

# Robust Aggregation for Federated Learning by Minimum $\gamma$ -Divergence Estimation

Cen-Jhih Li<sup>1</sup>, Pin-Han Huang<sup>2</sup>, Yi-Ting Ma<sup>1</sup>, Hung Hung<sup>2</sup>

and Su-Yun Huang<sup>1\*</sup>

<sup>1</sup>Academia Sinica, Taiwan

<sup>2</sup>National Taiwan University

April 1, 2022

## Abstract

Federated learning is a framework for multiple devices or institutions, called local clients, to collaboratively train a global model without sharing their data. For federated learning with a central server, an aggregation algorithm integrates model information sent from local clients to update the parameters for a global model. Sample mean is the simplest and most commonly used aggregation method. However, it is not robust for data with outliers or under the Byzantine problem, where Byzantine clients send malicious messages to interfere with the learning process. Some robust aggregation methods were introduced in literature including marginal median, geometric median and trimmed-mean. In this article, we propose an alternative robust aggregation method, named  $\gamma$ -mean, which is the minimum divergence estimation based on a robust density power divergence. This  $\gamma$ -mean aggregation mitigates the influence of Byzantine clients by assigning less weights. This weighting scheme is data-driven and controlled by the  $\gamma$  value. Robustness from the viewpoint of influence function is discussed and some numerical results are presented.

**Keywords:** Byzantine problem, density power divergence, federated learning, gamma-divergence, influence function, robustness.

---

\*Corresponding author: syhuang@stat.sinica.edu.tw

# 1 Introduction

Federated learning (FL), a distributed learning paradigm, was proposed by Google in 2016 [12] for distributing the optimization of model parameters over multiple machines. The distributed framework allows the data to be stored in local devices without sharing with each other. These machines, called clients, collaboratively train a global model. In this article, we consider the case that there is a central server, which works as a coordinator and aggregates the gradient information sent by the local clients to update the global model through an iterative process [16]. In real-world applications, clients, such as hospitals and clinical institutions, hold highly sensitive and private data such as electronic healthcare records, medical images, etc. [21]. However, FL still encounters several challenges in real world setting. In its training process, some unreliable clients may produce outlying data information, send malicious values, or even introduce an adversarial attack [1], which leads to a biased result or even a failure of the training process. These arbitrary clients are called Byzantine clients. With a properly designed aggregation scheme, the effect by Byzantine clients can be reduced or excluded. Several papers have made significant contribution to developing Byzantine-resilience methods with specific use of robust aggregators, including marginal median [3], geometric median [4], and trimmed-mean [20]. Another aggregation framework is Byzantine-robust stochastic aggregation method [13], which adds a regularization term in the objective function to handle the Byzantine and heterogeneous data problem.

In this article, we consider the case of FL, which consists of a central server and  $m$  clients. Assume that there is an  $\alpha$ -fraction of Byzantine clients. Our goal is to minimize the following objective function:

$$L(\theta) = \sum_{i=1}^m \pi_i \mathbb{E}_{\eta_i \sim \mathcal{D}} [l_{\eta_i}(\theta)], \quad (1)$$

where  $\theta \in \mathbb{R}^p$  is the global model parameter,  $l_{\eta_i}(\theta)$  is the loss function of client  $i$  with local data  $\eta_i$  having an unknown distribution  $\mathcal{D}$ , and  $\{\pi_i\}_{i=1}^m$  are the prior weights assigned to local clients. At the  $t^{\text{th}}$  round of FL iterative updates, the server broadcasts the current parameter value  $\theta^t \in \mathbb{R}^p$  to local clients. Normal clients faithfully compute an estimate of the local gradients  $\mathbb{E}_{\eta_i \sim \mathcal{D}} \nabla l_{\eta_i}(\theta^t)$  and send the gradient information back to the server. On the contrary, Byzantine clients send arbitrary erroneous messages to obstruct the optimization process. The central server aggregates the gradient information from clients and updates the model parameter from  $\theta^t$  to  $\theta^{t+1}$ . At the end, the server should output an estimate of the optimal  $\theta$ , which aims to minimize the objective loss (1). In this work we consider uniform client weights, i.e.,  $\pi_1 = \dots = \pi_m = 1$ . The main reason is that, under the Byzantine

problem, the local sample sizes claimed by Byzantine clients might not be reliable.

McMahan et al. introduced the *FedAvg* algorithm [14], which used the mean as its aggregator. The sample mean aggregation is a very popular FL framework [14, 5]. However, the sample mean is known to be vulnerable to outliers and heavy-tailed distribution. For instance, Byzantine clients [13, 20, 15] may send extreme values to strongly influence the aggregation of sample mean. There are some robust alternatives to the mean aggregator, such as marginal median, geometric median and trimmed mean. In contrast to mean, marginal median is a relatively robust aggregator by computing coordinate-wise medians. Another Byzantine-resilient aggregator is geometric median, which computes the geometric median of  $\{\mathbf{x}_i\}_{i=1}^m$  given by (see Weiszfeld’s algorithm [18] for solution computation)

$$\text{GeoMed}(\{\mathbf{x}_i\}_{i=1}^m) = \arg \min_{\mathbf{x} \in \mathbb{R}^p} \sum_{i=1}^m \|\mathbf{x} - \mathbf{x}_i\|_2. \quad (2)$$

The trimmed mean aggregator computes coordinate-wise trimmed means by removing the  $\beta$ -fraction of data points from each of the two tails [15]. Both the marginal median and geometric median, though are robust to outliers, but are not efficient, as they only use the most centrally-lying data for inference. As for trimmed mean, first it needs a proper selection of  $\beta$ . Second, it often does not work well for outliers lying on one side of a coordinate instead of on two sides. To have a better balance between robustness and efficiency, we will propose another robust aggregator based on a robust density power divergence, namely, the  $\gamma$ -divergence.

The rest of the article is organized as follows. In Section 2, we propose our main method for robust FL aggregation and provide some theoretical viewpoints by influence function. In Section 3, we conduct an extensive simulation study. In Section 4, we provide some further numerical examples using MNIST, fashion MNIST and chest x-ray images (pneumonia). In Section 5, we make some concluding remarks.

## 2 Proposed aggregator and its robustness

In this section, we introduce an aggregator, called “ $\gamma$ -mean”, which is based on the minimum  $\gamma$ -divergence estimation. Two versions of implementation algorithms are given and some robustness viewpoints based on influence function are provided.

## 2.1 Minimum $\gamma$ -divergence estimation

The  $\gamma$ -divergence [7, 9], which is also known as the density power divergence of type zero [10], is a robust divergence against outliers with a tuning parameter  $\gamma > 0$ . The  $\gamma$ -divergence between the data distribution with probability density function  $g$  and the model distribution  $f_\tau$  (indexed by  $\tau$ ) is defined by

$$D_\gamma(g, f_\tau) = \frac{\|g\|_{\gamma+1}}{\gamma(\gamma+1)} \left\{ 1 - \int \left( \frac{f_\tau(\mathbf{x})}{\|f_\tau\|_{\gamma+1}} \right)^\gamma \frac{g(\mathbf{x})}{\|g\|_{\gamma+1}} d\mathbf{x} \right\},$$

where  $\|f_\tau\|_{\gamma+1} = (\int f_\tau^{\gamma+1}(\mathbf{x}) d\mathbf{x})^{1/(\gamma+1)}$ . In the limiting case, it reduces to the Kullback-Leibler divergence, i.e.  $\lim_{\gamma \rightarrow 0} D_\gamma(g, f_\tau) = \int \ln\left(\frac{g(\mathbf{x})}{f_\tau(\mathbf{x})}\right) g(\mathbf{x}) d\mathbf{x}$ . In the population level, the minimum  $\gamma$ -divergence estimation of  $\tau$  is given by

$$\tau_\gamma^* = \arg \min_{\tau} D_\gamma(g, f_\tau) = \arg \max_{\tau} \int \left( \frac{f_\tau(\mathbf{x})}{\|f_\tau\|_{\gamma+1}} \right)^\gamma g(\mathbf{x}) d\mathbf{x}. \quad (3)$$

In the sample level with empirical data  $\{\mathbf{x}_i\}_{i=1}^m$ , the data distribution will be replaced by the empirical distribution to get the estimate  $\hat{\tau}_\gamma$ :

$$\hat{\tau}_\gamma = \arg \min_{\tau} D_\gamma(g, f_\tau) = \arg \max_{\tau} \frac{1}{m} \sum_{i=1}^m \left( \frac{f_\tau(\mathbf{x}_i)}{\|f_\tau\|_{\gamma+1}} \right)^\gamma.$$

In this article, we use a Gaussian working model, i.e.,  $f_\tau \sim N(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ . Suppose  $g$  takes a contaminated form:  $g(\mathbf{x}) = (1 - \epsilon)f_\tau(\mathbf{x}) + \epsilon h(\mathbf{x})$ , where  $h(\mathbf{x})$  is the probability density function of contamination distribution. It is assumed that  $\int f_\tau^\gamma(\mathbf{x}) h(\mathbf{x})$  is small for certain  $\gamma > 0$ , so that the contamination has little effect in the learning process. By taking derivatives with respect to  $\tau = (\boldsymbol{\mu}, \boldsymbol{\Sigma})$  and setting them to zero, the estimates  $(\hat{\boldsymbol{\mu}}, \hat{\boldsymbol{\Sigma}})$  have to satisfy the following estimating equations:

$$\begin{aligned} \frac{1}{m} \sum_{i=1}^m d_i^{(\gamma)}(\hat{\boldsymbol{\mu}}, \hat{\boldsymbol{\Sigma}})(\mathbf{x}_i - \hat{\boldsymbol{\mu}}) &= 0, \\ \frac{1}{m} \sum_{i=1}^m d_i^{(\gamma)}(\hat{\boldsymbol{\mu}}, \hat{\boldsymbol{\Sigma}}) \left\{ \frac{1}{\gamma+1} \hat{\boldsymbol{\Sigma}} - (\mathbf{x}_i - \hat{\boldsymbol{\mu}})(\mathbf{x}_i - \hat{\boldsymbol{\mu}})^\top \right\} &= 0, \end{aligned}$$

where  $d_i^{(\gamma)}(\hat{\boldsymbol{\mu}}, \hat{\boldsymbol{\Sigma}}) = \exp\left(\frac{-\gamma}{2}(\mathbf{x}_i - \hat{\boldsymbol{\mu}})^\top \hat{\boldsymbol{\Sigma}}^{-1}(\mathbf{x}_i - \hat{\boldsymbol{\mu}})\right)$ . The solution pair,  $(\hat{\boldsymbol{\mu}}, \hat{\boldsymbol{\Sigma}})$ , have to satisfy the following stationary equations:

$$\hat{\boldsymbol{\mu}} = \frac{\sum_{i=1}^m d_i^{(\gamma)}(\hat{\boldsymbol{\mu}}, \hat{\boldsymbol{\Sigma}}) \mathbf{x}_i}{\sum_{i=1}^m d_i^{(\gamma)}(\hat{\boldsymbol{\mu}}, \hat{\boldsymbol{\Sigma}})}, \quad (4)$$

$$\hat{\boldsymbol{\Sigma}} = (\gamma+1) \frac{\sum_{i=1}^m d_i^{(\gamma)}(\hat{\boldsymbol{\mu}}, \hat{\boldsymbol{\Sigma}}) (\mathbf{x}_i - \hat{\boldsymbol{\mu}})(\mathbf{x}_i - \hat{\boldsymbol{\mu}})^\top}{\sum_{i=1}^m d_i^{(\gamma)}(\hat{\boldsymbol{\mu}}, \hat{\boldsymbol{\Sigma}})}. \quad (5)$$

## 2.2 Robust aggregation by $\gamma$ -mean

In view of the stationary equations (4)-(5), we use a fixed-point iteration. Two algorithms are provided below. One is the usual  $\gamma$ -mean, which uses  $N(\boldsymbol{\mu}, \boldsymbol{\Sigma})$  as the working model, and the other is the simple  $\gamma$ -mean, which adopts  $N(\boldsymbol{\mu}, \sigma^2 \mathbf{I})$  as the working model. Note that, in the former case when the sample size is not sufficiently large enough to have a stable estimate of the covariance inverse, we will use  $\{\text{diag}(\hat{\boldsymbol{\Sigma}})\}^{-1}$  instead, where  $\hat{\boldsymbol{\Sigma}}$  is the minimum  $\gamma$ -divergence estimator for the covariance and  $\text{diag}(\hat{\boldsymbol{\Sigma}})$  denotes the diagonal matrix consisting of diagonal elements of  $\hat{\boldsymbol{\Sigma}}$ . Also note that, in the latter case of simple  $\gamma$ -mean,  $\sigma^2$  can be merged into the hyperparameter  $\gamma$ . Thus, we will simply use the standard Gaussian as the working model.

---

**Algorithm 1**  $\gamma$ -mean with a Gaussian working model

---

**Input:** Gradient information  $\mathbf{X} = [\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_m]$  and the maximum number of iterations  $S$

**Output:**  $\boldsymbol{\mu}_\gamma = \hat{\boldsymbol{\mu}}$

Start with initials  $\hat{\boldsymbol{\mu}} = \hat{\boldsymbol{\mu}}_{ini}$  and  $\hat{\boldsymbol{\Sigma}} = \hat{\boldsymbol{\Sigma}}_{ini}$ .

**for**  $s = 1, 2, \dots, S$  (while  $s \leq S$  and iterations not yet converge) **do**

**for**  $i = 1, 2, \dots, m$  **do**

        Calculate  $d_i^{(\gamma)}(\hat{\boldsymbol{\mu}}, \hat{\boldsymbol{\Sigma}}) \leftarrow \exp\left(-\frac{\gamma}{2}(\mathbf{x}_i - \hat{\boldsymbol{\mu}})^\top \hat{\boldsymbol{\Sigma}}^{-1}(\mathbf{x}_i - \hat{\boldsymbol{\mu}})\right)$  at the  $i^{\text{th}}$  local client.

**end for**

    Denote  $w_i^{(\gamma)} = \frac{d_i^{(\gamma)}(\hat{\boldsymbol{\mu}}, \hat{\boldsymbol{\Sigma}})}{\sum_{i=1}^m d_i^{(\gamma)}(\hat{\boldsymbol{\mu}}, \hat{\boldsymbol{\Sigma}})}$ .

$\hat{\boldsymbol{\mu}} \leftarrow \sum_{i=1}^m w_i^{(\gamma)} \cdot \mathbf{x}_i$ ,

$\hat{\boldsymbol{\Sigma}} \leftarrow (1 + \gamma) \cdot \sum_{i=1}^m w_i^{(\gamma)} \cdot (\mathbf{x}_i - \hat{\boldsymbol{\mu}})(\mathbf{x}_i - \hat{\boldsymbol{\mu}})^\top$ .

**end for**

---

For extremely large dimension  $p$ , such as in a deep neural network model, the simple  $\gamma$ -mean algorithm will be more feasible than the usual  $\gamma$ -mean for better numerical stability.

## 2.3 Robustness

### 2.3.1 Influence function

The influence function is a tool to evaluate the change of an estimator by a small perturbation to the data distribution. An estimator with a lower influence provides better resistance against outliers. Therefore, we analyze the robustness of different estimators by showing their influence functions. Let  $\{\mathbf{x}_i\}_{i=1}^m \subset \mathbb{R}^p$  be sampled from  $G$  and  $T$  be the statistical functional for estimation. The robustness of  $T(G)$  can be evaluated by the influence function

---

**Algorithm 2** simple  $\gamma$ -mean with the standard Gaussian as the working model

---

**Input:** Gradient information  $\mathbf{X} = [\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_m]$  and the maximum number of iterations  $S$

**Output:**  $\boldsymbol{\mu}_\gamma = \hat{\boldsymbol{\mu}}$

Start with initial  $\hat{\boldsymbol{\mu}} = \hat{\boldsymbol{\mu}}_{ini}$ .

**for**  $s = 1, 2, \dots, S$  (while  $s \leq S$  and iterations not yet converge) **do**

**for**  $i = 1, 2, \dots, m$  **do**

        Calculate  $d_i^{(\gamma)}(\hat{\boldsymbol{\mu}}) \leftarrow \exp\left(-\frac{\gamma}{2}(\mathbf{x}_i - \hat{\boldsymbol{\mu}})^\top(\mathbf{x}_i - \hat{\boldsymbol{\mu}})\right)$  at the  $i^{\text{th}}$  local client.

**end for**

    Denote  $w_i^{(\gamma)} = \frac{d_i^{(\gamma)}(\hat{\boldsymbol{\mu}})}{\sum_{i=1}^m d_i^{(\gamma)}(\hat{\boldsymbol{\mu}})}$ .

$\hat{\boldsymbol{\mu}} \leftarrow \sum_{i=1}^m w_i^{(\gamma)} \cdot \mathbf{x}_i$ .

**end for**

---

of  $T$ :

$$\text{IF}_T(\mathbf{x}_0; G) = \left. \frac{\partial}{\partial \epsilon} T\{(1 - \epsilon)G + \epsilon \delta_{\mathbf{x}_0}\} \right|_{\epsilon=0},$$

where  $\delta_{\mathbf{x}_0}$  is the point mass at  $\mathbf{x}_0$ . In other words, the influence function of  $T$  is the Gâteaux derivative of  $T(G)$  with respect to  $G$  along the direction  $\delta_{\mathbf{x}_0} - G$ . In this paper, we consider to evaluate the influence function at the Gaussian working model  $F_\tau$  with  $\tau = (\boldsymbol{\mu}, \boldsymbol{\Sigma})$  and to see the deviant effect when a point perturbation  $\mathbf{x}_0$  is added. The influence function of M-estimators have been well-studied in [8], which is related to the inverse of Hessian matrix and the first order derivative. The estimator based on  $\gamma$ -divergence is also an M-estimate. Its influence function derivation can be found in [9] and is given by

$$\text{IF}_{\boldsymbol{\mu}_\gamma}(\mathbf{x}_0; G)|_{G=F_\tau} = (\gamma + 1)^{\frac{p+2}{2}} \exp\left(\frac{-\gamma}{2}(\mathbf{x}_0 - \boldsymbol{\mu}_\gamma)^\top \boldsymbol{\Sigma}_\gamma^{-1}(\mathbf{x}_0 - \boldsymbol{\mu}_\gamma)\right)(\mathbf{x}_0 - \boldsymbol{\mu}_\gamma), \quad (6)$$

where  $\boldsymbol{\mu}_\gamma$  and  $\boldsymbol{\Sigma}_\gamma^{-1}$  are the minimum  $\gamma$ -divergence estimators for the location and covariance matrix of the working model  $F_\tau$ , respectively. The influence function of the  $\gamma$ -mean,  $\text{IF}_{\boldsymbol{\mu}_\gamma}(\mathbf{x}_0; F_\tau)$ , is similar to the influence function of the sample mean, which is given by  $(\mathbf{x}_0 - \boldsymbol{\mu})$ , except for an additional multiplicative weight  $d_{\mathbf{x}_0}^{(\gamma)}(\boldsymbol{\mu}_\gamma, \boldsymbol{\Sigma}_\gamma)$ . This weight focuses on the Mahalanobis distance between the perturbing point  $\mathbf{x}_0$  and  $\boldsymbol{\mu}_\gamma$ . As a result, when  $\mathbf{x}_0$  is an outlier away from the target mean of data distribution, this weighting factor produces a down-weighting effect and prevents the estimator from biased estimation caused by the outlier.

### 2.3.2 Comparison with other aggregators

In this subsection, we present influence functions for several robust aggregation methods including marginal median, geometric median and trimmed-mean.

**Geometric median  $\mu_{\text{geo}}$ :** The influence function of  $\mu_{\text{geo}}$  is given by [2]

$$\text{IF}_{\mu_{\text{geo}}}(\mathbf{x}_0; G) = \left( \mathbb{E}_X \left[ \|X - \mu_{\text{geo}}\|_2^{-1} \{I_p - \|X - \mu_{\text{geo}}\|_2^{-2} (X - \mu_{\text{geo}})(X - \mu_{\text{geo}})^\top\} \right] \right)^{-1} \frac{\mathbf{x}_0 - \mu_{\text{geo}}}{\|\mathbf{x}_0 - \mu_{\text{geo}}\|_2},$$

where the expectation  $E_X$  is taken with respect to the data distribution  $G$ .

**Trimmed-mean  $\mu_{\text{trm}}$ :** The trimmed-mean is an  $L$ -esimator, and its influence function is derived in [17]. For the influence function of the marginal trimmed-mean, we derive it coordinate-wise since the coordinates are independent to each other. Therefore, let  $G_i^{-1}$  be the quantile function of marginal data distribution on  $i$ -th component, and the influence function of  $\mu_{\text{trm}}$  is given by

$$\text{IF}_{\mu_{\text{trm}}}(\mathbf{x}_0; G) = \frac{1}{1 - 2\beta} \begin{pmatrix} w_1 \\ \vdots \\ w_p \end{pmatrix},$$

where

$$w_i = \begin{cases} G_i^{-1}(\beta) - W_i, & x_{0i} < G_i^{-1}(\beta), \\ \mathbf{x}_0 - W_i, & G_i^{-1}(\beta) < x_{0i} < G_i^{-1}(1 - \beta), \\ G_i^{-1}(1 - \beta) - W_i, & x_{0i} > G_i^{-1}(1 - \beta), \end{cases}$$

$\mathbf{x}_0 = (x_{01}, \dots, x_{0p})^\top$ , and

$$W_i = (1 - 2\beta) \int_{G_i^{-1}(\beta)}^{G_i^{-1}(1-\beta)} t dG_i(t) + \beta(G_i^{-1}(1 - \beta) + G_i^{-1}(\beta)).$$

It is obvious that an influential effect,  $\mathbf{x}_0 - W_i$ , still remains if  $\mathbf{x}_0$  is outside the trimmed range.

**Marginal median  $\mu^{(0.5)}$ :** The influence function of median is also derived in [17], and we use the same techniques as used in the trimmed-mean to derive the influence function of the marginal median. The influence function is given by

$$\text{IF}_{\mu^{(0.5)}}(\mathbf{x}_0; G) = \begin{pmatrix} \frac{x_{01} - \mu_1^{(0.5)}}{2g_1(\mu_1^{(0.5)})|x_{01} - \mu_1^{(0.5)}|} \\ \vdots \\ \frac{x_{0p} - \mu_p^{(0.5)}}{2g_p(\mu_p^{(0.5)})|x_{0p} - \mu_p^{(0.5)}|} \end{pmatrix}, \quad (7)$$

where  $\boldsymbol{\mu}^{(0.5)} = (\mu_1^{(0.5)}, \dots, \mu_p^{(0.5)})^\top$  is the marginal median and  $g_i$  is the density function of the marginal distribution  $G_i$ .

The influence functions given above reveal that different robust aggregators have different resistance ability against outliers. Trimmed mean can still be influenced by outliers outside the trimmed range. While marginal median and geometric median are fairly robust, they may lose too much efficiency due to the use of only the most centrally-lying data point. The  $\gamma$ -mean provides an adjustable control between the robustness and efficiency by the hyper-parameter  $\gamma$ .

### 3 Simulation study

In this section, we evaluate our robust  $\gamma$ -mean aggregator and compare it with existing aggregators, including mean, marginal median, geometric median and trimmed mean, through simulation. For computing geometric median, Weiszfeld's algorithm [18] is used. In this simulation study, we investigate the behavior and inspect the robustness of  $\gamma$ -mean and other aggregators as an estimator of the location parameter under multivariate Gaussian distribution and multivariate  $t$ -distribution with and without Byzantine interference. We build several simulation scenarios, including testing aggregators across the growth of dimension  $p$ , different fractions (i.e.  $\alpha$  values) of Byzantine attacks and for various settings of  $\gamma$  (which controls the degree of robustness).

#### 3.1 Simulation settings

- **Scenario 1.** We focus on the behavior of aggregators for increasing  $p$  from 20 to 1000. Other experimental setting is as follows: the number of clients  $m = 200$ , the fraction of Byzantine attacks  $\alpha = 0$  and  $\alpha = 0.1$ , and the hyper-parameter for controlling the robustness  $\gamma = 2/p$ .
- **Scenario 2.** We focus on the behavior of aggregators for increasing contamination fraction  $\alpha$  from 0 to 0.5. Other experimental setting is as follows:  $m = 200$ ,  $p = 1000$  and  $\gamma = 2/p$ .
- **Scenario 3.** We focus on the effect of  $\gamma$  values and set  $\gamma = c/p$  for various constants  $c$  ranging from 0.5 to 4. Other experimental setting is as follows:  $m = 200$ ,  $\alpha = 0.1$  and  $p$  ranges from 1 to 1000.



- **Scenario 4.** After comparison between the  $\gamma$ -mean and other aggregators, we focus on the comparison between two versions of our proposal, the  $\gamma$ -mean versus the simple  $\gamma$ -mean. Other experimental setting is as follows:  $m = 200$ ,  $\alpha = 0.1$ ,  $\gamma = 2/p$  and  $p$  ranges from 1 to 1000.

For regular clients, gradient vectors are generated from the standard Gaussian distribution and  $t$ -distribution with 5 degrees of freedom. For Byzantine clients, gradient vectors are generated from the same Gaussian and  $t$  distributions, except with a location shift  $\boldsymbol{\mu} = 100 \times \mathbf{1}_p$ , where  $\mathbf{1}_p$  is a vector with one in all entries. Each experimental scenario is implemented with 100 replicate runs.

## 3.2 Results

To compare the performance of different aggregators, the mean squared error (MSE) to the true target value is used as a performance metric. MSE can be further decomposed into the squared bias and variance.

### 3.2.1 Scenario 1

The results are shown in Figure 1. Without contamination (Figure 1a), all aggregators have MSE close to zero. For the Gaussian case, the MSE curves of the mean, geometric median and  $\gamma$ -mean are almost collapsed together. These 3 curves have the lowest MSE values followed by the trimmed-mean and marginal median. For the  $t$ -distribution, the  $\gamma$ -mean has the lowest MSE, followed by the geometric median, trimmed-mean, mean and then marginal median. With 10% of contamination, the mean-aggregator fails to perform well and has large values of MSE, verifying that it is not a robust aggregator. Due its large MSE value, the result of mean-aggregator is removed from Figure 1b. From Figure 1b we can see that  $\gamma$ -mean performs the best among the five aggregators in terms of MSE. In addition,  $\gamma$ -mean remains low and stable in variance under both Gaussian and  $t$  distributions.

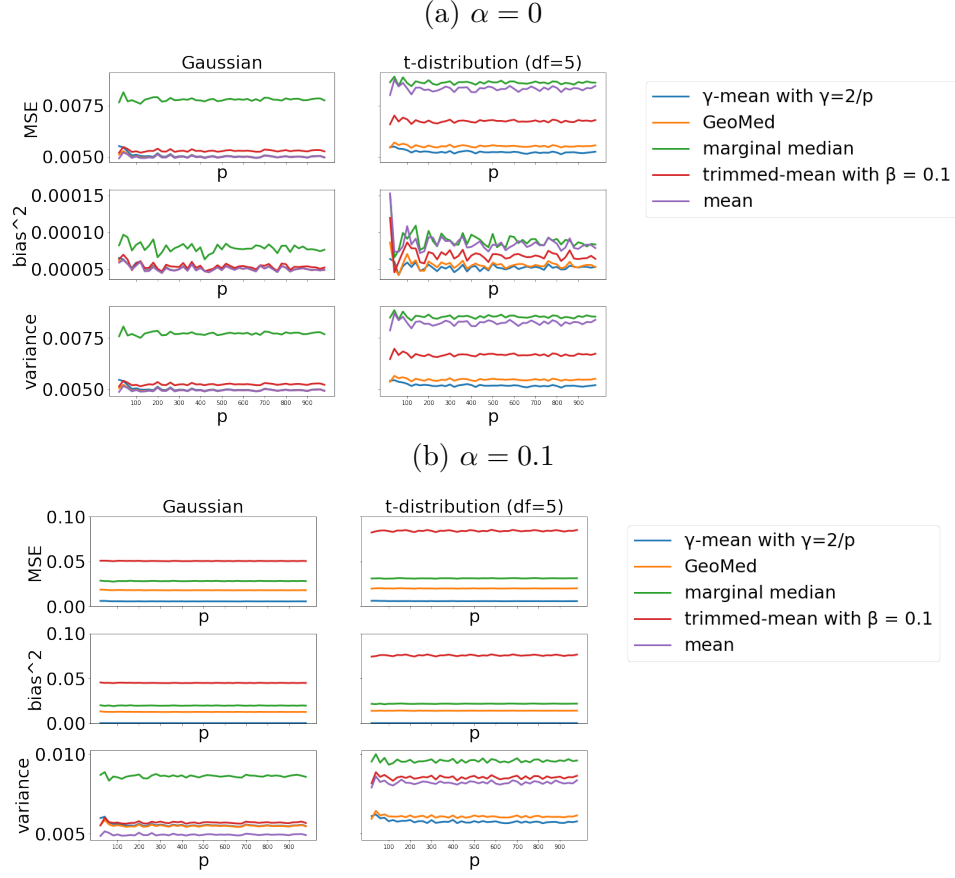


Figure 1: Comparison of different aggregators across different dimensions  $p$ . (a) Case  $\alpha = 0$  (no Byzantine client). (b) Case  $\alpha = 0.1$  (10% Byzantine clients).

### 3.2.2 Scenario 2

As Figure 2 shows, values of squared bias are close to zero when  $\alpha = 0$  and increase as  $\alpha$  grows for all aggregators. However, the increasing velocity of aggregators are quite different. Mean-aggregator surges with the highest velocity, while  $\gamma$ -mean progresses with the lowest velocity. Trimmed-mean slowly grows when  $\alpha \leq 0.1$ ; yet, its increasing pace is about the same as the mean after  $\alpha > 0.1$  due to the fixed 10% trimming percentage from both tails. The trends in MSE for different aggregators have similar patterns, and the only difference is in the scale with the  $\gamma$ -mean having the lowest MSE values. Also note that the variance of marginal median increases much faster than other aggregators as  $\alpha$  increases, followed by geometric median,  $\gamma$ -mean, trimmed mean and then mean. However, the trimmed mean and the mean have dominant bias leading to high MSE, even though they have small variance. As the values of squared bias are large for the mean and trimmed-mean, the bias curves for the  $\gamma$ -mean, geometric median and marginal median in the 2nd row of Figure 2 look collapsed

together. To have a better view of these three bias curves, zoomed-in views are provided in the 3rd row of Figure 2, and we can clearly see that  $\gamma$ -mean has the lowest bias values.

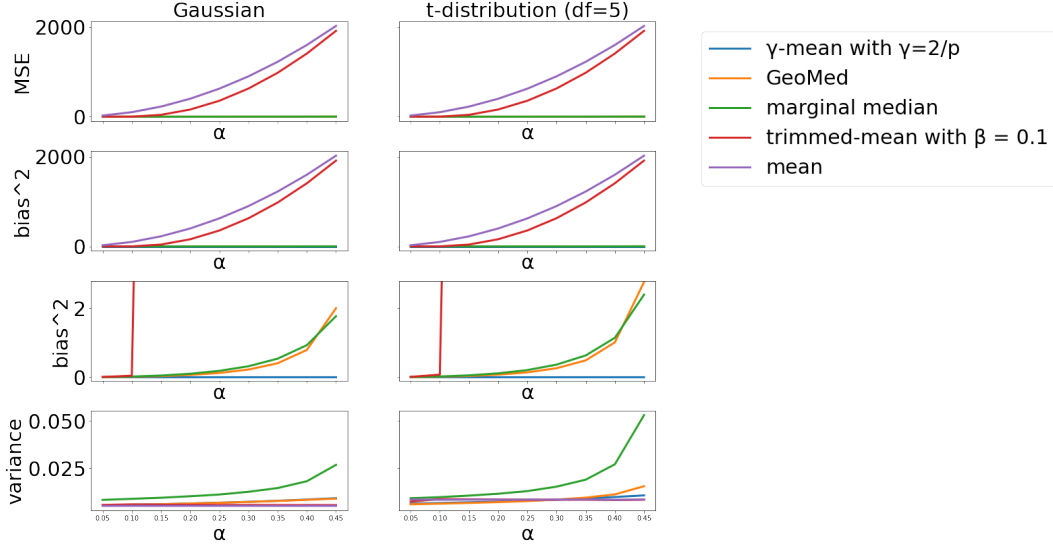


Figure 2: Comparison of different aggregators across different  $\alpha$  values with fixed  $p = 1000$ . The two subplots in the 3rd row are the zoomed-in views of subplots in the 2nd row.

### 3.2.3 Scenario 3

We conduct further experiments on  $\gamma$ -mean to see the effect of different  $\gamma$  values by setting  $\gamma = c/p$  with  $c = 0.5, 1, 2, 3, 4$ , where  $\alpha = 0.1$  and  $p \in [1, 1000]$ . Results are presented in Figure 3. The squared bias is ignorable, indicating that the  $\gamma$ -mean is quite robust with different  $\gamma$ -values. The main source of MSE comes from the variance. Larger  $\gamma$  values lead to larger stochastic variances, indicating that larger  $\gamma$  values result in lower estimation efficiency.

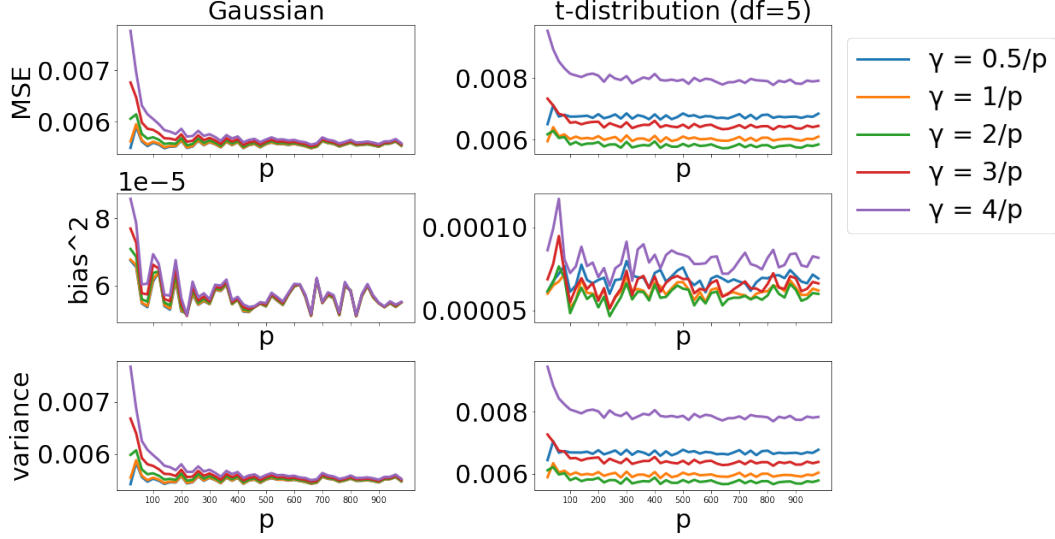


Figure 3: Effect of different  $\gamma$  values across different dimensions  $p \in [1, 1000]$  with  $\alpha = 0.1$ .

### 3.2.4 Scenario 4

For  $p \in [100, 1000]$  and under the Gaussian case, the MSE, bias and variance curves for  $\gamma$ -mean and simple  $\gamma$ -mean are collapsed together (Figure 4, left panel). For  $p \in [100, 1000]$  and under the  $t$ -distribution, these curves are nearly indistinguishable (Figure 4, right panel). When  $p$  is small ( $p \ll 100$ ), the MSE values of simple  $\gamma$ -mean are significantly larger than those of  $\gamma$ -mean. However, we did not display the results for  $p < 100$ , since the relative large MSE values for small  $p$  will make the MSE curves of  $\gamma$ -mean and simple  $\gamma$ -mean look lying on the  $x$ -axis for  $p \geq 100$ .

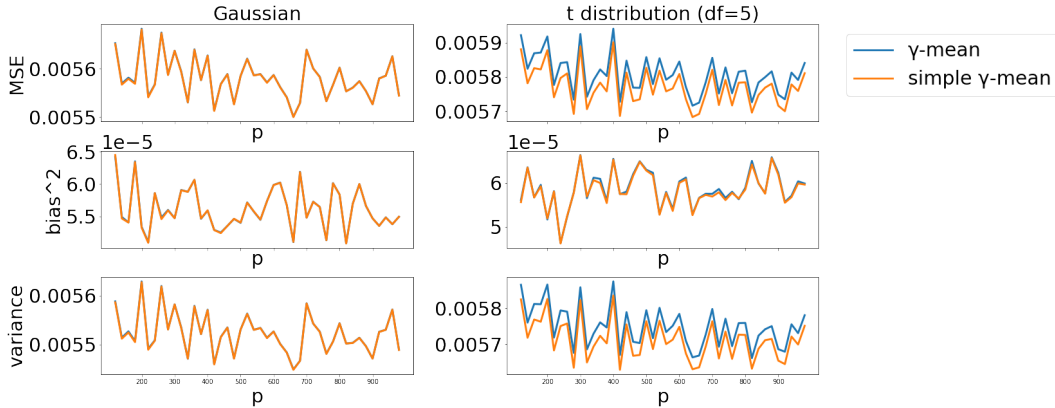


Figure 4: Comparison of  $\gamma$ -mean versus simple  $\gamma$ -mean with  $\alpha = 0.1$  and  $p \in [100, 1000]$ .

## 4 Real data examples

### 4.1 Datasets

- **MNIST** [6]. The MNIST database of handwritten digits has a training set of 60,000 examples, and a test set of 10,000 examples. The digits have been size-normalized and centered in a fixed-size,  $28 \times 28$  grayscale, images.
- **Fashion MNIST** [19]. Fashion-MNIST is a dataset of Zalando’s article images consisting of a training set of 60,000 examples and a test set of 10,000 examples. Each example is a  $28 \times 28$  grayscale image, associated with a label from 10 types of clothing, such as shoes, t-shirts, dresses, sandals, sneakers and more.
- **Chest x-ray images (pneumonia)** [11]. The dataset contains 5,856 x-ray images and 2 classes (pneumonia and normal). The 5,856 images consist of 5,232 training images (which we further split into 90% for model training and 10% for model validation to implement early stopping) and 624 testing images. Chest x-ray images (anterior-posterior) were selected from retrospective cohorts of pediatric patients of one to five years old from Guangzhou Women and Children’s Medical Center, Guangzhou.

### 4.2 Experimental setting

The experimental setting is described below.

- For MNIST and fashion MNIST, we set  $m = 20$  and the number of Byzantine clients is two. For chest x-ray images, we set  $m = 9$  and the number of Byzantine clients is one. Byzantine clients return random values from  $\text{Gaussian}(5, 1)$ .
- $\gamma$  is set to 0.5. This setting is different from the setting  $\gamma = c/p$  in simulation. The main reason is that there is a certain complicated relationship between  $\gamma$  value and the neural network model adopted, such as the dimensionality, gradient size, learning rate, etc. We have not yet fully understood this relationship, which might govern the selection of  $\gamma$ . We will leave it as a future study.
- We set  $\beta = 0.1$  for the trimmed mean.
- To allow imbalanced size of clients, we obtain the sample size of each client by sampling from the following steps [15].

1. Sampling a vector  $\mathbf{a}$  from Lognormal(1.5, 3.45<sup>2</sup>).
2. Sampling  $\boldsymbol{\eta}$  from Dirichlet( $\mathbf{a}$ ). The sum of vector  $\boldsymbol{\eta}$  will be 1 due to the property of Dirichlet distribution.
3. Obtain the sample sizes of clients from multinomial  $(n - m \times n_o, \boldsymbol{\eta})$ , where  $n$  is total sample size and  $n_o$  is the minimum sample size guaranteed for each client. We set  $n_o = 512$ .

In each round of parameter update, the server broadcasts the current parameter values to  $m$  clients. Each normal client computes the local gradient and local parameter update, and then returns the change of parameters after iterating  $k$  epochs using local data ( $k = 1$  for MNIST and fashion MNIST, and  $k = 20$  for chest x-ray images). Byzantine clients return random values from a Gaussian distribution  $N(5, 1)$ . Some further implementation details are given below.

- We run 1000 rounds of FL for MNIST and fashion MNIST, and 100 rounds for chest x-ray images.
- We apply stochastic gradient descent (SGD) with a cosine decay learning rate (decay over rounds), where the decay step is 1000 for MNIST and fashion MNIST, and 100 for chest x-ray images. In each epoch of local clients on MNIST and fashion MNIST, the SGD will go through only 10% of local data to save computing time. This implementation leads to some fluctuations in the early stage of training but the training process will be much faster than going through all local data.
- The initial learning rates are 0.1, 0.5, and  $10^{-4}$  for MNIST, fashion MNIST, and chest x-ray images, respectively. In each epoch in local iterates,  $10^{-5}$  is set as the decay constant.
- We apply gradient clipping to avoid exploding gradients on MNIST and fashion MNIST. If the 2-norm of aggregated gradient is larger than 1, the vector will be scaled to a new vector with norm 1.
- To handle the imbalanced class size, we use weighted cross-entropy as loss function in the chest x-ray example (pneumonia: 0.35, normal: 1.0), where chest x-ray training dataset contains 3,883 pneumonia cases and 1,349 normal cases. In addition to the classification accuracy, we also use ‘accuracy’, ‘sensitivity’ (also known as ‘recall’ and ‘true positive rate’) and ‘precision’ as our evaluation metrics. In particular, correctly predicting pneumonia is more important than predicting the normal case.

### 4.3 Models

We use a simple CNN model (Figure 8 in Appendix) for MNIST and fashion MNIST. For chest X-ray images (pneumonia), we use a pretrained ResNet50 (Figure 9 in Appendix) with input image size  $150 \times 150 \times 3$ . We further connect the model to 3 dense-block layers.

### 4.4 Results

Results for trimmed mean are not reported here, as they do not perform well due to its reliance on the specification of  $\beta$  and its non-robustness to non-symmetric outliers (note that the trimming procedure is symmetric in two tails of each coordinate).

#### 4.4.1 MNIST

Figure 5 shows the results of FL using different aggregators. In the case of no attack, the simple  $\gamma$ -mean performs as good as the mean-aggregator with 96% testing data classification accuracy. In the Byzantine case, simple  $\gamma$ -mean performs the best with testing accuracy 96% followed by the two median-based methods, geometric median and marginal median, while the mean-aggregator has failed. In addition, marginal median has some fluctuations in the early stage of training.

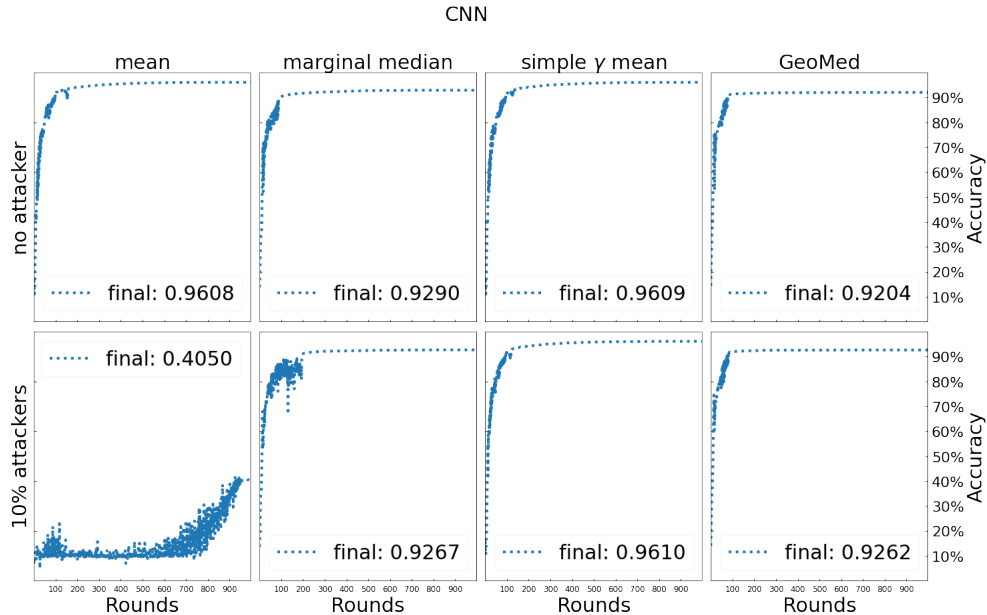


Figure 5: Testing process and comparison of testing accuracy for different aggregators on MNIST.

#### 4.4.2 Fashion MNIST

Figure 6 shows the results of fashion MNIST. The conclusion is similar to MNIST, that the simple  $\gamma$ -mean performs the best. Note that the fluctuations shown in the training process of marginal median look worse than those shown in Figure 5. Although we have applied gradient clipping, still there are some sudden surges in loss function values and they lead to sudden drops in classification accuracy.

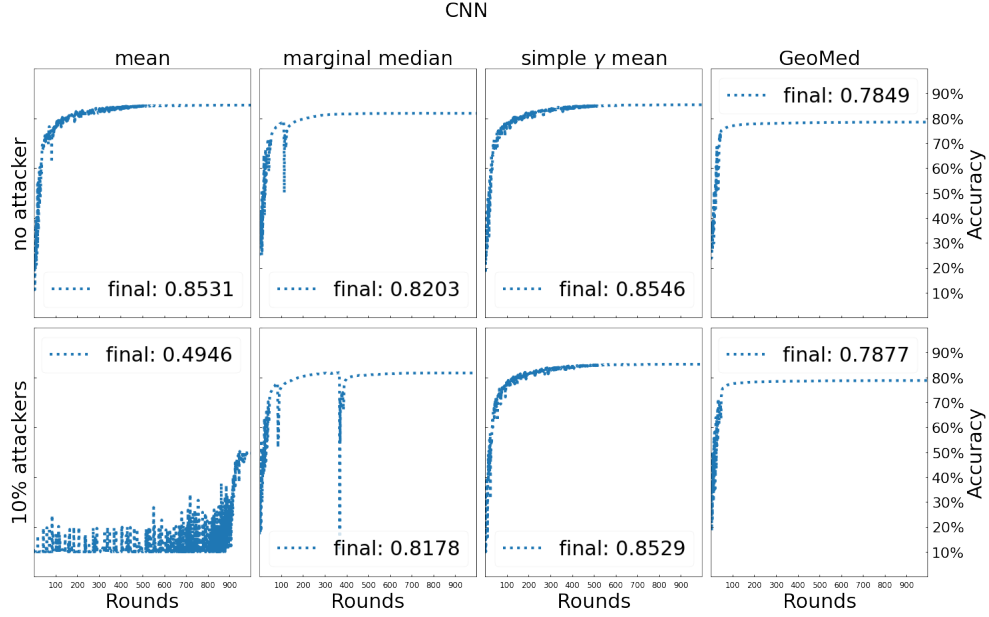


Figure 6: Testing process and comparison of testing accuracy for different aggregators on fashion MNIST.



#### 4.4.3 Chest x-ray images (pneumonia)

We train the model on the training dataset without any partitions (i.e. single machine with full training data) as a baseline. In this single machine case, the maximum number of training epochs is 1000 and the training process will stop early if the evaluation metric (loss + accuracy + precision + sensitivity on validation data) does not get better in 10 consecutive epochs. For FL case, the number of FL rounds is 100 and the maximum training epochs of each client are 20 at every FL round. At each round, the clients will stop early if the evaluation metric (loss + accuracy + precision + sensitivity on validation data) does not get better in 5 consecutive epochs. The simple  $\gamma$ -mean algorithm has adopted the marginal median as initial. All other settings are the same as those listed in Section 4.2.

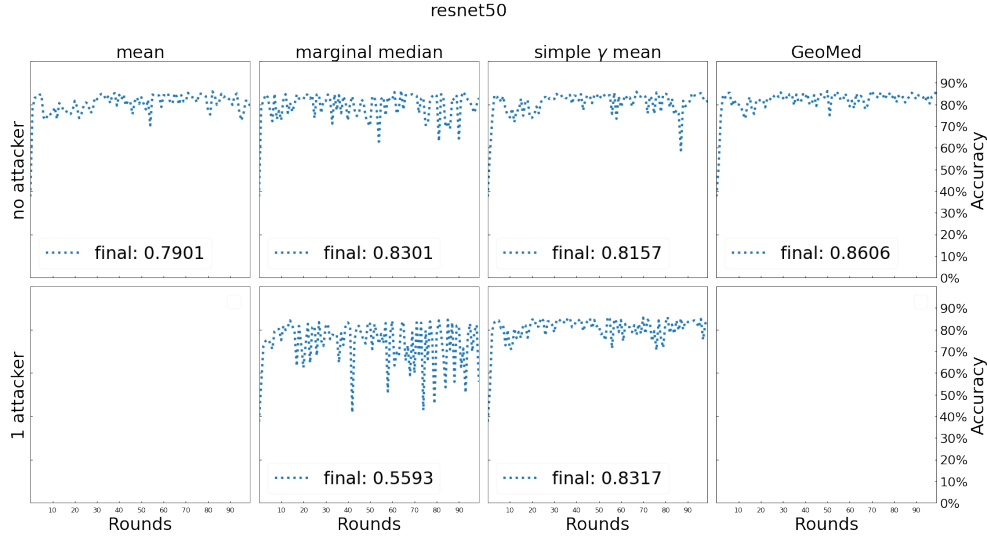


Figure 7: Testing process and comparison of testing accuracy for different aggregators on chest x ray. The aggregators by mean and geometric median cannot tolerate the Byzantine attack and both methods were crushed during model training. Thus, there are no results reported for these two methods.

We further report the following testing data results in Table 1: true positive (TP, which is predicted positive and actually positive), true negative (TN, which is predicted negative and actually negative), false positive (FP, which is predicted positive but actually negative), false negative (FN, which is predicted negative but actually positive), precision (Prec, which is  $\frac{TP}{TP+FP}$ ), sensitivity (Sens, which is  $\frac{TP}{TP+FN}$  as well as one minus type-II error) and classification accuracy (Acc). The aggregators by mean and geometric median cannot tolerate the Byzantine (Byz) attack and both methods were crushed during model training. Thus, there are no results reported for these two methods.

Table 1: Pneumonia prediction on test data.

Byz	Aggregator	TN	FN	FP	TP	Prec	Sens (Type II error)	Acc
	single machine	156	23	78	367	0.8247	0.9410 (0.0590)	0.8381
No	mean	212	103	22	287	0.9288	0.7359 (0.2661)	0.7997
	marginal median	190	63	44	327	0.8814	0.8385 (0.1615)	0.8285
	simple $\gamma$ -mean	126	8	108	382	0.7796	0.9795 (0.0205)	0.8141
	GeoMed	177	30	57	360	0.8633	0.9231 (0.0769)	<b>0.8606</b>
Yes	mean <sup>†</sup>	–	–	–	–	–	–	–
	marginal median	228	271	6	119	0.9520	0.3051 (0.6949)	0.5561
	simple $\gamma$ -mean	140	11	94	379	0.8013	0.9718 (0.0282)	<b>0.8317</b>
	GeoMed <sup>†</sup>	–	–	–	–	–	–	–

<sup>†</sup>The symbol “–” indicates model crushed during training.

## 5 Concluding remarks

Our major contribution in this article is to propose the  $\gamma$ -mean aggregation in federated learning, which is robust against outliers and Byzantine attacks. We have provided some theoretical discussions on influence functions and carried out numerical studies to justify our proposal. In both simulation and real data experiments,  $\gamma$ -mean based aggregation in general outperforms the existing robust aggregators such as marginal median, geometric median and trimmed mean. In addition, simple  $\gamma$ -mean does not require complex computation. It can be easily computed by fixed point iteration and works well for extremely high dimensional case such as in a deep neural network model.  $\gamma$ -mean (as well as simple  $\gamma$ -mean) takes a form of weighted average, where  $\gamma$  controls the trade-off between robustness and estimation efficiency. The selection of  $\gamma$  value is important, and a data-driven selection procedure is needed, which will be pursued in our future study.

# Appendix

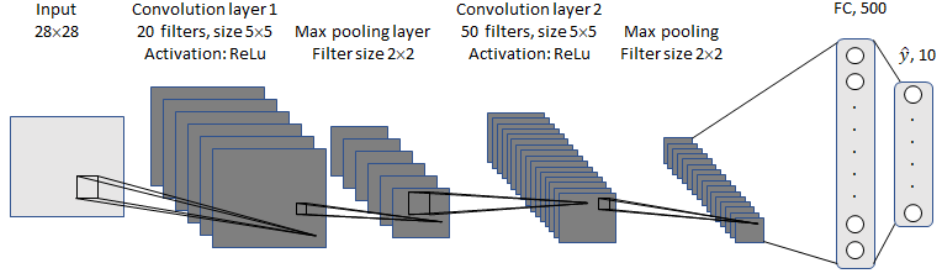


Figure 8: Model used in MNIST and fashion MNIST

For MNIST and fashion MNIST, we use a model similar to LeNet-5. All filters in the convolution layers are  $5 \times 5$  maps without padding. The first layer in fully connection is a dense layer with 500 neurons with ReLU activation. The last layer is an output layer with 10 neurons with softmax activation. The number of parameters is 931,080.

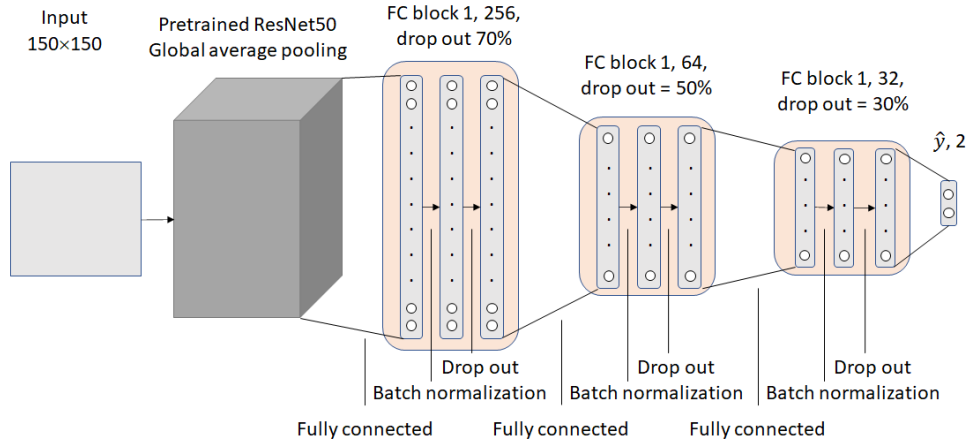


Figure 9: Model used in chest x-ray images.

For chest x-ray images, we use a pretrained ResNet50 with input image size is  $150 \times 150 \times 3$  and freeze all ResNet50 parameters. The output dimensionality is  $5 \times 5 \times 2,048$ . A global average pooling is adopted and it leads to the dimensionality 2,048. Next, we connect it to 3 dense-block layers. Dense block: dense layer  $\rightarrow$  batch normalization  $\rightarrow$  drop out. Number of neurons are 256, 64, and 32, respectively. The drop out probabilities are 0.7, 0.5, and 0.3, respectively. The number of trainable parameters is 543,809.

## References

- [1] Alistarh, D., Allen-Zhu, Z. and Li, J. (2018). Byzantine stochastic gradient descent. *Advances in Neural Information Processing Systems*, 31, Curran Associates, Inc.
- [2] Chaudhuri, P. (1996). On a geometric notion of quantiles for multivariate data. *Journal of the American Statistical Association*, 91(434): 862-872.
- [3] Chen, X., Chen, T, Sun, H., Wu, S. Z., and Hong, M. (2020). Distributed training with heterogeneous data: bridging median- and mean-based algorithms. *Advances in Neural Information Processing Systems*, 33, Curran Associates, Inc.
- [4] Chen, Y., Su, L., and X, J. (2017). Distributed statistical machine learning in adversarial settings: Byzantine gradient descent. *Proceedings of the ACM on Measurement and Analysis of Computing Systems*, 1(44): 1-25.
- [5] Dayan, I., Roth, H. R., Zhong, A., Harouni, A., Gentili, A., Abidin, A. Z., ..., and Li, Q. (2021). Federated learning for predicting clinical outcomes in patients with covid-19. *Nature Medicine*, 27: 1735-1743.
- [6] Deng, L. (2012). The MNIST database of handwritten digit images for machine learning research. *IEEE Signal Processing Magazine*, 29(6): 141-142.
- [7] Fujisawa, H. and Eguchi, S. (2008). Robust parameter estimation with a small bias against heavy contamination. *Journal of Multivariate Analysis*, 99: 2053-2081.
- [8] Huber, P. J. (2004). *Robust Statistics*. John Wiley & Sons.
- [9] Hung, H. (2019). A robust removing unwanted variation-testing procedure via  $\gamma$ -divergence. *Biometrics*, 75(2): 650-662.
- [10] Jones, M. C., Hjort, N. L., Harris, I. R., and Basu, A. (2001). A comparison of related density-based minimum divergence estimators. *Biometrika*, 88(3): 865-873.
- [11] Kermany, D. S., Goldbaum, M., Cai, W., Valentim, C., Liang, H., Baxter, S. L., McKeown, A., Yang, G., Wu, X., Yan, F., Dong, J., Prasadha, M. K., Pei, J., Ting, M., Zhu, J., Li, C., Hewett, S., Dong, J., Ziyar, I., Shi, A., ..., and Zhang, K. (2018). Identifying medical diagnoses and treatable diseases by image-based deep learning. *Cell*, 172(5): 1122-1131.e9.

- [12] Konecny, J., McMahan, H. B., Yu, F. X., Richtarik, P., Suresh, A. T., and Bacon, D. (2016). Federated learning: Strategies for improving communication efficiency. In *NeurIPS Workshop on Private Multi-Party Machine Learning*.
- [13] Li, L., Xu, W., Chen, T., Giannakis, G. B., and Ling, Q. (2019). RSA: Byzantine-robust stochastic aggregation methods for distributed learning from heterogeneous datasets. In *Proceedings of the AAAI Conference on Artificial Intelligence*, 33: 1544-1551.
- [14] McMahan, B., Moore, E., Ramage, D., Hampson, S., and y Arcas, B. A. (2017). Communication-efficient learning of deep networks from decentralized data. In *Artificial Intelligence and Statistics*, 1273-1282.
- [15] Portnoy, A., Tirosh, Y., and Hendler, D. (2021). Towards federated learning with byzantine-robust client weighting. In *International Workshop on Federated Learning for User Privacy and Data Confidentiality in Conjunction with ICML*.
- [16] So, J., Güler, B., and Avestimehr, A. S. (2020). Byzantine-resilient secure federated learning. In *IEEE Journal on Selected Areas in Communications*, 39(7): 2168-2181.
- [17] van der Vaart, A. W. (1998). *Asymptotic Statistics*. Cambridge Series in Statistical and Probabilistic Mathematics. Cambridge University Press.
- [18] Weiszfeld, E., and Plastria, F. (2009). On the point for which the sum of the distances to n given points is minimum. *Annals of Operations Research*, 167: 7-41.
- [19] Xiao, H., Rasul, K., and Vollgraf, R. (2017). Fashion-MNIST a novel image dataset for benchmarking machine learning algorithms. arXiv:1708.07747.
- [20] Xie, C., Koyejo, O., and Gupta, I. (2018). Generalized Byzantine-tolerant SGD. arXiv:1802.10116
- [21] Xu, J., Glicksberg, B. S., Su, C., Walker, P., Bian, J., and Wang, F. (2021). Federated learning for healthcare informatics. *Journal of Healthcare Informatics Research*, 5(1): 1-19.