

# Robust federated learning via voting mechanism

Xiang-Yu Liang  
*School of Computer Science*  
*Southwest Petroleum University*  
 Chengdu 610500, China  
 liangxyswpu@163.com  
 Fan Min  
*School of Computer Science*  
*Southwest Petroleum University*  
 Chengdu 610500, China  
 minfan@swpu.edu.cn

Heng-Ru Zhang\*  
*School of Computer Science*  
*Southwest Petroleum University*  
 Chengdu 610500, China  
 zhanghrswpu@163.com  
 Wei Tang  
*School of Computer Science*  
*Southwest Petroleum University*  
 Chengdu 610500, China  
 twei1123@163.com

## I. FORMAL SECURITY ANALYSIS

**Assumption 1:** The global model objective function  $F(\mathbf{w})$  is  $L$ -strongly convex and has an  $M$ -Lipschitz continuous gradient on  $\omega$ . For any  $\mathbf{w}, \mathbf{w}' \in \omega$ , we have the following:

$$F(\mathbf{w}) + \langle \nabla F(\mathbf{w}), \mathbf{w}' - \mathbf{w} \rangle + \frac{L}{2} \|\mathbf{w}' - \mathbf{w}\|^2 \leq F(\mathbf{w}'),$$

$$\|\nabla F(\mathbf{w}) - \nabla F(\mathbf{w}')\| \leq M \|\mathbf{w} - \mathbf{w}'\|,$$

where  $\langle \cdot, \cdot \rangle$  is the inner product of two vectors,  $\nabla$  is the gradient, and  $\|\cdot\|$  is the  $\ell_2$  norm.

**Assumption 2:** There exist positive constants  $\sigma_1$  and  $\gamma_1$  such that for every unit vector  $\mathbf{v} \in \mathbf{B}$ ,  $\langle \nabla f(D, \mathbf{w}^*), \mathbf{v} \rangle$  is sub-exponential with scaling parameters  $\sigma_1$  and  $\gamma_1$ , i.e.,

$$\sup_{\mathbf{v} \in \mathbf{B}} \mathbb{E}[\exp(\lambda \langle \nabla f(D, \mathbf{w}^*), \mathbf{v} \rangle)] \leq e^{\sigma_1^2 \lambda^2 / 2}, \quad \forall |\lambda| \leq \frac{1}{\gamma_1},$$

where  $\mathbf{B}$  denotes the unit sphere  $\{\mathbf{v} : \|\mathbf{v}\| = 1\}$ .

Assumption 2 is to ensure that the client uses the local dataset with high probability to find the optimal model  $\mathbf{w}^*$ . Specifically,  $(1/|D_i|) \sum_{X_j \in D_i} \nabla f(X_j, \mathbf{w}^*)$  is concentrated near  $\nabla F(\mathbf{w}^*) = 0$ , where  $|D_i|$  is represented as the number of elements of  $D_i$ .

Next, we define gradient difference:

$$h(D, \mathbf{w}) \triangleq \nabla f(D, \mathbf{w}) - \nabla f(D, \mathbf{w}^*), \quad (1)$$

which expresses the deviation of the empirical loss function from the optimal global model. Note that

$$\mathbb{E}[h(D, \mathbf{w})] = \nabla F(\mathbf{w}) - \nabla F(\mathbf{w}^*), \quad (2)$$

for each  $\mathbf{w}$ .

**Assumption 3:** There exist positive constant  $\sigma_2$  and  $\gamma_2$  such that for any  $\mathbf{w} \in \omega$  with  $\mathbf{w} \neq \mathbf{w}^*$  and any unit vector  $\mathbf{v} \in \mathbf{B}$ ,  $\langle h(D, \mathbf{w}) - \mathbb{E}[h(D, \mathbf{w})], \mathbf{v} \rangle / \|\mathbf{w} - \mathbf{w}^*\|$  is sub-exponential with scaling parameters  $\sigma_2$  and  $\gamma_2$ , i.e., for all  $|\lambda| < 1/\gamma_2$ ,

$$\sup_{\mathbf{w} \in \omega, \mathbf{v} \in \mathbf{B}} \mathbb{E}[\exp(\frac{\lambda \langle h(D, \mathbf{w}) - \mathbb{E}[h(D, \mathbf{w})], \mathbf{v} \rangle}{\|\mathbf{w} - \mathbf{w}^*\|})] \leq e^{\sigma_2^2 \lambda^2 / 2},$$

where  $\mathbf{B}$  denotes the unit sphere  $\{\mathbf{v} : \|\mathbf{v}\| = 1\}$ .

**Assumption 4:** For any  $\delta \in (0, 1)$ , there exists an  $M' = M'(|D_i|, \delta)$  that is non-increasing in  $|D_i|, \delta$  such that

$$\mathbb{P}\left\{ \sup_{\mathbf{w}, \mathbf{w}' \in \omega: \mathbf{w} \neq \mathbf{w}'} \frac{\|\nabla \bar{f}_{|D_i|}(\mathbf{w}) - \nabla \bar{f}_{|D_i|}(\mathbf{w}')\|}{\|\mathbf{w} - \mathbf{w}'\|} \leq M' \right\} \geq 1 - \frac{\delta}{3},$$

where  $\nabla \bar{f}_{|D_i|}(\mathbf{w}) = (\sum_{X_j \in D_i} \nabla f(X_j, \mathbf{w})) / |D_i|$ .

**Assumption 5:** Each local training dataset  $D_i$  ( $i = 1, 2, \dots, n$ ) is sampled from distribution  $\mathcal{X}$ .

**Theorem 1:** Suppose Assumptions 1-5 hold, learning rate  $\alpha = L/2M^2$ ,  $\delta \in (0, 1)$ ,  $\Delta_1 \geq \sigma_1^2/\gamma_1$ ,  $\Delta_2 \geq \sigma_2^2/\gamma_2$ , and  $\omega \subset \{\mathbf{w} : \|\mathbf{w} - \mathbf{w}^*\| \leq r\sqrt{d}\}$  for some positive parameter  $r$ , for any number of malicious clients, the difference between the global model aggregated by VSRFL and the optimal global model  $\mathbf{w}^*$  without attack is bounded. For any  $t \geq 1$ , we have:

$$\|\mathbf{w}^t - \mathbf{w}^*\| \leq (1 - \rho)^t \|\mathbf{w}^0 - \mathbf{w}^*\| + \frac{12\alpha\Delta_1}{\rho}.$$

where  $\mathbf{w}^t$  is the global model of the aggregation for each epoch,  $\rho = 1 - (\sqrt{1 - L^2/(4M^2)} + 24\alpha\Delta_2 + 2\alpha M)$ ,  $\Delta_1 = \sigma_1 \sqrt{2/|D_i|} \sqrt{d \log 6 + \log(3/\delta)}$ ,  $\Delta_2 = \sigma_2 \sqrt{\frac{2}{|D_i|}} \sqrt{d \log \frac{18M \vee M'}{\sigma_2} + \frac{1}{2} d \log \frac{|D_i|}{d} + \log(\frac{6\sigma_2^2 r \sqrt{|D_i|}}{\gamma_2 \sigma_1 \delta})}$ ,  $M \vee M' = \max(M, M')$ ,  $d$  is the dimension of  $\mathbf{w}$ . When  $|1 - \rho| < 1$ , we have  $\lim_{t \rightarrow \infty} \|\mathbf{w}^t - \mathbf{w}^*\| \leq 12\alpha\Delta_1/\rho$ .

## II. THE PROVING PROCESS

Recall that the optimal global model  $\mathbf{w}^*$  is a answer to the following optimization problem:  $\mathbf{w}^* = \arg \min_{\mathbf{w}} F(\mathbf{w})$ , where  $F(\mathbf{w}) = \mathbb{E}_{D \sim \mathcal{X}}[f(D, \mathbf{w})]$  is the expectation of the empirical loss  $f(D, \mathbf{w})$  on the joint training dataset  $D$ . We show that the difference between the global model trained by VSRFL and the optimal global model  $\mathbf{w}^*$  is bounded under certain assumptions. We denote the local update set filtered by the server in epoch  $t$  by  $\mathcal{S}$ . We let  $\hat{\mathbf{g}}_i = (\|\mathbf{g}_{median}\|/\|\mathbf{g}_i\|) \times \mathbf{g}_i$ , where  $i \in \mathcal{S}$  s.t.  $|\mathcal{S}| > 0$ . We let  $|D_i|$  is size of the local dataset. We first describe our lemmas and then state our theoretical results.

**Lemma 1:** For any number of abnormal local updates, the gap between the global update  $\mathbf{g}$  and the gradient  $\nabla F(\mathbf{w})$  is bounded:

$$\|\mathbf{g} - \nabla F(\mathbf{w})\| \leq 3\|\mathbf{g}_{\text{median}} - \nabla F(\mathbf{w})\| + 2\|\nabla F(\mathbf{w})\|,$$

where  $\mathbf{g}_{\text{median}}$  is the median updated of the server selection in each epoch.

*Proof:* We have the following equations:

$$\begin{aligned} & \|\mathbf{g} - \nabla F(\mathbf{w})\| \\ &= \left\| \frac{1}{|\mathcal{S}|} \sum_{i \in \mathcal{S}} \hat{\mathbf{g}}_i - \nabla F(\mathbf{w}) \right\| \\ &= \left\| \frac{1}{|\mathcal{S}|} \sum_{i \in \mathcal{S}} \hat{\mathbf{g}}_i - \mathbf{g}_{\text{median}} + \mathbf{g}_{\text{median}} - \nabla F(\mathbf{w}) \right\| \\ &\leq \left\| \frac{1}{|\mathcal{S}|} \sum_{i \in \mathcal{S}} \hat{\mathbf{g}}_i - \mathbf{g}_{\text{median}} \right\| + \|\mathbf{g}_{\text{median}} - \nabla F(\mathbf{w})\| \\ &= \left\| \frac{1}{|\mathcal{S}|} \sum_{i \in \mathcal{S}} \hat{\mathbf{g}}_i + (-\mathbf{g}_{\text{median}}) \right\| + \|\mathbf{g}_{\text{median}} - \nabla F(\mathbf{w})\| \\ &\leq \left\| \frac{1}{|\mathcal{S}|} \sum_{i \in \mathcal{S}} \hat{\mathbf{g}}_i \right\| + \|-\mathbf{g}_{\text{median}}\| + \|\mathbf{g}_{\text{median}} - \nabla F(\mathbf{w})\| \\ &\stackrel{(a)}{=} \left\| \frac{1}{|\mathcal{S}|} \sum_{i \in \mathcal{S}} \hat{\mathbf{g}}_i \right\| + \|\mathbf{g}_{\text{median}}\| + \|\mathbf{g}_{\text{median}} - \nabla F(\mathbf{w})\| \\ &\leq \frac{1}{|\mathcal{S}|} \sum_{i \in \mathcal{S}} \|\hat{\mathbf{g}}_i\| + \|\mathbf{g}_{\text{median}}\| + \|\mathbf{g}_{\text{median}} - \nabla F(\mathbf{w})\| \\ &\stackrel{(b)}{=} \frac{1}{|\mathcal{S}|} \sum_{i \in \mathcal{S}} \|\mathbf{g}_{\text{median}}\| + \|\mathbf{g}_{\text{median}}\| + \|\mathbf{g}_{\text{median}} - \nabla F(\mathbf{w})\| \\ &\stackrel{(c)}{=} 2\|\mathbf{g}_{\text{median}}\| + \|\mathbf{g}_{\text{median}} - \nabla F(\mathbf{w})\| \\ &= 2\|\mathbf{g}_{\text{median}} - \nabla F(\mathbf{w}) + \nabla F(\mathbf{w})\| + \|\mathbf{g}_{\text{median}} - \nabla F(\mathbf{w})\| \\ &\leq 2\|\mathbf{g}_{\text{median}} - \nabla F(\mathbf{w})\| + 2\|\nabla F(\mathbf{w})\| + \|\mathbf{g}_{\text{median}} - \nabla F(\mathbf{w})\| \\ &= 3\|\mathbf{g}_{\text{median}} - \nabla F(\mathbf{w})\| + 2\|\nabla F(\mathbf{w})\|, \end{aligned} \quad (3)$$

where (a) is because of the following equations:

$$\begin{aligned} \sqrt{g_1^2 + g_2^2 + \dots + g_i^2} &= \sqrt{(-g_1)^2 + (-g_2)^2 + \dots + (-g_i)^2}, \\ \text{s.t. } \mathbf{g} &= \{g_1, g_2, \dots, g_i\}; \end{aligned} \quad (4)$$

(b) is because VSRFL normalizes the filtered local updates to have the same magnitude as the median, i.e.,  $\|\hat{\mathbf{g}}_i\| = \|\mathbf{g}_{\text{median}}\|$ ; and (c) is because  $|\mathcal{S}|$  is represented as the number of elements of  $\mathcal{S}$ , e.g.,  $(\sum_{i \in \mathcal{S}} x)/|\mathcal{S}| = x$ .

**Lemma 2:** Suppose Assumption 1 holds. If we choose the learning rate  $\alpha = L/2M^2$ , there is the following inequality:

$$\begin{aligned} \|\mathbf{w}^{t-1} - \mathbf{w}^* - \alpha \nabla F(\mathbf{w}^{t-1})\| &\leq \sqrt{\left(1 - \frac{L^2}{4M^2}\right)} \|\mathbf{w}^{t-1} - \mathbf{w}^*\| \\ \text{s.t. } t &\geq 1. \end{aligned}$$

*Proof:* By Assumption 1, we have:

$$\|\nabla F(\mathbf{w}^{t-1}) - \nabla F(\mathbf{w}^*)\| \leq M\|\mathbf{w}^{t-1} - \mathbf{w}^*\|, \quad (5)$$

$$\begin{aligned} F(\mathbf{w}^{t-1}) &\geq F(\mathbf{w}^*) + \langle \nabla F(\mathbf{w}^*), \mathbf{w}^{t-1} - \mathbf{w}^* \rangle \\ &\quad + \frac{L}{2} \|\mathbf{w}^{t-1} - \mathbf{w}^*\|^2, \end{aligned} \quad (6)$$

$$F(\mathbf{w}^*) \geq F(\mathbf{w}^{t-1}) + \langle \nabla F(\mathbf{w}^{t-1}), \mathbf{w}^* - \mathbf{w}^{t-1} \rangle. \quad (7)$$

Combining equations 6 and 7, we have:

$$\begin{aligned} & \langle \nabla F(\mathbf{w}^*), \mathbf{w}^{t-1} - \mathbf{w}^* \rangle + \langle \nabla F(\mathbf{w}^{t-1}), \mathbf{w}^* - \mathbf{w}^{t-1} \rangle \\ &= \langle \nabla F(\mathbf{w}^*), \mathbf{w}^{t-1} - \mathbf{w}^* \rangle - \langle \nabla F(\mathbf{w}^{t-1}), \mathbf{w}^{t-1} - \mathbf{w}^* \rangle \\ &= \langle \nabla F(\mathbf{w}^*) - \nabla F(\mathbf{w}^{t-1}), \mathbf{w}^{t-1} - \mathbf{w}^* \rangle \\ &= -\langle \nabla F(\mathbf{w}^{t-1}) - \nabla F(\mathbf{w}^*), \mathbf{w}^{t-1} - \mathbf{w}^* \rangle \\ &\leq -\frac{L}{2} \|\mathbf{w}^{t-1} - \mathbf{w}^*\|^2. \end{aligned} \quad (8)$$

Due to the that  $\nabla F(\mathbf{w}^*) = \mathbf{0}$ , we have the following:

$$\begin{aligned} & \|\mathbf{w}^{t-1} - \mathbf{w}^* - \alpha \nabla F(\mathbf{w}^{t-1})\|^2 \\ &= \|\mathbf{w}^{t-1} - \mathbf{w}^* - \alpha (\nabla F(\mathbf{w}^{t-1}) - \nabla F(\mathbf{w}^*)) + \mathbf{w}^{t-1} - \mathbf{w}^*\|^2 \\ &= \alpha^2 \|\nabla F(\mathbf{w}^{t-1}) - \nabla F(\mathbf{w}^*)\|^2 + \|\mathbf{w}^{t-1} - \mathbf{w}^*\|^2 \\ &\quad - 2\alpha \langle \nabla F(\mathbf{w}^{t-1}) - \nabla F(\mathbf{w}^*), \mathbf{w}^{t-1} - \mathbf{w}^* \rangle \\ &\leq \alpha^2 M^2 \|\mathbf{w}^{t-1} - \mathbf{w}^*\|^2 + \|\mathbf{w}^{t-1} - \mathbf{w}^*\|^2 \\ &\quad - \alpha L \|\mathbf{w}^{t-1} - \mathbf{w}^*\|^2 \\ &= (1 + \alpha^2 M^2 - \alpha L) \|\mathbf{w}^{t-1} - \mathbf{w}^*\|^2. \end{aligned} \quad (9)$$

We let  $\alpha = L/2M^2$ , therefore:

$$\begin{aligned} & \|\mathbf{w}^{t-1} - \mathbf{w}^* - \alpha \nabla F(\mathbf{w}^{t-1})\|^2 \\ &\leq (1 + \alpha^2 M^2 - \alpha L) \|\mathbf{w}^{t-1} - \mathbf{w}^*\|^2 \\ &= \left(1 - \frac{L^2}{4M^2}\right) \|\mathbf{w}^{t-1} - \mathbf{w}^*\|^2, \end{aligned} \quad (10)$$

which concludes the proof.

**Lemma 3:** Suppose Assumption 2 holds. For any  $\delta \in (0, 1)$  and any positive integer  $|D_i|$ , we let

$$\Delta_1(|D_i|, d, \delta, \sigma_1) = \sqrt{2}\sigma_1 \sqrt{\frac{d \log 6 + \log(3/\delta)}{|D_i|}}.$$

We let  $\Delta_1 = \Delta_1(|D_i|, d, \delta, \sigma_1)$ . If  $\Delta_1 \leq \sigma_1^2/\gamma_1$ , we have

$$\mathbb{P}\left\{\left\|\frac{1}{|D_i|} \sum_{X_j \in D_i} \nabla f(X_j, \mathbf{w}^*) - \nabla F(\mathbf{w}^*)\right\| \geq 2\Delta_1\right\} \leq \frac{\delta}{3}.$$

For fixed  $\delta$  and  $\sigma_1$ , if  $d = o(|D_i|)$ ,

$$\Delta_1 = \sqrt{2}\sigma_1 \sqrt{\frac{d \log 6 + \log(3/\delta)}{|D_i|}} \rightarrow 0 \text{ as } |D_i| \rightarrow \infty.$$

So, if  $\gamma_1$  is fixed,  $\Delta_1 \leq \sigma_1^2/\gamma_1$  holds when  $l$  is large enough.

*Proof:* We let  $\mathcal{V} = \{v_1, \dots, v_{N_{1/2}}\}$  denote an  $\frac{1}{2}$ -cover of unit sphere  $\mathcal{B}$ . It is show in [1], [2] that  $\log N_{1/2} \leq d \log 6$ , and

$$\begin{aligned} & \left\| \frac{1}{|D_i|} \sum_{X_j \in D_i} \nabla f(X_j, \mathbf{w}^*) - \nabla F(\mathbf{w}^*) \right\| \\ &\leq 2 \sup_{v \in \mathcal{V}} \left\{ \left\langle \frac{1}{|D_i|} \sum_{X_j \in D_i} \nabla f(X_j, \mathbf{w}^*) - \nabla F(\mathbf{w}^*), v \right\rangle \right\}. \end{aligned} \quad (11)$$

By Assumption 2, the condition  $\Delta_1 \leq \sigma_1^2/\gamma_1$ , and the concentration inequalities for sub-exponential random variables, for  $v \in \mathcal{V}$  we have:

$$\begin{aligned} \mathbb{P}\left\{\left\langle \frac{1}{|D_i|} \sum_{X_j \in D_i} \nabla f(X_j, \mathbf{w}^*) - \nabla F(\mathbf{w}^*), v \right\rangle \geq \Delta_1\right\} \\ \leq \exp\left(-\frac{|D_i|\Delta_1^2}{2\sigma_1^2}\right). \end{aligned} \quad (12)$$

Recall that in  $\mathcal{V}$  contains at most  $6^d$  vectors. In view of the union bound, it further yields that

$$\begin{aligned} \mathbb{P}\left\{2 \sup_{v \in \mathcal{V}} \left\langle \frac{1}{|D_i|} \sum_{X_j \in D_i} \nabla f(X_j, \mathbf{w}^*) - \nabla F(\mathbf{w}^*), v \right\rangle \geq 2\Delta_1\right\} \\ \leq 6^d \exp\left(-\frac{|D_i|\Delta_1^2}{2\sigma_1^2}\right) \\ = \exp\left(-\frac{|D_i|\Delta_1^2}{2\sigma_1^2} + d \log 6\right). \end{aligned} \quad (13)$$

Therefore,

$$\begin{aligned} \mathbb{P}\left\{\left\| \frac{1}{|D_i|} \sum_{X_j \in D_i} \nabla f(X_j, \mathbf{w}^*) - \nabla F(\mathbf{w}^*) \right\| \geq 2\Delta_1\right\} \\ \leq \exp\left(-\frac{|D_i|\Delta_1^2}{2\sigma_1^2} + d \log 6\right). \end{aligned} \quad (14)$$

We conclude the proof by equation  $\Delta_1 = \sqrt{2}\sigma_1\sqrt{(d \log 6 + \log(3/\delta))/|D_i|}$ .

**Lemma 4:** Suppose Assumption 3 holds and fix any  $\mathbf{w} \in \omega$ . We let

$$\Delta'_1(|D_i|, d, \delta, \sigma_2) = \sqrt{2}\sigma_2\sqrt{\frac{d \log 6 + \log(3/\delta)}{|D_i|}}.$$

We let  $\Delta'_1 = \Delta'_1(|D_i|, d, \delta, \sigma_2)$ . If  $\Delta'_1 \leq \sigma_2^2/\gamma_2$ , then

$$\begin{aligned} \mathbb{P}\left\{\left\| \frac{1}{|D_i|} \sum_{X_j \in D_i} \nabla h(X_j, \mathbf{w}) - \mathbb{E}[h(X, \mathbf{w})] \right\| \geq 2\Delta'_1(\mathbf{w} - \mathbf{w}^*)\right\} \\ \leq \frac{\delta}{3}. \end{aligned}$$

Similar to  $\Delta_1$ , if  $\delta, \sigma_1$  and  $\sigma_2$  are fixed, and  $d = o(|D_i|)$ , then for all sufficiently large  $l$ , it holds that  $\Delta'_1(l, d, \delta, \sigma_2) \leq \sigma_2^2/\gamma_2$ .

*Proof:* It is similar to the proof of Lemma 3. Let  $\mathcal{V} = \{v_1, \dots, v_{N_{1/2}}\}$  denote an  $\frac{1}{2}$ -cover of unit sphere  $\mathcal{B}$ . This exists  $\log N_{1/2} \leq d \log 6$ , and

$$\begin{aligned} \left\| \frac{1}{|D_i|} \sum_{X_j \in D_i} \nabla h(X_j, \mathbf{w}) - \mathbb{E}[h(X, \mathbf{w})] \right\| \\ \leq 2 \sup_{v \in \mathcal{V}} \left\langle \frac{1}{|D_i|} \sum_{X_j \in D_i} \nabla h(X_j, \mathbf{w}) - \mathbb{E}[h(X, \mathbf{w})], v \right\rangle. \end{aligned} \quad (15)$$

By Assumption 3, the condition  $\Delta'_1 \leq \sigma_2^2/\gamma_2$ , and the concentration inequalities for sub-exponential random variables, for  $v \in \mathcal{V}$  we have:

$$\begin{aligned} \mathbb{P}\left\{\left\langle \frac{1}{|D_i|} \sum_{X_j \in D_i} \nabla h(X_j, \mathbf{w}) - \mathbb{E}[h(X, \mathbf{w})], v \right\rangle \geq \Delta'_1(\mathbf{w} - \mathbf{w}^*)\right\} \\ \leq \exp\left(-\frac{|D_i|(\Delta'_1)^2}{2\sigma_2^2}\right). \end{aligned} \quad (16)$$

Recall that in  $\mathcal{V}$  contains at most  $6^d$  vectors. In view of the union bound, it further yields that

$$\begin{aligned} \mathbb{P}\left\{2 \sup_{v \in \mathcal{V}} \left\langle \frac{1}{|D_i|} \sum_{X_j \in D_i} \nabla h(X_j, \mathbf{w}) - \mathbb{E}[h(X, \mathbf{w})], v \right\rangle \geq 2\Delta'_1(\mathbf{w} - \mathbf{w}^*)\right\} \\ \leq 6^d \exp\left(-\frac{|D_i|(\Delta'_1)^2}{2\sigma_2^2}\right) \\ = \exp\left(-\frac{|D_i|(\Delta'_1)^2}{2\sigma_2^2} + d \log 6\right). \end{aligned} \quad (17)$$

Therefore,

$$\begin{aligned} \mathbb{P}\left\{\left\| \frac{1}{|D_i|} \sum_{X_j \in D_i} \nabla h(X_j, \mathbf{w}) - \mathbb{E}[h(X, \mathbf{w})] \right\| \geq 2\Delta'_1(\mathbf{w} - \mathbf{w}^*)\right\} \\ \leq \exp\left(-\frac{|D_i|(\Delta'_1)^2}{2\sigma_2^2} + d \log 6\right). \end{aligned} \quad (18)$$

We conclude the proof by equation  $\Delta'_1 = \sqrt{2}\sigma_2\sqrt{(d \log 6 + \log(3/\delta))/|D_i|}$ .

**Lemma 5:** Given a real number  $r > 0$ , we let

$$\Delta_2(|D_i|) = \sigma_2\sqrt{\frac{2}{|D_i|}}\sqrt{K_1 + K_2 + K_3},$$

where  $K_1 = d \log \frac{18M \vee M'}{\sigma_2}$ ,  $K_2 = \frac{1}{2}d \log \frac{|D_i|}{d}$ ,  $K_3 = \log\left(\frac{6\sigma_2^2 r \sqrt{|D_i|}}{\gamma_2 \sigma_1 \delta}\right)$ , and  $|D_i|$  is size of the local dataset.

Suppose Assumption 2 - Assumption 5 hold, and  $\omega \subset \{\mathbf{w} : \|\mathbf{w} - \mathbf{w}^*\| \leq r\sqrt{d}\}$  for some positive parameter  $r$ . For any  $\delta \in (0, 1)$  and any integer  $|D_i|$ , if  $\Delta_1 \leq \sigma_1^2/\gamma_1$  and  $\Delta_2 \leq \sigma_2^2/\gamma_2$ , we have:

$$\begin{aligned} \mathbb{P}\{\forall \mathbf{w} \in \omega : \left\| \frac{1}{|D_i|} \sum_{X_j \in D_i} \nabla f(X_j, \mathbf{w}) - \nabla F(\mathbf{w}) \right\| \\ \leq 8\Delta_2\|\mathbf{w} - \mathbf{w}^*\| + 4\Delta_1\} \geq 1 - \delta. \end{aligned}$$

*Proof:* Our proof is mainly based on the  $\varepsilon$ -net argument [1], [3]. We let  $\tau = \frac{\gamma_2 \sigma_1}{2\sigma_2^2} \sqrt{\frac{d}{|D_i|}}$  and  $\ell^* = \lceil r\sqrt{d}/\tau \rceil$ . For any integer  $1 \leq \ell \leq \ell^*$ , we let  $\omega_\ell \triangleq \{\mathbf{w} : \|\mathbf{w} - \mathbf{w}^*\| \leq r\sqrt{d}\}$ . Given an integer  $\ell$ , we let  $\mathbf{w}_1, \dots, \mathbf{w}_{N_{\varepsilon_\ell}}$  be an  $\varepsilon$ -cover of  $\omega_\ell$ , where  $\varepsilon_\ell = (\sigma_2 \tau \ell \sqrt{d}/|D_i|)/(M \vee M')$ , and  $M \vee M' = \max\{M, M'\}$ . We know  $\log N_{\varepsilon_\ell} \leq d \log \left(\frac{3\tau \ell}{\varepsilon_\ell}\right)$  from [2]. For any  $\mathbf{w} \in \omega$ , there exists a  $k_\ell$  ( $1 \leq k_\ell \leq N_{\varepsilon_\ell}$ ) such that  $\|\mathbf{w} - \mathbf{w}_{k_\ell}\| \leq \varepsilon_\ell$ . By triangle's inequality, we have:

$$\begin{aligned} \left\| \frac{1}{|D_i|} \sum_{X_j \in D_i} \nabla f(X_j, \mathbf{w}) - \nabla F(\mathbf{w}) \right\| &\leq \|\nabla F(\mathbf{w}) - \nabla F(\mathbf{w}_{k_\ell})\| \\ &+ \left\| \frac{1}{|D_i|} \sum_{X_j \in D_i} (\nabla f(X_j, \mathbf{w}) - \nabla f(X_j, \mathbf{w}_{k_\ell})) \right\| \\ &+ \left\| \frac{1}{|D_i|} \sum_{X_j \in D_i} \nabla f(X_j, \mathbf{w}_{k_\ell}) - \nabla F(\mathbf{w}_{k_\ell}) \right\|. \end{aligned} \quad (19)$$

In view of Assumption 1, we have:

$$\|\nabla F(\mathbf{w}) - \nabla F(\mathbf{w}_{k_\ell})\| \leq M\|\mathbf{w} - \mathbf{w}_{k_\ell}\| \leq M\varepsilon_\ell. \quad (20)$$

We define event

$$\mathcal{E}_1 = \left\{ \sup_{\mathbf{w}, \mathbf{w}' \in \omega: \mathbf{w} \neq \mathbf{w}'} \frac{\|\nabla \bar{f}_{|D_i|}(\mathbf{w}) - \nabla \bar{f}_{|D_i|}(\mathbf{w}')\|}{\|\mathbf{w} - \mathbf{w}'\|} \leq M' \right\}, \quad (21)$$

where  $\nabla \bar{f}_{|D_i|}(\mathbf{w}) = (\sum_{X_j \in D_i} \nabla f(X_j, \mathbf{w})) / |D_i|$ .

By Assumption 4, we have  $\mathbb{P}\{\mathcal{E}_1\} \geq 1 - \delta/3$ . On event  $\mathcal{E}_1$ , we have the following:

$$\sup_{\mathbf{w}, \mathbf{w}' \in \omega: \mathbf{w} \neq \mathbf{w}'} \left\| \frac{1}{|D_i|} \sum_{X_i \in D_i} (f(X_i, \mathbf{w}) - f(X_i, \mathbf{w}_{k_\ell})) \right\| \leq M' \varepsilon_\ell. \quad (22)$$

By triangle's inequality, we have:

$$\begin{aligned} & \left\| \frac{1}{|D_i|} \sum_{X_j \in D_i} \nabla f(X_j, \mathbf{w}_{k_\ell}) - \nabla F(\mathbf{w}_{k_\ell}) \right\| \\ & \leq \left\| \frac{1}{|D_i|} \sum_{X_j \in D_i} \nabla f(X_j, \mathbf{w}^*) - \nabla F(\mathbf{w}^*) \right\| \\ & + \left\| \frac{1}{|D_i|} \sum_{X_j \in D_i} (\nabla f(X_j, \mathbf{w}_{k_\ell}) - \nabla f(X_j, \mathbf{w}^*)) \right. \\ & \quad \left. - (\nabla F(\mathbf{w}_{k_\ell}) - \nabla F(\mathbf{w}^*)) \right\| \\ & \stackrel{(a)}{\leq} \left\| \frac{1}{|D_i|} \sum_{X_j \in D_i} \nabla f(X_j, \mathbf{w}^*) - \nabla F(\mathbf{w}^*) \right\| \\ & + \left\| \frac{1}{|D_i|} \sum_{X_j \in D_i} h(X_j, \mathbf{w}_{k_\ell}) - \mathbb{E}[h(X, \mathbf{w}_{k_\ell})] \right\|, \quad (23) \end{aligned}$$

where (a) is because Equations 1 and 2.

We define events as:

$$\mathcal{E}_2 = \left\{ \left\| \frac{1}{|D_i|} \sum_{X_j \in D_i} \nabla f(X_j, \mathbf{w}^*) - \nabla F(\mathbf{w}^*) \right\| \leq 2\Delta_1 \right\}, \quad (24)$$

$$\mathcal{F}_\ell = \left\{ \sup_{1 \leq k \leq N_\ell} \left\| \frac{1}{|D_i|} \sum_{X_j \in D_i} h(X_j, \mathbf{w}_k) - \mathbb{E}[h(X, \mathbf{w}_k)] \right\| \leq 2\tau\ell\Delta_2 \right\}. \quad (25)$$

Since  $\Delta_1 \leq \sigma_1^2/\gamma_1$ , it follows from Lemma 3 that  $\mathbb{P}\{\mathcal{E}_2\} \geq 1 - \delta/3$ . For  $\Delta_2 \leq \sigma_2^2/\gamma_2$  from Lemma 4, we have:

$$\begin{aligned} \mathbb{P}\{\mathcal{F}_\ell^c\} &= \mathbb{P}\left\{ \sup_{1 \leq k \leq N_\ell} \left\| \frac{1}{|D_i|} \sum_{X_j \in D_i} h(X_j, \mathbf{w}_k) \right. \right. \\ & \quad \left. \left. - \mathbb{E}[h(X, \mathbf{w}_k)] \right\| > 2\tau\ell\Delta_2 \right\} \\ &\leq \sum_{k=1}^{N_\ell} \mathbb{P}\left\{ \left\| \frac{1}{|D_i|} \sum_{X_j \in D_i} h(X_j, \mathbf{w}_k) - \mathbb{E}[h(X, \mathbf{w}_k)] \right\| > 2\tau\ell\Delta_2 \right\} \\ &\leq \frac{\delta}{3\ell^*} \frac{1}{\left(\frac{3\tau\ell}{\varepsilon_\ell}\right)^d} \left(\frac{3\tau\ell}{\varepsilon_\ell}\right)^d = \frac{\delta}{3\ell^*}. \quad (26) \end{aligned}$$

In conclusion, by combining Equations 19, 20, 22, and 23, on event  $\mathcal{E}_1 \cap \mathcal{E}_2 \cap \mathcal{F}_\ell$ , we have:

$$\begin{aligned} & \sup_{\mathbf{w} \in \omega_\ell} \left\| \frac{1}{|D_i|} \sum_{X_j \in D_i} \nabla f(X_j, \mathbf{w}) - \nabla F(\mathbf{w}) \right\| \\ & \leq (M + M')\varepsilon_\ell + 2\Delta_1 + 2\Delta_2\tau\ell \\ & \stackrel{(a)}{\leq} 4\Delta_2\tau\ell + 2\Delta_1, \quad (27) \end{aligned}$$

where (a) is due to  $(M \vee M')\varepsilon_\ell \leq \Delta_2\tau\ell$ . We let event  $\mathcal{E} = \mathcal{E}_1 \cap \mathcal{E}_2 \cap (\cap_{\ell=1}^{\ell^*} \mathcal{F}_\ell)$ . By the union bound, we have  $\mathbb{P}\{\mathcal{E}\} \geq 1 - \delta$ . Moreover, suppose event  $\mathcal{E}$  holds. For any  $\mathbf{w} \in \omega_\ell$ , there exists an  $1 \leq \ell \leq \ell^*$  such that  $(\ell - 1)\tau < \|\mathbf{w} - \mathbf{w}^*\| \leq \ell\tau$ . If  $2 \leq \ell \leq 2(\ell - 1)$ , we have:

$$\begin{aligned} & \left\| \frac{1}{|D_i|} \sum_{X_j \in D_i} \nabla f(X_j, \mathbf{w}) - \nabla F(\mathbf{w}) \right\| \leq 4\Delta_2\tau\ell + 2\Delta_1 \\ & \leq 8\Delta_2\|\mathbf{w} - \mathbf{w}^*\| + 2\Delta_1. \quad (28) \end{aligned}$$

If  $\ell = 1$ , we have:

$$\left\| \frac{1}{|D_i|} \sum_{X_j \in D_i} \nabla f(X_j, \mathbf{w}) - \nabla F(\mathbf{w}) \right\| \leq 4\Delta_2\tau\ell + 2\Delta_1 \stackrel{(a)}{\leq} 4\Delta_1, \quad (29)$$

where (a) is due to that  $\Delta_2 \leq \sigma_2^2/\gamma_2$  and  $\Delta_1 \geq \sigma_1\sqrt{d/|D_i|}$ . Combining inequalities 5 and 29, we have:

$$\begin{aligned} & \sup_{\mathbf{w} \in \omega_{\ell^*}} \left\| \frac{1}{|D_i|} \sum_{X_j \in D_i} \nabla f(X_j, \mathbf{w}) - \nabla F(\mathbf{w}) \right\| \\ & \leq 8\Delta_2\|\mathbf{w} - \mathbf{w}^*\| + 4\Delta_1. \quad (30) \end{aligned}$$

The proposition follows by the Assumption that  $\omega \subset \omega_{\ell^*}$ . In addition, we let  $\mathbf{g}_{median} = 1/|D_i| \sum_{X_j \in D_i} \nabla f(X_j, \mathbf{w})$ , that is because the median update selected by the server at each epoch is a local update uploaded by a certain client.

**Proof of Theorem 1:** With the help of the above lemma, we can prove the theory that the difference between the aggregated result of the global model at epoch  $t$  and the optimal solution

is bounded. We have:

$$\begin{aligned}
& \|\mathbf{w}^t - \mathbf{w}^*\| \\
&= \|\mathbf{w}^{t-1} - \alpha \mathbf{g}^{t-1} - \mathbf{w}^*\| \\
&= \|\mathbf{w}^{t-1} - \alpha \nabla F(\mathbf{w}^{t-1}) - \mathbf{w}^* + \alpha \nabla F(\mathbf{w}^{t-1}) - \alpha \mathbf{g}^{t-1}\| \\
&\leq \|\mathbf{w}^{t-1} - \alpha \nabla F(\mathbf{w}^{t-1}) - \mathbf{w}^*\| + \alpha \|\nabla F(\mathbf{w}^{t-1}) - \mathbf{g}^{t-1}\| \\
&\stackrel{(a)}{=} \|\mathbf{w}^{t-1} - \alpha \nabla F(\mathbf{w}^{t-1}) - \mathbf{w}^*\| + \alpha \|\mathbf{g}^{t-1} - \nabla F(\mathbf{w}^{t-1})\| \\
&\stackrel{(b)}{\leq} \|\mathbf{w}^{t-1} - \alpha \nabla F(\mathbf{w}^{t-1}) - \mathbf{w}^*\| + 2\alpha \|\nabla F(\mathbf{w}^{t-1})\| \\
&\quad + 3\alpha \|\mathbf{g}_{median}^{t-1} - \nabla F(\mathbf{w}^{t-1})\| \\
&\stackrel{(c)}{=} \underbrace{\|\mathbf{w}^{t-1} - \alpha \nabla F(\mathbf{w}^{t-1}) - \mathbf{w}^*\|}_{A_1} + 2\alpha \underbrace{\|\nabla F(\mathbf{w}^{t-1}) - \nabla F(\mathbf{w}^*)\|}_{A_2} \\
&\quad + 3\alpha \underbrace{\|\mathbf{g}_{median}^{t-1} - \nabla F(\mathbf{w}^{t-1})\|}_{A_3} \\
&\stackrel{(d)}{\leq} \sqrt{1 - L^2/(4M^2)} \|\mathbf{w}^{t-1} - \mathbf{w}^*\| + 2\alpha M \|\mathbf{w}^{t-1} - \mathbf{w}^*\| \\
&\quad + 3\alpha (8\Delta_2 \|\mathbf{w}^{t-1} - \mathbf{w}^*\| + 4\Delta_1) \\
&= (\sqrt{1 - L^2/(4M^2)} + 24\alpha\Delta_2 + 2\alpha M) \|\mathbf{w}^{t-1} - \mathbf{w}^*\| + 12\alpha\Delta_1, \\
&\tag{31}
\end{aligned}$$

where (a) has the same reason as Lemma 1(a); (b) is obtained according to Lemma 1, (c) is due to  $\nabla F(\mathbf{w}^*) = 0$ ; and  $A_1$ ,  $A_2$ , and  $A_3$  in (c) are Lemma 2, Assumption 1, and Lemma 5, respectively.

By letting  $\rho = 1 - (\sqrt{1 - L^2/(4M^2)} + 24\alpha\Delta_2 + 2\alpha M)$ , we have:

$$\begin{aligned}
& \|\mathbf{w}^t - \mathbf{w}^*\| \\
&\leq (1 - \rho) \|\mathbf{w}^{t-1} - \mathbf{w}^*\| + 12\alpha\Delta_1 \\
&\leq (1 - \rho) [(1 - \rho) \|\mathbf{w}^{t-2} - \mathbf{w}^*\| + 12\alpha\Delta_1] + 12\alpha\Delta_1 \\
&= (1 - \rho)^2 \|\mathbf{w}^{t-2} - \mathbf{w}^*\| + 12[1 + (1 - \rho)]\alpha\Delta_1 \\
&\leq (1 - \rho)^2 [(1 - \rho) \|\mathbf{w}^{t-3} - \mathbf{w}^*\| + 12\alpha\Delta_1] + 12[1 + (1 - \rho)]\alpha\Delta_1 \\
&= (1 - \rho)^3 \|\mathbf{w}^{t-3} - \mathbf{w}^*\| + 12[1 + (1 - \rho) + (1 - \rho)^2]\alpha\Delta_1 \\
&\dots \\
&\leq 12[1 + (1 - \rho) + (1 - \rho)^2 + \dots + (1 - \rho)^{t-1}]\alpha\Delta_1 \\
&\quad + (1 - \rho)^t \|\mathbf{w}^0 - \mathbf{w}^*\| \\
&= (1 - \rho)^t \|\mathbf{w}^0 - \mathbf{w}^*\| + 12 \frac{1 - (1 - \rho)^t}{1 - (1 - \rho)} \alpha\Delta_1 \\
&= (1 - \rho)^t \|\mathbf{w}^0 - \mathbf{w}^*\| + \frac{12\alpha\Delta_1}{\rho} - \frac{12(1 - \rho)^t \alpha\Delta_1}{\rho} \\
&\leq (1 - \rho)^t \|\mathbf{w}^0 - \mathbf{w}^*\| + \frac{12\alpha\Delta_1}{\rho}. \tag{32}
\end{aligned}$$

Thus, we conclude the proof.

## REFERENCES

- [1] Y. Chen, L. Su, and J. Xu, "Distributed statistical machine learning in adversarial settings: Byzantine gradient descent," *Proc. ACM Meas. Anal. Comput. Syst.*, 2017.
- [2] R. Vershynin, "Introduction to the non-asymptotic analysis of random matrices," *arXiv preprint arXiv:1011.3027*, 2010.
- [3] X. Cao, M. Fang, J. Liu, and N. Z. Gong, "Fltrust: Byzantine-robust federated learning via trust bootstrapping," in *ISOC Network and Distributed System Security Symposium*, 2021.