMNIST datasets, we normalize the inputs to the interval $[0, 1]$. In this scenario, for the selected noisy clients we add Gaussian noise to all the pixels, so that $x \leftarrow x + \epsilon$, with $\epsilon \sim N(0, 0.7)$. Then we normalize the resulting values again to the interval $[0, 1]$. For the Shakespeare dataset, we randomly select half of the characters and shuffle them so that the input sentence might be disordered. For each corruption scenario, we set 30% and 50% of the clients to be corrupted clients (i.e. providing corrupted data).

We empirically tune the hyper-parameters on `ARFL` and use the same values in all experiments of each dataset. We use the parameter setups in Table 1, unless specified otherwise. Following the standard setup, we use SGD and train for $E$ local epochs with local learning rate $l_r$. A shared global model is trained by all participants, a subset $S_t$ is randomly selected in each round of local training, and $|S_t|$ is the size of $S_t$. By default, we use a large $\lambda$ ($\lambda = 10000 \times M$) for clean data, and use a relatively small $\lambda$ ($\lambda = M$) for all corruptions, where $M$ is the total amount of training samples. We repeat every experiment five times with different random seeds for data corruption and client selection, and evaluate the accuracy of the learned model with the clean test set.

We compare the performance of `ARFL` with the following state-of-the-art solutions:

- **FedAvg** (McMahan et al., 2017). The standard Federated Averaging aggregation approach that just calculates the weighted average of the parameters from local clients.

- **RFA** (Pillutla et al., 2019). A robust aggregation approach that minimizes the weighted Geometric Median (GM) of the parameters from local clients. A smoothed Weiszfeld's algorithm is used to compute the approximate GM.

- **MKrum (Multi-Krum)** (Blanchard et al., 2017). A Byzantine tolerant aggregation rule. Note that this approach tolerates some Byzantine failures such as completely arbitrary behaviors from local updates.

---

[2]https://www.tensorflow.org/tutorials/images/cnn

Table 2: Averaged test accuracy over five random seeds for `FedAvg`, `RFA`, `MKRUM`, `CFL` and `ARFL` in four different scenarios. In the *shuffling* and *flipping* scenarios, `ARFL` significantly outperforms the others. In the *clean* and *noisy* scenario, `FedAvg`, `RFA`, `MKRUM` and `ARFL` achieve similar accuracy.

| CIFAR-10 | Clean | Shuffling | | Flipping | | Noisy | |
|---|---|---|---|---|---|---|---|
| Corr. Per. | - | 30% | 50% | 30% | 50% | 30% | 50% |
| FedAvg (McMahan et al., 2017) | $73.59 \pm 0.44$ | $61.17 \pm 1.81$ | $47.00 \pm 7.51$ | $65.01 \pm 2.38$ | $51.75 \pm 7.75$ | $\mathbf{73.75 \pm 0.49}$ | $73.61 \pm 0.53$ |
| RFA (Pillutla et al., 2019) | $71.36 \pm 0.47$ | $57.86 \pm 3.22$ | $40.26 \pm 9.14$ | $55.47 \pm 5.17$ | $40.91 \pm 11.06$ | $73.74 \pm 0.52$ | $\mathbf{73.69 \pm 0.63}$ |
| MKRUM (Blanchard et al., 2017) | $67.03 \pm 0.93$ | $59.27 \pm 9.34$ | $52.32 \pm 14.90$ | $60.21 \pm 5.73$ | $47.96 \pm 10.25$ | $73.41 \pm 0.69$ | $73.49 \pm 0.49$ |
| CFL (Sattler et al., 2020) | $71.68 \pm 0.36$ | $52.54 \pm 1.71$ | $50.29 \pm 1.95$ | $52.87 \pm 1.07$ | $51.67 \pm 0.92$ | $54.97 \pm 1.14$ | $55.26 \pm 1.96$ |
| ARFL (ours) | $73.42 \pm 0.40$ | $\mathbf{71.68 \pm 1.01}$ | $\mathbf{69.66 \pm 0.73}$ | $\mathbf{71.78 \pm 0.53}$ | $\mathbf{70.25 \pm 0.56}$ | $73.48 \pm 0.56$ | $73.29 \pm 0.79$ |
| **FEMNIST** | **Clean** | **Shuffling** | | **Flipping** | | **Noisy** | |
| Corr. Per. | - | 30% | 50% | 30% | 50% | 30% | 50% |
| FedAvg (McMahan et al., 2017) | $82.12 \pm 0.20$ | $61.91 \pm 21.33$ | $39.69 \pm 20.80$ | $70.19 \pm 10.17$ | $48.53 \pm 23.49$ | $79.94 \pm 0.36$ | $78.27 \pm 0.47$ |
| RFA (Pillutla et al., 2019) | $82.11 \pm 0.32$ | $74.36 \pm 7.52$ | $52.02 \pm 22.51$ | $73.80 \pm 7.49$ | $50.75 \pm 19.91$ | $80.45 \pm 0.30$ | $79.21 \pm 0.41$ |
| MKRUM (Blanchard et al., 2017) | $79.38 \pm 0.41$ | $57.51 \pm 21.17$ | $42.40 \pm 24.84$ | $78.57 \pm 4.83$ | $67.10 \pm 7.35$ | $\mathbf{81.52 \pm 0.53}$ | $\mathbf{79.80 \pm 0.22}$ |
| CFL (Sattler et al., 2020) | $82.18 \pm 0.30$ | $81.24 \pm 0.47$ | $36.03 \pm 36.38$ | $81.22 \pm 0.36$ | $65.54 \pm 26.94$ | $80.13 \pm 0.70$ | $79.21 \pm 0.64$ |
| ARFL (ours) | $\mathbf{82.32 \pm 0.19}$ | $\mathbf{81.60 \pm 0.31}$ | $\mathbf{81.35 \pm 0.43}$ | $\mathbf{81.87 \pm 0.22}$ | $\mathbf{81.30 \pm 0.24}$ | $80.71 \pm 0.28$ | $79.40 \pm 0.45$ |
| **Shakespeare** | **Clean** | **Shuffling** | | **Flipping** | | **Noisy** | |
| Corr. Per. | - | 30% | 50% | 30% | 50% | 30% | 50% |
| FedAvg (McMahan et al., 2017) | $53.80 \pm 0.33$ | $51.98 \pm 0.48$ | $47.70 \pm 4.96$ | $52.08 \pm 0.39$ | $41.85 \pm 16.18$ | $51.85 \pm 0.56$ | $50.43 \pm 1.19$ |
| RFA (Pillutla et al., 2019) | $\mathbf{54.27 \pm 0.41}$ | $50.16 \pm 1.28$ | $32.49 \pm 13.81$ | $50.50 \pm 1.02$ | $23.84 \pm 21.78$ | $\mathbf{52.17 \pm 0.50}$ | $50.69 \pm 1.04$ |
| MKRUM (Blanchard et al., 2017) | $50.81 \pm 0.85$ | $40.38 \pm 7.44$ | $24.46 \pm 6.88$ | $44.95 \pm 2.43$ | $16.11 \pm 15.46$ | $48.19 \pm 0.40$ | $45.67 \pm 0.46$ |
| CFL (Sattler et al., 2020) | $54.01 \pm 0.34$ | $49.76 \pm 4.47$ | $43.68 \pm 12.68$ | $51.09 \pm 1.36$ | $37.30 \pm 19.76$ | $51.98 \pm 1.03$ | $50.38 \pm 1.39$ |
| ARFL (ours) | $53.52 \pm 0.32$ | $\mathbf{52.85 \pm 0.49}$ | $\mathbf{51.61 \pm 0.68}$ | $\mathbf{52.82 \pm 0.48}$ | $\mathbf{51.74 \pm 0.69}$ | $52.09 \pm 1.27$ | $\mathbf{50.98 \pm 0.75}$ |



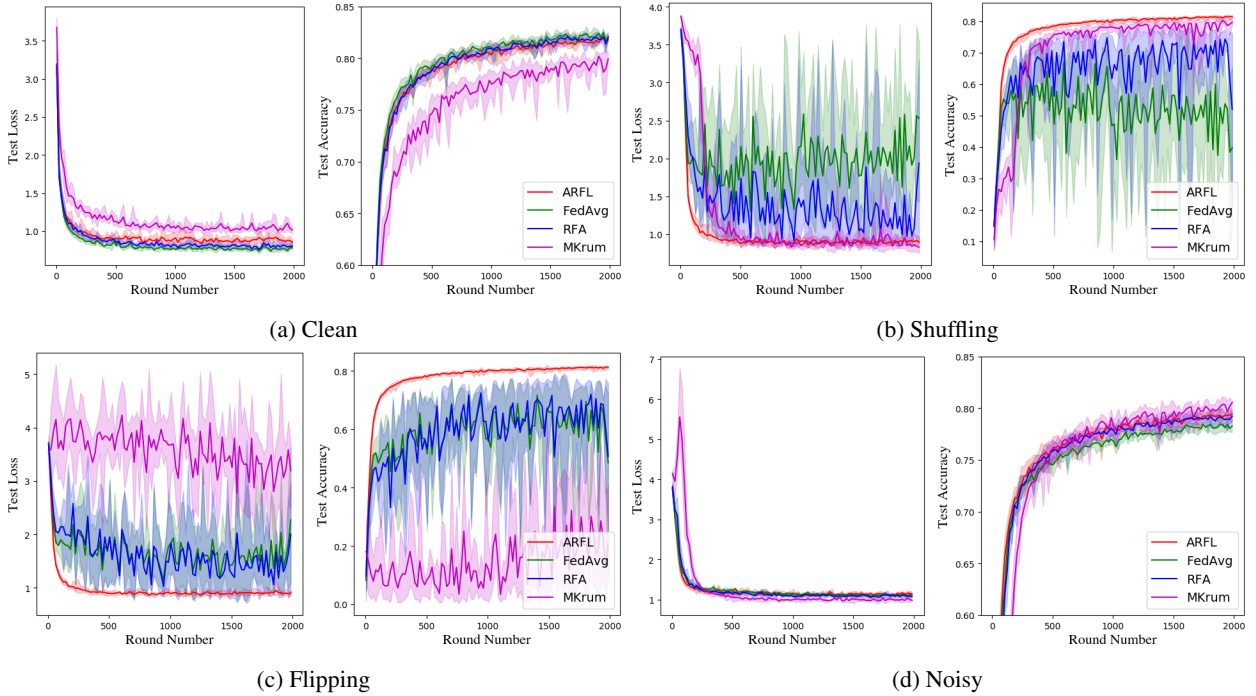(a) Clean

(b) Shuffling

(c) Flipping

(d) Noisy

Figure 2: Test loss and accuracy vs. round number for `FedAvg`, `RFA`, `MKRUM` and `ARFL` on the FEMNIST dataset with $50\%$ clients with different corruption scenarios. In the *shuffling* and *flipping* corruption scenarios, `ARFL` converges to the highest accuracy among the four approaches.

- **CFL** (Sattler et al., 2020). A Clustered Federated Learning (CFL) approach that separates the client population into different groups based on the pairwise cosine similarities between their parameter updates, where the clients are partitioned into two groups, i.e., benign clients and corrupted clients.

## 5.2 Robustness and Convergence

Firstly, we show that `ARFL` is a more robust solution by comparing the average test accuracy in Table 2. The results show that `ARFL` is robust in all data corruption scenarios and corruption levels. It achieves the highest test accuracy in most scenarios compared with the other approaches. `ARFL` achieves significantly higher test accuracy in the *shuffling* and *flipping* corruption scenarios compared with all the existing methods, which is especially noticeable in the case
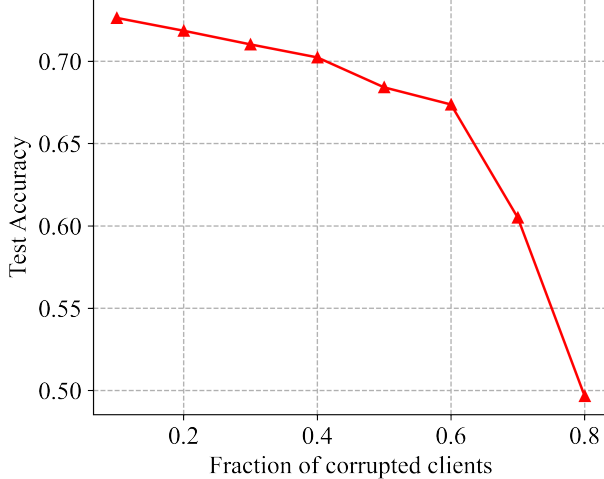
Figure 3: Impact of fraction of corrupted clients on test accuracy of `ARFL` using CIFAR-10. The accuracy decreases as the fraction of corrupted clients increases, especially when 60% clients are corrupted.
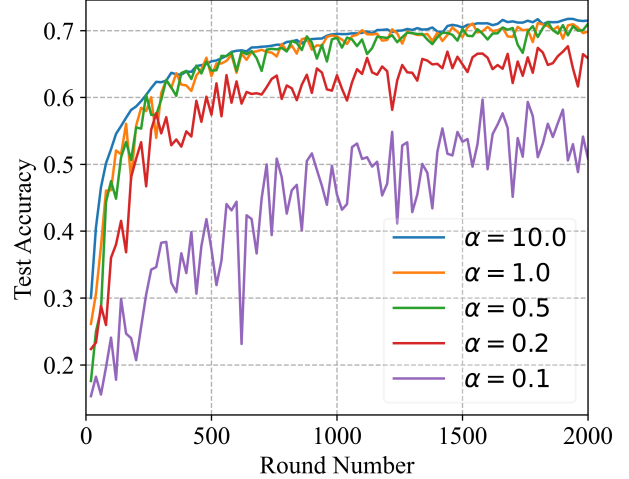


Figure 4: Learning curves of `ARFL` on CIFAR-10 with different non-i.i.d settings, where smaller $\alpha$ indicates stronger non-i.i.d. data distribution. The curves fluctuates as we increase the degree of non-i.i.d.

of *flipping* corruption with the level of $50\%$, where the test accuracy for the CIFAR-10, FEMNIST and Shakespeare datasets are $70.25\%, 81.30\%$ and $51.74\%$, which are $18.5\%, 14.2\%$, and $9.89\%$ higher than that of the best of existing methods, respectively. In the *clean* and *noisy* scenarios, `ARFL`'s test accuracy is very close to the best method in the comparison.

As expected, `FedAvg`'s performance is significantly affected by the presence of corrupted clients, especially in *shuffling* and *flipping* scenarios. Furthermore, `MKRUM` also shows poor performance in *shuffling* and *flipping* scenarios of all datasets. `RFA` works well for the FEMNIST dataset, but worse than `FedAvg` in the *shuffling* and *flipping* scenarios for the CIFAR-10 and Shakespere datasets. It is also interesting to observe that `CFL` works well for the FEMNIST and Shakespeare datasets under $30\%$ corruption level, but the accuracy decreases significantly when the corruption level is $50\%$. The reason is that when half of the clients are corrupted, `CFL` fails to identify which group of clients are corrupted. These results demonstrate that `ARFL` offers better performance than the existing approaches across the corruption scenarios we consider. Note that our approach can handle even higher corruption rates in those scenarios. For example, using the FEMNIST data set with $70\%$ corrupted clients we still achieve an accuracy above $79\%$. However, it is also noticeable that if multiple clients try to bias their data to the same distribution (i.e., colluding corruption), our approach is unable to handle such a high corruption rate.

Next, we study the convergence of the approaches by comparing the test loss and accuracy of the global model versus the number of training rounds in Figure 2, where $50\%$ of the clients are corrupted. As discussed before, `CFL` is unable to handle such a high corruption level. Therefore we only compare the remaining four approaches. The shaded areas denote the minimum and maximum values over five repeated runs. The figure shows that all approaches converge to a good solution in the *clean* and *noisy* scenarios. However, in the *clean* scenario, the test loss of `ARFL` is slightly higher than the others. The reason is that the global model could bias towards some of the local updates, which can be avoided by increasing $\lambda$. Furthermore, all the `RFA`, `FedAvg` and `MKRUM` approaches diverge in the *shuffling* and *flipping* corruption scenarios, which indicates that the local data corruption harms the global model during their training process. On the contrary, our `ARFL` method is able to converge with high accuracy, since it is able to lower the contribution of the corrupted clients. In the *clean* and *noisy* scenarios, we observe that `MKRUM` converges slower, as it only uses a subset of selected updates to aggregate the global model.

## 5.3 Impact of the Fraction of Corrupted Clients

To demonstrate the effect of the fraction of corrupted clients to the robustness of `ARFL`, we run an experiment on CIFAR-10 with different number of corrupted clients and apply *shuffling* to corrupt the clients. As shown in Fig. 3, the performance decreases with the fraction of corrupted clients. Noticeably, the accuracy does not significantly change until more than 60% clients are corrupted. Similar phenomenon has been shown by Li et al. in a previous work (Li et al., 2021), which achieves robustness by personalized federated training. However, different from the solution in (Li
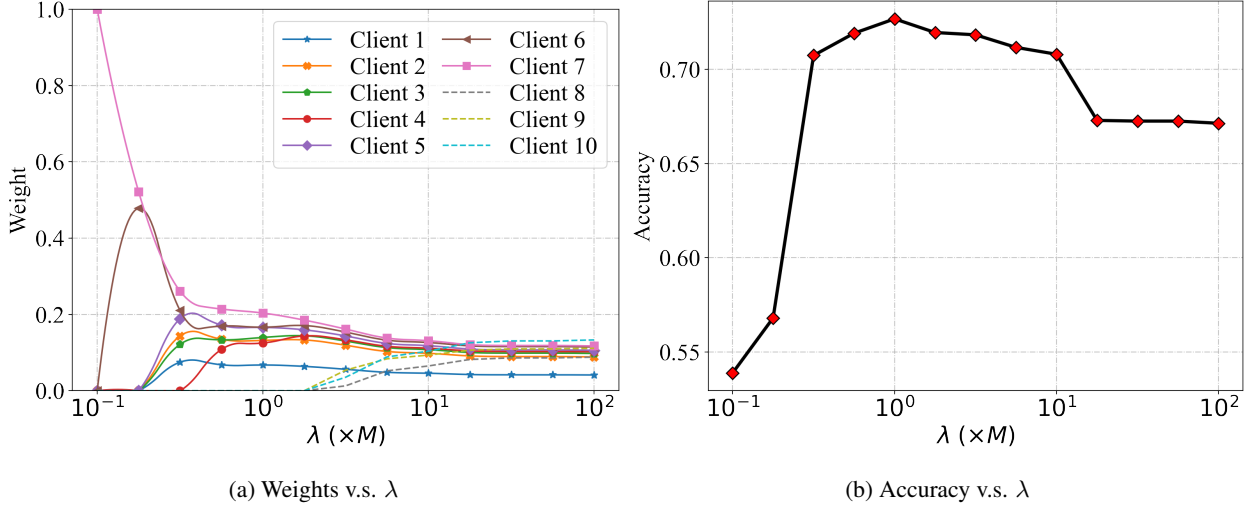
(a) Weights v.s. $\lambda$

(b) Accuracy v.s. $\lambda$

Figure 5: Effects of $\lambda$ on weights and model accuracy on the CIFAR-10 dataset. The corrupted clients (dashed lines) are zero-weighted when $\lambda \leq 2.5 \times M$, and the accuracy reaches its peak when $\lambda = 1 \times M$.

et al., 2021), our `ARFL` only learns a single global model, without the need of extra training overhead for personalized models.

## 5.4   Impact of Data Heterogeneity

In order to study the impact of non-i.i.d. data, i.e., the extent of data heterogeneity among clients, we simulate a non-i.i.d. partition of CIFAR-10, for which the number of data points and class proportions are unbalanced, and 30% of the clients are corrupted with *shuffling*. Following prior arts (Hsu et al., 2019; Lin et al., 2020), we model non-i.i.d. data distributions by using a Dirichlet distribution $\boldsymbol{P}_k \sim Dir_N(\alpha)$ and by allocating a $\boldsymbol{P}_{k,i}$ proportion of the training instances of class $k$ to local client $i$, in which a smaller $\alpha$ indicates stronger non-i.i.d. data distributions. As shown in Fig. 4, when the data is less heterogeneous (e.g., $\alpha = 10.0$), the learning curve increases more steadily, and therefore results in higher accuracy at the end of the training. On the contrary, the learning curve fluctuates over the training rounds when $\alpha$ is small. One explanation for this is that it becomes hard to induce a consensus model for the benign clients if their data distributions are significantly different, thus the corrupted clients can damage the global model more easily as they are not down-weighted properly.

## 5.5   Tuning $\lambda$

As mentioned before, $\lambda$ in the objective should be tuned in order to provide a trade-off between robustness and average accuracy. Here we conduct further experiments on the CIFAR-10 dataset to study the effects of $\lambda$. The dataset is partitioned into $N = 10$ clients by sampling $\boldsymbol{P}_k \sim Dir_N(0.5)$ and allocating a $\boldsymbol{P}_{k,i}$ proportion of the training instances of class $k$ to local client $i$, where three clients are corrupted by shuffling their labels randomly (Client 8-10). We use the original test set in CIFAR-10 as the global test set. We train the model for 1000 rounds where each client runs one epoch of SGD on their training set before each aggregation, where $\lambda$ is set in the range of $[10^{-1} \times M, 10^2 \times M]$. All the other settings are the same as in Sec. 5.1.

Fig. 5(a) shows the optimized weights as a function of $\lambda$ on the CIFAR-10 dataset. It is readily apparent that $\boldsymbol{\alpha}$ has only one non-zero element (Client 7) for small $\lambda$ and all elements of $\boldsymbol{\alpha}$ come to certain non-zero values for large $\lambda$. In between these two extremes, we obtain sparse solutions of $\boldsymbol{\alpha}$ in which only a part of elements have non-zero values. It is also noticeable that all the corrupted clients (dashed lines) are zero-weighted when $\lambda \leq 2.5 \times M$, which means that they make no contribution when the server aggregates the updated local models.

In Fig. 5(b), we plot the model accuracy as a function of $\lambda$, showing that the model achieves relatively low accuracy for small $\lambda$, which demonstrates that extremely sparse weights are not favorable under this non-i.i.d. data setting. The reason behind this is that the model finally fits only one local dataset without considering data from other clients. The accuracy increases as more benign clients are upweighted and contribute their local updates to the global model. The optimal value for $\lambda$ is $1.0 \times M$ in this experiment, where all the benign clients have non-zero weights, while all the
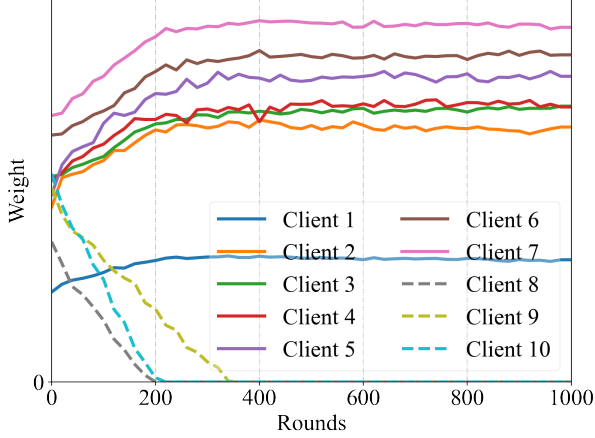
11

Figure 6: Visualization of the weights v.s. round number of `ARFL`. All corrupted clients (dashed lines) are zero-weighted after 350 rounds of training.
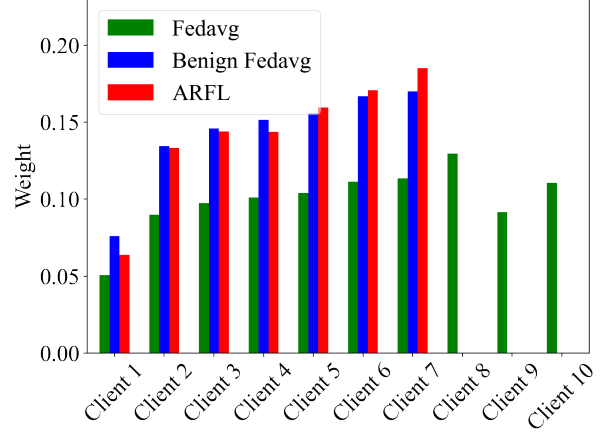
Figure 7: Weights comparison between ARFL and FedAvg. ARFL's weights distributed similar to those of FedAvg with only benign clients

corrupted clients have zero weights. However, as $\lambda$ further increases, the global model is harmed by the corrupted clients as they gain adequate weights, which leads to lower accuracy.

### 5.6 Auto-weighting Analysis

In this part, we investigate the auto-weighting process during training in `ARFL` when $\lambda = 1.0 \times M$, where we keep the setup the same as the previous above. As shown in Fig. 6, our approach successfully downweights all the three corrupted clients, for which the weights become zero after 350 rounds of training.

We next compare the reallocated weights with the standard `FedAvg`, where the weights are fixed as $\alpha_i = \frac{m_i}{M}$, and `FedAvg` with only benign clients (Benign `FedAvg`), where $\alpha_i = \frac{m_i}{M_\mathcal{B}}$ for benign clients, and $\alpha_i = 0$ for corrupted clients. Fig. 7 shows that the learned weight distribution (in red) of `ARFL` is approximated to the distribution considering only benign clients (in blue). We conclude that our approach can automatically reweight the clients and approximate the mixture distribution to the benign uniform distribution during the model training process, even when the centralized server is agnostic to the local corruptions.

## 6 Limitations

The first limitation is the fairness. When a small $\lambda$ is chosen, the clients are treated unfairly, e.g., a client with data that is difficult to learn (but not due to corruption) would carry less weight. This happens in some strongly non-i.i.d cases where the distributions among sources are inherently different. However, we argue that robustness and fairness are two competing targets (Li et al., 2021), it is technically difficult to distinguish between "corrupted" and "just different" (but not corrupted) data sources if the data is strongly non-i.i.d., which we leave it for future work. Nevertheless, we show in our experiments (Sec. 5) that our approach performs well in some general non-i.i.d. settings.

Another limitation is on the assumption of corruption strategies. Colluding corruptions among clients, e.g., multiple clients are trying to bias their data distribution to the same target, may bring risk to the robustness of the objective. In addition, our approach is not designed to handle malicious attacks on the model parameters, i.e., only data corruption scenarios are considered. Our work is based on the assumption that all clients are honest during local training and communication steps, which means that the training loss should reflect the real empirical loss of the global model on the local data. If some corrupted clients cheat the server by giving extremely small fake losses, then they will be allocated with high weights and dominate the training accordingly.

## 7 Conclusions and Future work

In this article, we proposed Auto-weighted Robust Federated Learning (`ARFL`), a novel approach that automatically re-weights the local updates to lower the contribution of corrupted clients who provide low-quality updates to the global

model. Experimental results on benchmark datasets corroborate the competitive performance of `ARFL` compared to the state-of-the-art methods under different data corruption scenarios. Our future work will focus on robust aggregation without the losses provided by clients since evaluating training loss from local clients could also be a potential overhead. Also, we will explore the possibility of further improving the robustness of federated learning and coping with higher fraction of corrupted clients. Furthermore, we plan to extend `ARFL` to a more general federated learning approach and make it robust against both model poisoning and data corruption.

# References

Martín Abadi, Ashish Agarwal, Paul Barham, Eugene Brevdo, Zhifeng Chen, Craig Citro, Greg S Corrado, Andy Davis, Jeffrey Dean, Matthieu Devin, et al. 2016. Tensorflow: Large-scale machine learning on heterogeneous distributed systems. *arXiv preprint arXiv:1603.04467* (2016).

Arjun Nitin Bhagoji, Supriyo Chakraborty, Prateek Mittal, and Seraphin Calo. 2019. Analyzing federated learning through an adversarial lens. In *International Conference on Machine Learning*. PMLR, 634–643.

Peva Blanchard, Rachid Guerraoui, Julien Stainer, et al. 2017. Machine learning with adversaries: Byzantine tolerant gradient descent. In *Advances in Neural Information Processing Systems*. 119–129.

Olivier Bousquet, Stéphane Boucheron, and Gábor Lugosi. 2003. Introduction to statistical learning theory. In *Summer School on Machine Learning*. Springer, 169–207.

Stephen Boyd, Stephen P Boyd, and Lieven Vandenberghe. 2004. *Convex optimization*. Cambridge university press.

Sebastian Caldas, Peter Wu, Tian Li, Jakub Konečný, H Brendan McMahan, Virginia Smith, and Ameet Talwalkar. 2018. Leaf: A benchmark for federated settings. *arXiv preprint arXiv:1812.01097* (2018).

Xiaoyu Cao, Minghong Fang, Jia Liu, and Neil Zhenqiang Gong. 2021. FLTrust: Byzantine-robust Federated Learning via Trust Bootstrapping. In *ISOC Network and Distributed System Security Symposium (NDSS)*.

Mingqing Chen, Ananda Theertha Suresh, Rajiv Mathews, Adeline Wong, Cyril Allauzen, Françoise Beaufays, and Michael Riley. 2019. Federated Learning of N-Gram Language Models. In *Proceedings of the 23rd Conference on Computational Natural Language Learning (CoNLL)*. 121–130.

Yudong Chen, Lili Su, and Jiaming Xu. 2017. Distributed statistical machine learning in adversarial settings: Byzantine gradient descent. *Proceedings of the ACM on Measurement and Analysis of Computing Systems* 1, 2 (2017), 1–25.

Gregory Cohen, Saeed Afshar, Jonathan Tapson, and Andre Van Schaik. 2017. EMNIST: Extending MNIST to handwritten letters. In *2017 International Joint Conference on Neural Networks (IJCNN)*. IEEE, 2921–2926.

Minghong Fang, Xiaoyu Cao, Jinyuan Jia, and Neil Gong. 2020. Local model poisoning attacks to byzantine-robust federated learning. In *29th {USENIX} Security Symposium ({USENIX} Security 20)*. 1605–1622.

Luigi Grippo and Marco Sciandrone. 2000. On the convergence of the block nonlinear Gauss–Seidel method under convex constraints. *Operations research letters* 26, 3 (2000), 127–136.

Rachid Guerraoui, Sébastien Rouault, et al. 2018. The Hidden Vulnerability of Distributed Learning in Byzantium. In *International Conference on Machine Learning*. 3521–3530.

Dhruv Guliani, Francoise Beaufays, and Giovanni Motta. 2020. Training Speech Recognition Models with Federated Learning: A Quality/Cost Framework. *arXiv preprint arXiv:2010.15965* (2020).

Jenny Hamer, Mehryar Mohri, and Ananda Theertha Suresh. 2020. FedBoost: A Communication-Efficient Algorithm for Federated Learning. In *International Conference on Machine Learning*. PMLR, 3973–3983.

Yufei Han and Xiangliang Zhang. 2020. Robust Federated Learning via Collaborative Machine Teaching.. In *AAAI*. 4075–4082.

Tzu-Ming Harry Hsu, Hang Qi, and Matthew Brown. 2019. Measuring the effects of non-identical data distribution for federated visual classification. *arXiv preprint arXiv:1909.06335* (2019).

Peter Kairouz, Ziyu Liu, and Thomas Steinke. 2021. The Distributed Discrete Gaussian Mechanism for Federated Learning with Secure Aggregation. In *Proceedings of the 38th International Conference on Machine Learning (Proceedings of Machine Learning Research, Vol. 139)*, Marina Meila and Tong Zhang (Eds.). 5201–5212.

Jakub Konečný, H Brendan McMahan, Daniel Ramage, and Peter Richtárik. 2016a. Federated optimization: Distributed machine learning for on-device intelligence. *arXiv preprint arXiv:1610.02527* (2016).

Jakub Konečný, H Brendan McMahan, Felix X Yu, Peter Richtárik, Ananda Theertha Suresh, and Dave Bacon. 2016b. Federated learning: Strategies for improving communication efficiency. *arXiv preprint arXiv:1610.05492* (2016).

Nikola Konstantinov and Christoph Lampert. 2019. Robust learning from untrusted sources. In *International Conference on Machine Learning*. PMLR, 3488–3498.

Suyi Li, Yong Cheng, Wei Wang, Yang Liu, and Tianjian Chen. 2020a. Learning to Detect Malicious Clients for Robust Federated Learning. *arXiv preprint arXiv:2002.00211* (2020).

Tian Li, Shengyuan Hu, Ahmad Beirami, and Virginia Smith. 2021. Ditto: Fair and robust federated learning through personalization. In *International Conference on Machine Learning*.

Tian Li, Maziar Sanjabi, Ahmad Beirami, and Virginia Smith. 2019. Fair Resource Allocation in Federated Learning. In *International Conference on Learning Representations*.

Xiaoli Li, Nan Liu, Chuan Chen, Zibin Zheng, Huizhong Li, and Qiang Yan. 2020b. Communication-Efficient Collaborative Learning of Geo-Distributed JointCloud from Heterogeneous Datasets. In *2020 IEEE International Conference on Joint Cloud Computing*. IEEE, 22–29.

Tao Lin, Lingjing Kong, Sebastian U Stich, and Martin Jaggi. 2020. Ensemble Distillation for Robust Model Fusion in Federated Learning. In *NeurIPS*.

Xiuming Liu and Edith Ngai. 2019. Gaussian Process Learning for Distributed Sensor Networks under False Data Injection Attacks. In *2019 IEEE Conference on Dependable and Secure Computing (DSC)*. IEEE, 1–6.

Jiahuan Luo, Xueyang Wu, Yun Luo, Anbu Huang, Yunfeng Huang, Yang Liu, and Qiang Yang. 2019. Real-world image datasets for federated learning. *arXiv preprint arXiv:1910.11089* (2019).

Lingjuan Lyu, Han Yu, and Qiang Yang. 2020. Threats to federated learning: A survey. *arXiv preprint arXiv:2003.02133* (2020).

Saeed Mahloujifar, Mohammad Mahmoody, and Ameer Mohammed. 2019. Data Poisoning Attacks in Multi-Party Learning. In *International Conference on Machine Learning*. 4274–4283.

Brendan McMahan, Eider Moore, Daniel Ramage, Seth Hampson, and Blaise Aguera y Arcas. 2017. Communication-efficient learning of deep networks from decentralized data. In *Artificial Intelligence and Statistics*. 1273–1282.

Mehryar Mohri, Gary Sivek, and Ananda Theertha Suresh. 2019. Agnostic Federated Learning. In *International Conference on Machine Learning*. 4615–4625.

Xingchao Peng, Zijun Huang, Yizhe Zhu, and Kate Saenko. 2019. Federated Adversarial Domain Adaptation. In *International Conference on Learning Representations*.

Krishna Pillutla, Sham M Kakade, and Zaid Harchaoui. 2019. Robust aggregation for federated learning. *arXiv preprint arXiv:1912.13445* (2019).

Nuria Rodríguez-Barroso, Eugenio Martínez-Cámara, M Luzón, Gerardo González Seco, Miguel Ángel Veganzones, and Francisco Herrera. 2020. Dynamic Federated Learning Model for Identifying Adversarial Clients. *arXiv preprint arXiv:2007.15030* (2020).

Felix Sattler, Klaus-Robert Müller, Thomas Wiegand, and Wojciech Samek. 2020. On the Byzantine Robustness of Clustered Federated Learning. In *ICASSP 2020-2020 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 8861–8865.

William Shakespeare. [n.d.]. *The complete works of William Shakespeare.* Publicly available at http://www.gutenberg.org/ebooks/100/.

Shai Shalev-Shwartz and Shai Ben-David. 2014. Understanding Machine Learning: From Theory to Algorithms.

Jinhyun So, Başak Güler, and A Salman Avestimehr. 2020. Byzantine-resilient secure federated learning. *IEEE Journal on Selected Areas in Communications* 39, 7 (2020), 2168–2181.

Sebastian U Stich. 2018. Local SGD Converges Fast and Communicates Little. In *International Conference on Learning Representations*.

Dianbo Sui, Yubo Chen, Jun Zhao, Yantao Jia, Yuantao Xie, and Weijian Sun. 2020. FedED: Federated Learning via Ensemble Distillation for Medical Relation Extraction. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*. 2118–2128.

David A van Dyk and Xiao-Li Meng. 2001. The Art of Data Augmentation. *Journal of Computational and Graphical Statistics* 10, 1 (2001), 1–50.

Paul Wais, Shivaram Lingamneni, Duncan Cook, Jason Fennell, Benjamin Goldenberg, Daniel Lubarov, David Marin, and Hari Simons. 2010. Towards Building a High-Quality Workforce with Mechanical Turk. In *In Proc. NIPS Workshop on Computational Social Science and the Wisdom of Crowds*. Citeseer.

Hongyi Wang, Mikhail Yurochkin, Yuekai Sun, Dimitris Papailiopoulos, and Yasaman Khazaeni. 2019. Federated Learning with Matched Averaging. In *International Conference on Learning Representations*.

Jianyu Wang and Gauri Joshi. 2019. Cooperative SGD: A Unified Framework for the Design and Analysis of Communication-Efficient SGD Algorithms. In *ICML Workshop on Coding Theory for Machine Learning*.

Qiang Yang, Yang Liu, Tianjian Chen, and Yongxin Tong. 2019a. Federated machine learning: Concept and applications. *ACM Transactions on Intelligent Systems and Technology (TIST)* 10, 2 (2019), 1–19.

Qiang Yang, Yang Liu, Yong Cheng, Yan Kang, Tianjian Chen, and Han Yu. 2019b. Federated learning. *Synthesis Lectures on Artificial Intelligence and Machine Learning* 13, 3 (2019), 1–207.

Timothy Yang, Galen Andrew, Hubert Eichner, Haicheng Sun, Wei Li, Nicholas Kong, Daniel Ramage, and Françoise Beaufays. 2018. Applied federated learning: Improving google keyboard query suggestions. *arXiv preprint arXiv:1812.02903* (2018).

Dong Yin, Yudong Chen, Ramchandran Kannan, and Peter Bartlett. 2018. Byzantine-Robust Distributed Learning: Towards Optimal Statistical Rates. In *International Conference on Machine Learning*. 5650–5659.

Dong Yin, Ramchandran Kannan, and Peter Bartlett. 2019. Rademacher complexity for adversarially robust generalization. In *International Conference on Machine Learning*. PMLR, 7085–7094.

Hao Yu, Sen Yang, and Shenghuo Zhu. 2019. 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*, Vol. 33. 5693–5700.

Shuai Zheng, Ziyue Huang, and James Kwok. 2019. Communication-efficient distributed blockwise momentum SGD with error-feedback. In *Advances in Neural Information Processing Systems*. 11450–11460.

## A  Proof of Theorem 1

**Theorem 1.** *We denote $\hat{\mathcal{L}}_i(h)$ as the corresponding empirical counterparts of $\mathcal{L}_{\mathcal{D}_i}(h)$. Assume that the loss function $\ell_h(z)$ is bounded by a constant $\mathcal{M} > 0$. Then, for any $\delta \in (0, 1]$, with probability at least $1 - \delta$ over the data, the following inequality holds:*

$$\mathcal{L}_{\mathcal{D}_\alpha}(h) \leq \sum_{i=1}^{N} \alpha_i \hat{\mathcal{L}}_i(h) + 2 \sum_{i=1}^{N} \alpha_i \mathcal{R}_i(\mathcal{H}) + 3\sqrt{\frac{\log\left(\frac{4}{\delta}\right)\mathcal{M}^2}{2}} \sqrt{\sum_{i=1}^{n} \frac{\alpha_i^2}{m_i}},$$

*where, for each client $i = 1, \ldots, N$,*

$$\mathcal{R}_i(\mathcal{H}) = \mathbb{E}_\sigma \left( \sup_{f \in \mathcal{H}} \left( \frac{1}{m_i} \sum_{j=1}^{m_i} \sigma_{i,j} \ell_f(z_{i,j}) \right) \right)$$

*and $\sigma_{i,j}$s are independent uniform random variables taking values in $\{-1, +1\}$. These variables are called Rademacher random variables (Yin et al., 2019).*

*Proof.* Write:

$$\mathcal{L}_{\mathcal{D}_\alpha}(h) \leq \hat{\mathcal{L}}_{\mathcal{D}_\alpha}(h) + \sup_{f \in \mathcal{H}} \left( \mathcal{L}_{\mathcal{D}_\alpha}(f) - \hat{\mathcal{L}}_{\mathcal{D}_\alpha}(f) \right) \tag{12}$$

To link the second term to its expectation, we prove the following:

**Lemma 1.** *Define the function $\phi : (\mathcal{X} \times \mathcal{Y})^m \to \mathbb{R}$ by:*

$$\phi(\{x_{1,1}, y_{1,1}\}, \ldots, \{x_{N,m_N}, y_{N,m_N}\}) = \sup_{f \in \mathcal{H}} \left( \mathcal{L}_{\mathcal{D}_\alpha}(f) - \hat{\mathcal{L}}_{\mathcal{D}_\alpha}(f) \right).$$

*Denote for brevity $z_{i,j} = \{x_{i,j}, y_{i,j}\}$. Then, for any $i \in \{1, 2, \ldots, N\}, j \in \{1, 2, \ldots, m_i\}$:*

$$\sup_{z_{1,1}, \ldots, z_{N,m_N}, z'_{i,j}} |\phi(z_{1,1}, \ldots, z_{i,j}, \ldots, z_{N,m_N}) - \phi(z_{1,1}, \ldots, z'_{i,j}, \ldots, z_{N,m_N})| \leq \frac{\alpha_i}{m_i}\mathcal{M} \tag{13}$$

*Proof.* Fix any $i, j$ and any $z_{1,1}, \ldots, z_{N,m_N}, z'_{i,j}$. Denote the $\alpha$-weighted empirical average of the loss with respect to the sample $z_{1,1}, \ldots, z'_{i,j}, \ldots, z_{N,m_N}$ by $\hat{\mathcal{L}}'_{\mathcal{D}_\alpha}$. Then we have that:

$$
\begin{aligned}
|\phi(\ldots, z_{i,j}, \ldots) - \phi(\ldots, z'_{i,j}, \ldots)| &= |\sup_{f \in \mathcal{H}} \left( \mathcal{L}_{\mathcal{D}_\alpha}(f) - \hat{\mathcal{L}}_{\mathcal{D}_\alpha}(f) \right) - \sup_{f \in \mathcal{H}} (\mathcal{L}_{\mathcal{D}_\alpha}(f) - \hat{\mathcal{L}}'_{\mathcal{D}_\alpha}(f))| \\
&\leq |\sup_{f \in \mathcal{H}} (\hat{\mathcal{L}}'_{\mathcal{D}_\alpha}(f) - \hat{\mathcal{L}}_{\mathcal{D}_\alpha}(f))| \\
&= \frac{\alpha_i}{m_i} |\sup_{f \in \mathcal{H}} \left( \ell_f(z'_{i,j}) - \ell_f(z_{i,j}) \right)| \\
&\leq \frac{\alpha_i}{m_i} \mathcal{M}
\end{aligned}
$$

Note: the inequality we used above holds for bounded functions inside the supremum. $\qquad\square$

Let $S$ denote a random sample of size $m$ drawn from a distribution as the one generating out data (i.e. $m_i$ samples from $\mathcal{D}_i$ for each $i$). Now, using Lemma 1, McDiarmid's inequality gives:

$$
\begin{aligned}
\mathbb{P}\left( \phi(S) - \mathbb{E}(\phi(S)) \geq t \right) &\leq \exp\left( -\frac{2t^2}{\sum_{i=1}^N \sum_{j=1}^{m_i} \frac{\alpha_i^2}{m_i^2} \mathcal{M}^2} \right) \\
&= \exp\left( -\frac{2t^2}{\mathcal{M}^2 \sum_{i=1}^N \frac{\alpha_i^2}{m_i}} \right)
\end{aligned}
$$

For any $\delta > 0$, setting the right-hand side above to be $\delta/4$ and using (12), we obtain that with probability at least $1 - \delta/4$:

$$
\mathcal{L}_{\mathcal{D}_\alpha}(h) \leq \hat{\mathcal{L}}_{\mathcal{D}_\alpha}(h) + \mathbb{E}_S \left( \sup_{f \in \mathcal{H}} \left( \mathcal{L}_{\mathcal{D}_\alpha}(f) - \hat{\mathcal{L}}_{\mathcal{D}_\alpha}(f) \right) \right) + \sqrt{\frac{\log\left(\frac{4}{\delta}\right) \mathcal{M}^2}{2}} \sqrt{\sum_{i=1}^N \frac{\alpha_i^2}{m_i}} \tag{14}
$$

To deal with the expected loss inside the second term, introduce a ghost sample (denoted by $S'$), drawn from the same distributions as our original sample (denoted by $S$). Denoting the weighted empirical loss with respect to the ghost sample by $\hat{\mathcal{L}}'_{\mathcal{D}_\alpha}$, $\beta_i = m_i/m$ for all $i$, and using the convexity of the supremum, we obtain:

$$
\begin{aligned}
\mathbb{E}_S \left( \sup_{f \in \mathcal{H}} \left( \mathcal{L}_{\mathcal{D}_\alpha}(f) - \hat{\mathcal{L}}_{\mathcal{D}_\alpha}(f) \right) \right) &= \mathbb{E}_S \left( \sup_{f \in \mathcal{H}} \left( \mathbb{E}_{S'} \left( \hat{\mathcal{L}}'_{\mathcal{D}_\alpha}(f) \right) - \hat{\mathcal{L}}_{\mathcal{D}_\alpha}(f) \right) \right) \\
&\leq \mathbb{E}_{S,S'} \left( \sup_{f \in \mathcal{H}} \left( \hat{\mathcal{L}}'_{\mathcal{D}_\alpha}(f) - \hat{\mathcal{L}}_{\mathcal{D}_\alpha}(f) \right) \right) \\
&= \mathbb{E}_{S,S'} \left( \sup_{f \in \mathcal{H}} \left( \frac{1}{m} \sum_{i=1}^N \sum_{j=1}^{m_i} \frac{\alpha_i}{\beta_i} \left( \ell_f(z'_{i,j}) - \ell_f(z_{i,j}) \right) \right) \right)
\end{aligned}
$$

Introducing $m$ independent Rademacher random variables and noting that $(\ell_f(z') - \ell_f(z))$ and $\sigma(\ell_f(z') - \ell_f(z))$ have the same distribution, as long as $z$ and $z'$ have the same distribution:

$$
\begin{aligned}
\mathbb{E}_S \left( \sup_{f \in \mathcal{H}} \left( \mathcal{L}_{\mathcal{D}_\alpha}(f) - \hat{\mathcal{L}}_{\mathcal{D}_\alpha}(f) \right) \right) &\leq \mathbb{E}_{S,S',\sigma} \left( \sup_{f \in \mathcal{H}} \left( \frac{1}{m} \sum_{i=1}^N \sum_{j=1}^{m_i} \frac{\alpha_i}{\beta_i} \sigma_{i,j} \left( \ell_f(z_{i,j})' - \ell_f(z_{i,j}) \right) \right) \right) \\
&\leq \mathbb{E}_{S',\sigma} \left( \sup_{f \in \mathcal{H}} \left( \frac{1}{m} \sum_{i=1}^N \sum_{j=1}^{m_i} \frac{\alpha_i}{\beta_i} \sigma_{i,j} \ell_f(z_{i,j}) \right) \right) \\
&\quad + \mathbb{E}_{S,\sigma} \left( \sup_{f \in \mathcal{H}} \left( \frac{1}{m} \sum_{i=1}^N \sum_{j=1}^{m_i} \frac{\alpha_i}{\beta_i} (-\sigma_{i,j}) \ell_f(z_{i,j}) \right) \right) \\
&= 2\mathbb{E}_{S,\sigma} \left( \sup_{f \in \mathcal{H}} \left( \frac{1}{m} \sum_{i=1}^N \sum_{j=1}^{m_i} \frac{\alpha_i}{\beta_i} \sigma_{i,j} \ell_f(z_{i,j}) \right) \right).
\end{aligned}
$$

We can now link the last term to the empirical analog of the Rademacher complexity, by using the McDiarmid Inequality (with an observation similar to Lemma 1). Putting this together, we obtain that for any $\delta > 0$ with probability at least $1 - \delta/2$:

$$\mathcal{L}_{\mathcal{D}_{\boldsymbol{\alpha}}}(h) \leq \hat{\mathcal{L}}_{\mathcal{D}_{\boldsymbol{\alpha}}}(h) + 2\mathbb{E}_{\sigma}\left(\sup_{f \in \mathcal{H}}\left(\frac{1}{m}\sum_{i=1}^{N}\sum_{j=1}^{m_i}\frac{\alpha_i}{\beta_i}\sigma_{i,j}\ell_f(z_{i,j})\right)\right) + 3\sqrt{\frac{\log\left(\frac{4}{\delta}\right)M^2}{2}}\sqrt{\sum_{i=1}^{N}\frac{\alpha_i^2}{m_i}} \tag{15}$$

Finally, note that:

$$\mathbb{E}_{\sigma}\left(\sup_{f \in \mathcal{H}}\left(\frac{1}{m}\sum_{i=1}^{N}\sum_{j=1}^{m_i}\frac{\alpha_i}{\beta_i}\sigma_{i,j}\ell_f(z_{i,j})\right)\right) \leq \mathbb{E}_{\sigma}\left(\sum_{i=1}^{N}\alpha_i\sup_{f \in \mathcal{H}}\left(\frac{1}{m_i}\sum_{j=1}^{m_i}\sigma_{i,j}\ell_f(z_{i,j})\right)\right)$$

$$= \sum_{i=1}^{N}\alpha_i\mathbb{E}_{\sigma}\left(\sup_{f \in \mathcal{H}}\left(\frac{1}{m_i}\sum_{j=1}^{m_i}\sigma_{i,j}\ell_f(z_{i,j})\right)\right)$$

$$= \sum_{i=1}^{N}\alpha_i\mathcal{R}_i(\mathcal{H})$$

Bounding $\hat{\mathcal{L}}_{\mathcal{D}_{\boldsymbol{\alpha}}}(h) - \mathcal{L}_{\mathcal{D}_{\boldsymbol{\alpha}}}(h)$ with the same quantity and with probability at least $1 - \delta/2$ follows by a similar argument. The result then follows by applying the union bound. □

## B    Proof of Theorem 2

**Theorem 2.** *For any $\boldsymbol{w}$, when $\lambda > 0$ and $\{\hat{\mathcal{L}}_i(\boldsymbol{w})\}_{i=1}^{N}$ are sorted in increasing order: $\hat{\mathcal{L}}_1(\boldsymbol{w}) \leq \hat{\mathcal{L}}_2(\boldsymbol{w}) \leq ... \leq \hat{\mathcal{L}}_N(\boldsymbol{w})$, by setting:*

$$p = \underset{k}{\operatorname{argmax}}\{1 + \frac{M_k(\overline{\mathcal{L}}_k(\boldsymbol{w}) - \hat{\mathcal{L}}_k(\boldsymbol{w}))}{\lambda} > 0\}, \tag{7}$$

*where $M_k = \sum_{i=1}^{k}m_i$,*

$$\overline{\mathcal{L}}_k(h) = \frac{\sum_{i=1}^{k}m_i\hat{\mathcal{L}}_i(\boldsymbol{w})}{M_k} \tag{8}$$

*is the average loss over the first $k$ clients that have the smallest empirical risks. Then the optimal $\boldsymbol{\alpha}$ to the problem (6) is given by:*

$$\alpha_i(\boldsymbol{w}) = \frac{m_i}{M_p}[1 + \frac{M_p(\overline{\mathcal{L}}_p(\boldsymbol{w}) - \hat{\mathcal{L}}_i(\boldsymbol{w}))}{\lambda}]_+, \tag{9}$$

*where $[\cdot]_+ = max(0, \cdot)$.*

*Proof.* The Lagrangian function of Eq. (6) is

$$\mathbb{L} = \boldsymbol{\alpha}^{\top}\hat{\mathcal{L}}(\boldsymbol{w}) + \frac{\lambda}{2}||\boldsymbol{\alpha}^{\top}\boldsymbol{m}^{\circ-\frac{1}{2}}||_2^2 - \boldsymbol{\alpha}^{\top}\boldsymbol{\beta} - \eta(\boldsymbol{\alpha}^{\top}\boldsymbol{1} - 1), \tag{16}$$

where $\hat{\mathcal{L}}(\boldsymbol{w}) = [\hat{\mathcal{L}}_1(\boldsymbol{w}), \hat{\mathcal{L}}_2(\boldsymbol{w}), ..., \hat{\mathcal{L}}_N(\boldsymbol{w})]^{\top}$, $\circ$ is the Hadamard root operation, $\boldsymbol{\beta}$ and $\eta$ are the Lagrangian multipliers. Then the following Karush-Kuhn-Tucker (KKT) conditions hold:

$$\partial_{\boldsymbol{\alpha}}\mathbb{L}(\boldsymbol{\alpha}, \boldsymbol{\beta}, \eta) = 0, \tag{17}$$
$$\boldsymbol{\alpha}^{\mathsf{T}}\boldsymbol{1} - 1 = 0, \tag{18}$$
$$\boldsymbol{\alpha} \geq 0, \tag{19}$$
$$\boldsymbol{\beta} \geq 0, \tag{20}$$
$$\alpha_i\beta_i = 0, \forall i = 1, 2, ...N. \tag{21}$$

According to Eq. (17), we have:

$$\alpha_i = \frac{m_i(\beta_i + \eta - \hat{\mathcal{L}}_i(\boldsymbol{w}))}{\lambda}. \tag{22}$$

Since $\beta_i \geq 0$, we discuss the following cases:

1. When $\beta_i = 0$, we have $\alpha_i = \frac{m_i(\eta - \hat{\mathcal{L}}_i(\boldsymbol{w}))}{\lambda} \geq 0$. Note that we further have $\eta - \hat{\mathcal{L}}_i(\boldsymbol{w}) \geq 0$.

2. When $\beta_i > 0$, from the condition $\alpha_i \beta_i = 0$, we have $\alpha_i = 0$.

Therefore, the optimal solution to Eq. (6) is given by:

$$\alpha_i(\boldsymbol{w}) = [\frac{m_i(\eta - \hat{\mathcal{L}}_i(\boldsymbol{w}))}{\lambda}]_+, \tag{23}$$

where $[\cdot]_+ = max(0, \cdot)$.

We notice that $\sum_{i=1}^{p} \alpha_i = 1$, thus we can get:

$$\eta = \frac{\sum_{i=1}^{p} m_i \hat{\mathcal{L}}_i(\boldsymbol{w}) + \lambda}{\sum_{i=1}^{p} m_i}. \tag{24}$$

According to $\eta - \hat{\mathcal{L}}_i(\boldsymbol{w}) \geq 0$, we have Eq. (7) and Eq. (8). Finally, plugging Eq. (24) into Eq. (23) yields Eq. (9).  □