

GENERALIZATION OF FEDAVG UNDER CONSTRAINED POLYAK-ŁOJASIEWICZ TYPE CONDITIONS: A SINGLE HIDDEN LAYER NEURAL NETWORK ANALYSIS

Anonymous authors

Paper under double-blind review

ABSTRACT

In this work, we study the optimization and the generalization performance of the widely used FedAvg algorithm for solving Federated Learning (FL) problems. We analyze the generalization performance of FedAvg by handling the optimization error and the Rademacher complexity. Towards handling optimization error, we propose novel constrained Polyak-Łojasiewicz (PL)-type conditions on the objective function that ensure existence of a global optimal to which FedAvg converges linearly after $\mathcal{O}(\log(1/\epsilon))$ rounds of communication, where ϵ is the desired optimality gap. Importantly, we demonstrate that a class of single hidden layer neural networks satisfies the proposed constrained PL-type conditions required to establish the linear convergence of FedAvg as long as $m > nK/d$, where m is the width of the neural network, K is the number of clients, n is the number of samples at each client, and d is the feature dimension. We then bound the Rademacher complexity for this class of neural networks and establish that both Rademacher complexity and the generalization error of FedAvg decrease at an optimal rate of $\mathcal{O}(1/\sqrt{n})$. We further show that increasing the number of clients K decreases the generalization error at the rate of $\mathcal{O}(1/\sqrt{n} + 1/\sqrt{nK})$.

1 INTRODUCTION

Federated learning (FL) is a distributed learning paradigm where multiple client devices collaborate with the help of a server to solve a joint problem while keeping the data of each client private (Kairouz et al., 2021). A typical FL problem aims to solve $\min_{\mathbf{w}} \sum_{k=1}^K \Phi_k(\mathbf{w})$, where $\Phi_k(\mathbf{w})$ is the loss at the k^{th} client and \mathbf{w} refers to the joint model the clients aim to learn. A standard and most widely adopted algorithm to solve the FL problem is the Federated Averaging (FedAvg) algorithm first proposed in (McMahan et al., 2017). Consequently, the study of the convergence performance of FedAvg has received wide attention (Konečný et al., 2015; Stich, 2018; McMahan et al., 2017; Li et al., 2020; Zhou & Cong, 2017b). However, when it comes to ensuring generalization guarantees for FedAvg, the problem has not received significant attention, partially because of the challenging nature of the problem (Mohri et al., 2019; Sun et al., 2023; Hu et al., 2022). To prove the generalization guarantees for FedAvg, we need to bound (a) the **optimization error** (on empirical loss) achieved by FedAvg, and (b) the **complexity measure** such as the Rademacher complexity of the model (Arora et al., 2019; Mohri et al., 2019; 2018). The major challenge in guaranteeing good generalization performance is to bound both (a) and (b) above, which are often contradictory, i.e., proving optimization guarantees usually rely on restrictive assumptions on the loss landscape like (strong)-convexity or Polyak-Łojasiewicz (PL) inequality to be satisfied over the entire parameter space Haddadpour et al. (2019); Haddadpour & Mahdavi (2019) while the Rademacher complexity is large for an unbounded parameter space [see Theorem 5.10 (Mohri et al., 2018)]. Therefore, bounding both (a) and (b) simultaneously is challenging, thereby making it difficult to provide satisfactory generalization guarantees for FedAvg. To address these challenges in this work:

➤ We first analyze the convergence of FedAvg and establish **linear convergence** under a **new set of assumptions** that are only required to be satisfied **locally**. Importantly, to highlight the practicality of the assumptions, we establish that the proposed **assumptions are naturally satisfied by a single hidden-layer Neural Network (NN)**.

➤ We then study the generalization guarantees of FedAvg for the single hidden-layer NN and show

that the proposed local assumptions lead to a **Rademacher complexity that goes down with the number of samples n as $\mathcal{O}(1/\sqrt{n})$** . Specifically, our analysis captures the effects of local samples, the number of clients, and model sizes on the performance of the FedAvg algorithm.

In the following, we discuss specific challenges and the drawbacks of the current state-of-the-art with respect to challenges (a) and (b) discussed above.

Convergence of FedAvg. As discussed earlier, several works have analyzed the convergence performance of FedAvg under various settings. In the non-convex regime, multiple works have established the convergence of FedAvg to a stationary point (local optimal) (Konečný et al., 2015; Stich, 2018; McMahan et al., 2017; Li et al., 2020; Zhou & Cong, 2017b). However, the local optimal does not guarantee a small empirical loss, and hence cannot be used to provide generalization guarantees. Some works have shown convergence of FedAvg to global optimal but under restrictive assumptions of (strong) convexity (Stich, 2018; Qu et al., 2020). In Haddadpour et al. (2019), the authors provide convergence of FedAvg to the global optimal by imposing the PL condition on the objective function, which is unfortunately not satisfied by several loss functions (e.g., log-logistic loss) over the whole parameter space. Importantly, assuming that the PL inequality is satisfied globally (without any restriction on the parameter space Haddadpour & Mahdavi (2019)) leads to a large Rademacher complexity, thus leading to worse generalization guarantees. This leads to the following question:

Q1: *Can we develop conditions that are satisfied locally (on a restricted parameter space) rather than globally and provide convergence guarantees for FedAvg? Are there models that satisfy such a condition?*

To address Q1, we provide *new weaker* conditions (a constrained variant of the PL-inequality) on the global and local loss functions. Importantly, we prove that there exists a globally optimal point within a ball of radius ρ around initialization to which FedAvg converges linearly. Moreover, we also establish that there exist NN architectures that satisfy the conditions proposed in our work.

Generalization guarantees for FedAvg: The generalization performance of centralized machine learning algorithms has been extensively studied (Mohri et al., 2018; Bousquet & Elisseeff, 2002; Emami et al., 2020). However, the study of generalization guarantees of FL algorithms is rather limited (Mohri et al., 2019; Hu et al., 2022; Yuan et al., 2021a). Notably, these studies often overlook the impact of the optimization algorithm Sun et al. (2023), and often rely on assumptions like Binary loss Hu et al. (2022); Mohri et al. (2019) and the Bernstein condition (Yuan et al., 2021a). Additionally, generalization bounds for meta-learning and FL are established in Fallah et al. (2021); Chen et al. (2021) under stringent assumptions such as strong convexity and bounded loss functions. Recently, Sun et al. (2023) has investigated the generalization of FedAvg via the lens of uniform stability. We note that these analyses impose strong assumptions such as bounded gradient and heterogeneity on the data, which are usually not satisfied by many problems of practical interest. Moreover, the optimization guarantees provided in Sun et al. (2023) are weaker compared to the linear convergence established in our work. Based on the above observations, we ask the following main question:

Q2: *Can we provide generalization guarantees for FedAvg? If so, what is the impact of (a) the number of samples per client, (b) the model size, and (c) the number of clients on the generalization performance?*

We address Q2 by deriving Rademacher complexity when each client employs a single hidden-layer NN for FedAvg implementation. We show that the local assumptions developed to address Q1 play an important role in bounding the Rademacher complexity for FedAvg. Importantly, our analysis captures the effect of data samples and NN size, and the number of clients on the generalization performance of FedAvg. It is worth mentioning that to address both Q1 and Q2, we *do not* make some standard assumptions that are typically used in many existing works Li et al. (2019); Stich (2018); Yu et al. (2019); Haddadpour et al. (2019); Qu et al. (2020); Woodworth et al. (2020a;b); Hu et al. (2022); Mohri et al. (2019) such as: (i) (strongly) convex loss, (ii) bounded loss, (iii) bounded gradients (iv) bounded heterogeneity, and (v) interpolation¹. In this work, we have not assumed the existence of a global optimal point; rather, it is part of our conclusion.

¹Interpolation refers to the existence of a \mathbf{w}^* such that $\Phi_{k,i}(\mathbf{w}^*) = 0$ for all $k \in [K]$ and $i \in [n]$.

Contributions. The major contributions of our work include:

➤ **Answer to Q1:** For the first time, we show that FedAvg converges linearly to the optimal solution (see Corollary 3.2) if the local loss functions at each client and the global loss function satisfy a novel local PL-type assumption introduced in Assumption 2.4. It is important to note that the existence of a global optimal in our analysis is a part of our conclusion, *not* an assumption. To the best of our knowledge, both conditions introduced in Assumption 2.4 are new. It is also worth noting that these conditions do not follow from any of the existing results, even in the special case of centralized setting, i.e., for $K = 1$ (Chatterjee, 2022; An & Lu, 2023). In addition, we also establish that a single hidden-layer NN satisfies the two conditions proposed in Assumption 2.4. Specifically, we establish the conditions on the width of the NN as a function of the number of samples, number of clients, and the feature dimension, and on the eigenvalues of the Jacobian of the loss functions (or the scaling factor of the final output layer) such that the proposed conditions are satisfied. To our knowledge, these results are novel (see Theorems 4.5).

➤ **Answer to Q2:** To address Q2, we derive an upper bound on the Rademacher complexity for a class of single hidden layer NNs by utilizing the fact that the FedAvg iterates stay within a ρ -ball around the initialization. We point out that this is made possible by the conditions provided in Assumption 2.4. In particular, we show that the Rademacher complexity approaches zero if the radius $\rho = \mathcal{O}(\sqrt{n})^2$ and $m = \mathcal{O}(n^3)$, where n is the number of samples at each client and m is the width of the NN. We show that the generalization error regardless of the data heterogeneity diminishes as $\mathcal{O}(1/\sqrt{n})$. We finally corroborate our theoretical findings through numerical experiments.

2 FEDAVG: ALGORITHM AND ASSUMPTIONS

As discussed in Section 1, FL aims to solve the following optimization problem:

$$\min_{\mathbf{w}} \left\{ \Phi(\mathbf{w}) := \frac{1}{K} \sum_{k=1}^K \Phi_k(\mathbf{w}) \right\}, \quad (1)$$

where $\Phi_k(\mathbf{w}) := \mathbb{E}_{(\mathbf{x}, y) \sim \mathcal{D}_k} l_k(f_{\mathbf{w}}(\mathbf{x}), y)$ is the loss function at client $k \in [K]$. Here, $y \in \mathcal{Y}$ is the true label, and $f_{\mathbf{w}}(\mathbf{x})$ is the output of model $\mathbf{w} \in \mathbb{R}^{d'}$ for an input feature $\mathbf{x} \in \mathbb{R}^d$, and $l_k : \mathcal{Y} \times \mathcal{Y} \rightarrow \mathbb{R}^+$ is the loss function at the client $k \in [K]$. In the above, d' is the dimension of the parameter space. The following algorithm captures the main steps of FedAvg (McMahan et al., 2017). In Algorithm 1, $\Phi_{k,i}(\mathbf{w}_k^{r,t})$ denotes the empirical loss function at client $k \in [N]$ computed using sample $i \in [n]$.

In this and the subsequent section, we answer Q1 posed in Sec. 1. In particular, we provide a general condition for the above algorithm to converge to a global optimum and for the model parameters to stay within a closed ball of radius ρ . In the later sections, we show that this condition is, in fact, satisfied for a single hidden layer NN. Specifically, this constraint imposes a natural regularization of the NN which provides better generalization, as discussed later. To prove our claim, we make the following standard assumptions on the loss function Ji & Telgarsky (2018).

Assumption 2.1 (*L*-Smoothness). *The loss functions Φ_k and Φ are assumed to be L_k -smooth and L -smooth, respectively, i.e., $\|\nabla \Phi_k(\mathbf{u}) - \nabla \Phi_k(\mathbf{v})\| \leq L_k \|\mathbf{u} - \mathbf{v}\|$ for all $k \in [K]$ and $\|\nabla \Phi(\mathbf{u}) - \nabla \Phi(\mathbf{v})\| \leq L \|\mathbf{u} - \mathbf{v}\|$ for all \mathbf{u} and \mathbf{v} .*

Assumption 2.2 (Samplewise Smoothness). *The loss functions $\Phi_{k,i}(\mathbf{w})$ are assumed to be $l_{k,i}$ -sample-wise smooth, i.e., $\|\nabla \Phi_{k,i}(\mathbf{v})\|^2 \leq 2l_{k,i} \Phi_{k,i}(\mathbf{v})$ for all $k \in [K]$ and $i \in [n]$.*

To define the major assumptions required for the convergence of FedAvg Algorithm 1, we need the following definition (Chatterjee, 2022).

Definition 2.3. *Let $f : \mathbb{R}^d \rightarrow \mathbb{R}^+$ be continuously differentiable function on closed ball $\mathbb{B}[\underline{\mathbf{w}}^0, \rho]$ with center at initialization $\underline{\mathbf{w}}^0 \in \mathbb{R}^d$ and radius $\rho > 0$. Define*

$$\alpha(\underline{\mathbf{w}}^0, \rho) := \inf_{\mathbf{w} \in \mathbb{B}[\underline{\mathbf{w}}^0, \rho]} \frac{\|\nabla f(\mathbf{w})\|^2}{f(\mathbf{w})} > 0. \quad (2)$$

Next, we state an important assumption that leads to linear convergence within a ball around initialization.

²This is the radius over which our new condition should be satisfied.

Algorithm 1 FedAvg McMahan et al. (2017)

```

1: Initialize:  $\{\mathbf{w}_k^{0,0} = \underline{\mathbf{w}}^0\}, \mathbf{w}_k \in \mathbb{R}^d$  for  $k = 1, 2, \dots, K$ 
2: for  $r = 0, 1, \dots, R - 1$  do
3:   Broadcast  $\underline{\mathbf{w}}^r$  to all the clients  $k \in [K]$ 
4:   for  $\tau = 0, 1, \dots, T - 1$  do
5:     for each client  $k \in [K]$  do
6:       Sample a batch  $\mathcal{B}_k^{r,t}$  of size  $|\mathcal{B}_k^{r,t}| = b$ 
7:       SGD step on  $\mathbf{w}_k^{r,t}$  for  $k \in [K]$ :
8:          $\mathbf{w}_k^{r,t+1} = \mathbf{w}_k^{r,t} - \eta \widehat{\nabla} \Phi_k(\mathbf{w}_k^{r,t})$ 
9:          $\widehat{\nabla} \Phi_k(\mathbf{w}_k^{r,t}) := \frac{1}{b} \sum_{i \in \mathcal{B}_k^{r,t}} \nabla \Phi_{k,i}(\mathbf{w}_k^{r,t})$ 
10:    end for
11:  end for
12:  Receive  $\mathbf{w}_k^{r,T}$  from nodes  $k \in [K]$ 
13:  Aggregation step :  $\underline{\mathbf{w}}^{r+1} = \frac{1}{K} \sum_{k \in [K]} \mathbf{w}_k^{r,T}$ 
14: end for

```

Assumption 2.4. For some initialization $\underline{\mathbf{w}}^0$ and radius $\rho > 0$, we make the following assumptions on the local and global loss functions:

1. The loss function at each client is assumed to satisfy (see Theorem E.1)

$$32\Phi_k(\underline{\mathbf{w}}^0) < \rho^2 \alpha_k(\underline{\mathbf{w}}^0, \rho). \quad (3)$$

Here, $\alpha_k(\underline{\mathbf{w}}^0, \rho)$ is as defined in equation 2 but with $f(\cdot)$ replaced by $\Phi_k(\cdot)$.

2. The global loss function is assumed to satisfy the following condition

$$\sqrt{128e l'_{\max} K \Phi(\underline{\mathbf{w}}^0)} < (1 - \zeta_\rho) \rho \alpha_g(\underline{\mathbf{w}}^0, \rho), \quad (4)$$

for some $\zeta_\rho \in (0, 1)$. Here, $\alpha_g(\underline{\mathbf{w}}^0, \rho)$ is as defined in equation 2 but with $f(\cdot)$ replaced by $\Phi(\cdot)$.

Remark 1. In general, two very critical assumptions are made in the literature while proving linear convergence: (i) interpolation, i.e., there exists \mathbf{w}^* such that $\Phi_i(\mathbf{w}^*) = 0$ for all samples $i \in [n]$ Liu et al. (2022); Li et al. (2019), and (ii) strongly convex loss Li et al. (2019); Karimireddy et al. (2020) or loss function satisfying the PL-inequality Fan et al. (2023). Later, a relaxed version of PL-inequality called local PL or PL^* -inequality was proposed where the PL-inequality needs to be satisfied over a small ball around the initialization (see Liu et al. (2022); Oymak & Soltanolkotabi (2019)). Despite this relaxation, it makes a critical assumption on the existence of the optimal \mathbf{w}^* such that the loss $\Phi_i(\mathbf{w}^*) = 0$ for all samples $i \in [n]$ -the interpolation regime. In our work, we argue that this assumption can be relaxed with our novel condition shown in Assumption 2.4. It is important to note that our condition is fundamentally different from the PL^* -inequality in the following way:

- There is a stark difference between our proposed condition and the the PL-condition (or PL^* condition), which is defined as $\|\nabla \Phi(\mathbf{w})\|^2 \geq \mu(\Phi(\mathbf{w}) - \Phi(\mathbf{w}^*))$ for all $\mathbf{w} \in \mathbb{R}^d$ (and $\mathbf{w} \in \mathbb{B}[\underline{\mathbf{w}}^0, \rho]$ for PL^* condition). In the PL-condition (and local PL), the constants do not depend on the initialization and radius as the condition is universally satisfied. Another important assumption made in the local/global PL-condition is the existence of a global optimal point \mathbf{w}^* . In contrast, our proposed condition does not require this assumption; instead, we prove the existence of a global optimal point under our novel condition.
- It is important to note that the PL-condition must be satisfied over the entire parameter space, which can restrict its applicability to certain loss functions such as logistic loss Karimi et al. (2016). On the other hand, our novel condition is assumed only over a small neighborhood around the initialization, making it more broadly applicable. Later we show that parameters such as initialization and the radius ρ can be chosen so that the condition is easily (compared to the PL inequality) satisfied.

In this work, we have shown that the proposed condition is satisfied for at least a single hidden layer neural network. In Chatterjee (2022), the authors have shown that the wide neural network satisfies

the constrained PL inequality for a single client setting. Therefore, we strongly believe that the proposed condition in our work will also be satisfied for wide neural networks.

3 CONVERGENCE ANALYSIS

In this section, we establish that the FedAvg Algorithm 1 achieves linear convergence to a global optimum under the set of assumptions introduced in Sec. 2. Importantly, note that the existence of this global optimum is established as a conclusion rather than an assumption. Moreover, unlike other works, we do not explicitly assume interpolation to establish linear convergence of FedAvg (Haddadpour et al., 2019; Stich, 2018). In particular, we establish a proof that the sufficient conditions stated in equation 2.4 not only guarantee the linear convergence of Algorithm 1 but also ensure the existence of an optimal point denoted as \mathbf{w}^* within the closed ball $\mathbb{B}[\mathbf{w}^0, \rho]$. The following theorem is a precise statement whose proof can be found in Appendix 3.1.

Theorem 3.1. *Assuming that there exists an initialization $\mathbf{w}^0 \in \mathbb{R}^d$, and a radius $\rho > 0$ such that Assumptions 2.1 and 2.4 are satisfied by loss functions Φ and Φ_k for $k \in [K]$, then FedAvg ensures that there exists a $\mathbf{w}^* \in \mathbb{B}[\mathbf{w}^0, \rho]$ such that $\lim_{R \rightarrow \infty} \Phi(\mathbf{w}^R) = \Phi(\mathbf{w}^*) = 0$ provided the learning rate*

$$\eta \leq \min \left\{ \frac{2}{\alpha_{\min}}, \frac{\alpha_{\min}}{4L_{\max}l'_{\max}}, \frac{\alpha_{\min}}{2L_{\max}l'_{\max}}, \frac{1}{T\sqrt{\Psi_0}}, \frac{8}{\alpha_g T}, \frac{\zeta_\rho \rho}{T\sqrt{\Psi_0}}, \Psi_1, \Psi_2 \right\},$$

where $l'_{\max} := \max_k l'_k := \max_i l_{k,i}$; $L_{\max} := \max_k L_k$; $\alpha_{\min} := \min_{k \in [K]} \alpha_k$; $\Psi_0 := 2el'_{\max} K \Phi(\mathbf{w}^0)$; $\Psi_1 := \sqrt{\frac{3}{L_{\max}l'_{\max}}}$ and $\Psi_2 := \min \left\{ \frac{\alpha_g \alpha_{\min}}{4T(4L_{\max}^2 l'_{\max} + L'_{\max} \alpha_{\min})}, \frac{1}{3L_{\max} T} \right\}$. More precisely, after $R > 0$ communication rounds, the FedAvg Algorithm 1 satisfies

$$\Phi(\mathbf{w}^R) \leq \left(1 - \frac{\eta T \alpha_g (\mathbf{w}^0, \rho)}{4} \right)^R \Phi(\mathbf{w}^0). \quad (5)$$

Essence of the Proof of Theorem 3.1: Assumptions 2.1 and 2.4 lead to an exponential relation, specifically $\Phi(\mathbf{w}^{r+1}) \leq \gamma^r \Phi(\mathbf{w}^0)$, where $\gamma \in (0, 1)$, (refer to Lemma F.4). To prove the existence of global optima \mathbf{w}^* within the ball $\mathbb{B}[\mathbf{w}^0, \rho]$, we have used the method of induction on two variables: global communication round r and local updates t . By doing so, we conclude that the sequence $\{\mathbf{w}_k^{r,\tau}\}_{r,\tau \geq 0}$ remains confined within the ball $\mathbb{B}[\mathbf{w}^0, \rho]$ (refer to Lemma F.6), which ensures that the sequence $\{\mathbf{w}^r\}_{r=1}^\infty$ remains within the ball $\mathbb{B}[\mathbf{w}^0, \rho]$ for all r . Further, we have shown that the sequence $\{\mathbf{w}^r\}_{r=1}^\infty$ is Cauchy sequence in the closed subset $\mathbb{B}[\mathbf{w}^0, \rho]$ of complete space. Therefore, it guarantees the limit of the sequence $\{\mathbf{w}^r\}_{r=1}^\infty$, denoted by \mathbf{w}^* belongs to the ball. A complete proof is provided in Appendix F. \square

Note that Chatterjee (2022) required one condition to be satisfied for the linear convergence since their work considered a centralized setting. In contrast, our work requires two conditions for both global and local loss functions as stated in Assumptions 2.4 to guarantee linear convergence of FedAvg. Later we show that as the number of clients, K , increases, the requirement becomes more stringent. The above theorem leads to the following corollary.

Corollary 3.2. *By choosing η as in Theorem 3.1, for any error $\epsilon > 0$, Algorithm 1 achieves a loss of $\Phi(\mathbf{w}^R) < \epsilon$ after $R \geq \mathcal{O} \left(\left\lceil 2 \log \left(\frac{\Phi(\mathbf{w}^0)}{\epsilon} \right) \right\rceil \right)$ communication rounds.*

Our next goal is to show that it is possible to initialize a NN such that it satisfies the conditions provided in Assumption 2.4. However, note that this does not provide any guarantees on the generalization error. To fill this gap, in the following sections, we consider a single hidden-layer NN and show that (a) there exist an initialization and radius ρ such that it results in a linear convergence leading to zero training loss (i.e., assumptions stated in Sec. 2 are satisfied), and (b) prove that the generalization error can be made small by choosing large enough training samples and performing FedAvg for a sufficiently large number of communication rounds.

4 ASSUMPTION 2.4 FOR SINGLE HIDDEN LAYER NN WITH SQUARED ERROR LOSS

In this section, we show that there exist NNs such that Assumption 2.4 is satisfied, and hence leads to linear convergence of FedAvg (see Theorem 3.1). Towards this, we consider the following NN with a single hidden layer. In particular, we assume that the first layer has m neurons followed by a smooth activation function. The output of this NN is given by Arora et al. (2019)

$$f_{\mathbf{w}}(\mathbf{x}) = \frac{1}{\sqrt{m}} \sum_{j=1}^m v_j \sigma(\mathbf{w}_j^\top \mathbf{x}), \quad (6)$$

where $\mathbf{x} \in \mathbb{R}^d$ is the input feature vector. With a slight abuse of notation, we have used $\mathbf{w} = \text{vec}([\mathbf{w}_1, \mathbf{w}_2, \dots, \mathbf{w}_m]) \in \mathbb{R}^{dm \times 1}$ to denote the aggregated weight vectors in the first layer and $\mathbf{v} = (v_1, v_2, \dots, v_m)^\top$ to denote the weight in the second layer, where $v_j \stackrel{\text{i.i.d.}}{\sim} \{-1, 1\}$. Now, we make the following assumption on the activation function.

Assumption 4.1. We assume that $\sigma : \mathbb{R} \rightarrow \mathbb{R}$ is a smooth non-decreasing activation function such that $\sigma(0) = 0$. Further, first and second order derivatives of σ are bounded i.e., $|\sigma'(x)| \leq D_\sigma$ and $|\sigma''(x)| \leq \Delta_\sigma$.

Note that the above condition is satisfied by the tanh activation function, i.e., $\sigma(x) = \tanh(x)$. The condition $\sigma(0) = 0$ is assumed for the sake of simplicity and ease of notation. It turns out that, with random initialization, this can be relaxed without changing the main result of the paper. With $\sigma(x) \neq 0$, many activation functions such as Softmax, tanh to name a few (see Xu et al. (2015)) satisfy the conditions mentioned in Assumption 4.1. It is worth noting that the well-known ReLU activation does not satisfy the smoothness condition, but it can be well approximated by a smooth proxy function (see (Xu et al., 2015)).

Assumption 4.2. Each node $k \in [K]$ samples n i.i.d. data points denoted $\mathcal{X}_k = \{(\mathbf{x}_{k,1}, y_{k,1}), \dots, (\mathbf{x}_{k,n}, y_{k,n})\}$ from a continuous and possibly different distributions $p_k(\mathbf{x})$, $k \in [K]$ with $y_{k,i} \leq y_{\max}$ for all $i \in [n]$.

We consider the average loss function $\Phi(\mathbf{w}) := \frac{1}{K} \sum_{k=1}^K \Phi_k(\mathbf{w})$, where $\Phi_k : \mathbb{R}^{md} \rightarrow \mathbb{R}$ is the squared loss function for each client $k \in [K]$ and is defined as $\Phi_k(\mathbf{w}) = \sum_{i=1}^n [f_{\mathbf{w}}(\mathbf{x}_{k,i}) - y_{k,i}]^2 = \|\mathbf{e}_k\|_2^2$, where the i^{th} entry of the error vector $\mathbf{e}_k := [f_{\mathbf{w}}(\mathbf{x}_{k,i}) - y_{k,i}]$. Using $\mathbf{e} = [\mathbf{e}_1, \mathbf{e}_2, \dots, \mathbf{e}_K]$, the global loss can be written as $\Phi(\mathbf{w}) := \frac{1}{K} \|\mathbf{e}\|^2$. Next, we discuss the conditions under which a single hidden layer neural network satisfies Assumption 2.4. It turns out that these conditions are dependent on the following Jacobian matrix:

$$\mathbf{J}_k(\mathbf{w}) = \frac{1}{\sqrt{m}} \times \mathbf{H}_k(\mathbf{w}), \quad (7)$$

where each entry of $\mathbf{J}_k(\mathbf{w})$ is a d -dimensional row vector, and $\mathbf{H}_k(\mathbf{w})$ is defined as follows

$$\mathbf{H}_k(\mathbf{w}) := \begin{bmatrix} v_1 \sigma'(\mathbf{w}_1^\top \mathbf{x}_{k,1}) \mathbf{x}_{k,1}^\top & v_2 \sigma'(\mathbf{w}_2^\top \mathbf{x}_{k,1}) \mathbf{x}_{k,1}^\top & \dots & v_m \sigma'(\mathbf{w}_m^\top \mathbf{x}_{k,1}) \mathbf{x}_{k,1}^\top \\ \vdots & \vdots & \ddots & \vdots \\ v_1 \sigma'(\mathbf{w}_1^\top \mathbf{x}_{k,n}) \mathbf{x}_{k,n}^\top & v_2 \sigma'(\mathbf{w}_2^\top \mathbf{x}_{k,n}) \mathbf{x}_{k,n}^\top & \dots & v_m \sigma'(\mathbf{w}_m^\top \mathbf{x}_{k,n}) \mathbf{x}_{k,n}^\top \end{bmatrix}, \quad (8)$$

where $k \in [K]$ and the size of the matrix $\mathbf{H}_k(\mathbf{w})$ is $n \times md$, i.e., $\mathbf{H}_k(\mathbf{w}) \in \mathbb{R}^{n \times md}$. We define a global Jacobian matrix $\mathbf{J}(\mathbf{w})$ by stacking $\mathbf{H}_k^\top(\mathbf{w})$ row-wise as $\mathbf{J}(\mathbf{w}) = \frac{1}{\sqrt{m}} \times [\mathbf{H}_1^\top(\mathbf{w}), \mathbf{H}_2^\top(\mathbf{w}), \dots, \mathbf{H}_K^\top(\mathbf{w})] \in \mathbb{R}^{md \times Kn}$. The following lemma provides a condition under which $\mathbf{J}_k(\mathbf{w}^0)$ and $\mathbf{J}(\mathbf{w}^0)^\top$ are full rank matrices. Note that the full rank requirement is only at the initialization. The size of the NN scales as n/d as opposed to n in (Chatterjee, 2022). This result is similar to the results of Zhang et al. (2021) but for an FL setting.

Algorithm 2 FedAvg Algorithm for single hidden layer NN

- 1: **Initialization:** Initialize using $\underline{w}^0 \sim \mathcal{N}(\mathbf{0}, \frac{1}{d} I_{md \times md})$ and $v_i \stackrel{i.i.d.}{\sim} \{-1, 1\} \forall i \in [m]$.
- 2: Broadcast \underline{w}^r to all the clients $k \in [K]$
- 3: Run the FedAvg Algorithm 1

Lemma 4.3. *At the random initialization $\underline{w}^0 \sim \mathcal{N}(\mathbf{0}, \frac{1}{d} I_{md \times md})$, and $v_i \stackrel{i.i.d.}{\sim} \{-1, 1\}$ for all $i \in [m]$, the matrices $\mathbf{J}_k(\underline{w}^0)$ and $\mathbf{J}(\underline{w}^0)^\top$ have full column ranks almost surely provided $m \geq n/d$ and $m \geq nK/d$, respectively.*

Proof: The result follows by following the proof of Lemma E.1 of Zhang et al. (2021) for the matrices $\mathbf{H}_k(\underline{w}^0)$ and $\mathbf{H}(\underline{w}^0)^\top$. One main difference is that Zhang et al. (2021) uses mirrored Le-cun. However, the proof does not change for our initialization. \square

Towards stating the condition for neural network, we need the following definitions

$$\lambda_{k,\rho}^-(m) := \inf_{\mathbf{w} \in \mathbb{B}[\underline{w}^0, \rho]} \frac{\mathbf{e}_k^\top \mathbf{H}_k(\underline{w}^0) \mathbf{H}_k(\underline{w}^0)^\top \mathbf{e}_k}{\|\mathbf{e}_k\|^2}, \quad (9)$$

where \mathbf{e}_k and $\mathbf{H}_k(\mathbf{w})$ are as defined earlier.³ The following is an extension of the above definition to K clients

$$\lambda_\rho^-(m) := \inf_{\mathbf{w} \in \mathbb{B}[\underline{w}^0, \rho]} \frac{\mathbf{e}^\top \mathbf{H}(\underline{w}^0)^\top \mathbf{H}(\underline{w}^0) \mathbf{e}}{\|\mathbf{e}\|^2} \quad (10)$$

where $\mathbf{e} = [\mathbf{e}_1, \mathbf{e}_2, \dots, \mathbf{e}_K]^\top \in \mathbb{R}^{nK}$ and $\mathbf{H}(\underline{w}^0)$ is defined earlier. Similarly, $\tilde{\lambda}_{k,\rho}^-(m)$ and $\tilde{\lambda}_\rho^-(m)$ are defined by replacing $\mathbf{H}_k(\underline{w}^0)$ by $\mathbf{H}_k(\mathbf{w})$ and $\mathbf{H}(\underline{w}^0)$ by $\mathbf{H}(\mathbf{w})$ in equations 9 and equation 10, respectively. In addition, $\lambda_{\max}(\rho) := \sup_{\mathbf{w} \in \mathbb{B}(\underline{w}^0, \rho)} \lambda_{\max}(\mathbf{H}(\mathbf{w})\mathbf{H}(\mathbf{w})^\top)$. These notations will be used in Theorem 4.5. Since we know from the above Lemma that the matrices $\mathbf{H}(\underline{w}^0)\mathbf{H}(\underline{w}^0)^\top$ and $\mathbf{H}_k(\underline{w}^0)^\top \mathbf{H}_k(\underline{w}^0)$, $k \in [K]$ are full rank, we next ask if the above terms scale with m . Recall that we are looking at the Jacobian to state the condition under which Assumption 2.4 is satisfied. Thus, the following assumption is important, whose analytical justification is provided in App. G.

Assumption 4.4. *We assume that both $\lambda_{k,\rho}^-(m)$ and $\lambda_\rho^-(m)$ scale linearly with m .*

Experimental Justification of Assumption

4.4: An observation similar to the above assumption was also made in (Telgarsky, 2021, page 39). We verify the above assumption via experiments in Fig. 1, where we have plotted the minimum eigenvalue of the Jacobian versus m for different numbers of clients K using the MNIST data set (LeCun & Cortes, 2010). We can observe from the figure that the variation is almost linear, and the slope increases with decreasing K .

4.1 CONDITION ON NEURAL NETWORK (NN)

To prove the linear convergence of Algorithm 1 for single hidden layer NN, we need the definitions stated in equations 10 and 9. The following theorem provides a condition under which the Algorithm 1 converges linearly to a global optimal point, and the proof can be found in the Appendix I.

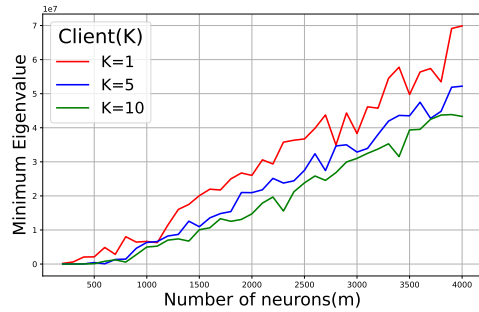


Figure 1: Plot of $\lambda_{\min}(m)$ versus m for $K = 1, 5, 10$. Here, $K = 1$ corresponds to $\lambda_{1,\min}(m)$. This shows that Assumption 4.4 is valid in the real-world setting as well, i.e., the minimum eigenvalue scales linearly with m .

³Here, \mathbf{e}_k and \mathbf{e} depend on \mathbf{w} .

Theorem 4.5. Let $\Psi_{m,K,n,\rho} := \sqrt{bn \left(\frac{\lambda_{\rho}^{+}(m)}{m} + \frac{d\Delta_{\sigma}^2\rho^2}{m} \right)}$ and $b := \frac{2D_{\sigma}^2\rho^2 d \log(2n/\delta)}{m} + 2y_{\max}^2$,

where $\lambda_{\rho}^{+}(m) := \sup_{\mathbf{w} \in \mathbb{B}[\mathbf{w}^0, \rho]} \frac{\|\mathbf{H}(\mathbf{w}^0)\mathbf{e}\|^2}{\|\mathbf{e}\|^2}$. The loss functions for single hidden layer NN satisfy equation 3 and equation 4 of Assumption 2.4 with a probability of at least $1 - \delta/2$, for any $\delta > 0$ provided the following holds:

$$\frac{\lambda_{k,\rho}^{-}(m)}{m} > 2 \times \left[\frac{\Delta_{\sigma}^2 d \rho^2}{m} + \frac{8bn}{\rho^2} \right], \text{ and } \frac{\lambda_{\rho}^{-}(m)}{m} > \frac{8K\Psi_{m,K,n,\rho}}{(1 - \zeta_{\rho})\rho} + \frac{2d\Delta_{\sigma}^2\rho}{m}, \quad (11)$$

where $\lambda_{k,\rho}^{-}(m)$ and $\lambda_{\rho}^{-}(m)$ are as defined in equation 9 and equation 10, respectively.

To the best of our knowledge, these conditions are the first of their kind. First, note that the terms $\lambda_{k,\rho}^{-}(m)/m$ and $\lambda_{\rho}^{-}(m)/m$ are less sensitive to ρ since they are sandwiched between the smallest and the largest eigenvalues of $\mathbf{H}(\mathbf{w}^0)^{\top} \mathbf{H}(\mathbf{w}^0)$ and $\mathbf{H}_k(\mathbf{w}^0)^{\top} \mathbf{H}_k(\mathbf{w}^0)$, respectively. In particular, these eigenvalues depend on the initialization \mathbf{w}^0 while the original condition is in terms of the ball around the initialization. Hence, using the eigenvalues in place of $\lambda_{k,\rho}^{-}(m)$ and $\lambda_{\rho}^{-}(m)$ in the new conditions makes it easy to verify (see Fig. 1). Secondly, the larger values of ρ make the right-hand sides in the equation 11 large, and hence the conditions may not be satisfied, as expected. On the other hand, the same can be observed for smaller values of ρ as well. Thus, a critical ρ is necessary. By choosing $\rho = c \times \mathcal{O}(\sqrt{n})$ and $m = \mathcal{O}(n^3)$ in Theorem 4.5 ensures that the right hand sides scale down with c . Thus, the right-hand side is small for a large enough c . However, by Assumption 4.4, the left-hand sides, i.e., $\lambda_{k,\rho}^{-}(m)/m$ and $\lambda_{\rho}^{-}(m)/m$ are constants that depend only on the initialization (not on ρ), and do not scale with m or n or c . Hence, the conditions are satisfied for large enough c .

Corollary 4.6. Choosing $\rho = c \times \mathcal{O}(\sqrt{n})$ and $m = \mathcal{O}(n^3)$ in Theorem 4.5 ensure that the conditions in equation 11 are satisfied for sufficiently large c .

The above corollary shows that by choosing a large radius of ρ and a large number of nodes in the second layer, linear convergence can be guaranteed. This brings in several challenges while proving the generalization guarantee, especially while proving a bound on the Rademacher complexity.

5 GENERALIZATION PERFORMANCE: SINGLE HIDDEN LAYER NN

In this section, we show that single hidden layer NN architectures exhibit impressive generalization guarantees. To state the generalization result, we need the following notion of Rademacher complexity of the single hidden layer NN.

Definition 5.1 (See Mohri et al. (2019)). The Rademacher complexity of a class of single hidden layer NN constrained to a ball of radius ρ at client $k \in [K]$ is defined as

$$\text{Rad}_k(\mathbf{w}^0, \rho) := \mathbb{E}_{\mathbf{v} \in \mathcal{G}_{\mathbf{v}}} \left[\sup_{\mathbf{w} \in \mathbb{B}[\mathbf{w}^0, \rho]} \frac{1}{n} \sum_{i=1}^n \zeta_i f_{\mathbf{w}; \mathbf{v}}(\mathbf{x}_{k,i}) \right],$$

where the expectation is with respect to $\zeta := (\zeta_1, \zeta_2, \dots, \zeta_n) \stackrel{i.i.d.}{\sim} \{-1, +1\}^n$, conditioned on $\mathbf{v} := (v_1, v_2, \dots, v_m) \in \mathcal{G}_{\mathbf{v}} := \{\mathbf{v} \in \{-1, 1\}^m : |\sum_{i=1}^n \zeta_i f_{\mathbf{w}; \mathbf{v}}(\mathbf{x})| < \Delta\}$. Here, $\Delta := \sqrt{2}D_{\sigma}d\sqrt{\frac{\rho^2+m}{m}} \log 4$ and \mathbf{x} is any data point sampled from $p_k(\mathbf{x})$.

For a FL setting, the generalization guarantee is provided in Mohri et al. (2019), and the result requires the loss to be bounded. However, in our case, the loss can potentially be unbounded. We handle this by focusing on the class of “good” NNs, i.e., $\mathbf{v} \in \mathcal{G}_{\mathbf{v}}$, whose output is bounded. In Appendix H, using the fact that the weight vector lies within a ball of radius ρ around \mathbf{w}^0 , we show that there exists such NNs with bounded output. Subsequently, we show that for such NNs, the generalization is guaranteed. We use this result along with the result of Mohri et al. (2019) to show the following Theorem whose proof can be found in Appendix J.

Theorem 5.2. Let $\Psi := \left((\rho^2 + 3m) \frac{2D_\sigma^2 d^2 \log 4}{m} + y_{max}^2 \right) \sqrt{2 \log(\frac{1}{\delta})}$. For the single hidden layer NN with the initialization as in Algorithm 2 satisfying Assumptions 4.4 with $m \geq nK/d$, and the conditions of Theorem 4.5, with a probability of at least $1 - \delta$, the following inequality holds

$$\Phi(\mathbf{w}; \mathbf{v}) \leq \Phi_S(\mathbf{w}; \mathbf{v}) + \frac{2n}{K} \sum_{k=1}^K \text{Rad}_k(\underline{\mathbf{w}}^0, \rho) + \Psi \sqrt{\frac{n}{K}}. \quad (12)$$

Recall that the loss function is defined as the sum of the loss on individual training samples. Thus, defining $\mathcal{L}(\mathbf{w}; \mathbf{v}) := \frac{\Phi(\mathbf{w}, \mathbf{v})}{n}$ and $\mathcal{L}_S(\mathbf{w}; \mathbf{v}) := \frac{\Phi_S(\mathbf{w}; \mathbf{v})}{n}$, and using this in the above theorem leads to the following.

Corollary 5.3. For the single hidden layer NN with initialization as in Algorithm 2, with probability at least $1 - 2\delta$ over the draw of the samples $X_k \sim \mathcal{D}_k^n$, the following inequality holds

$$\mathcal{L}(\mathbf{w}; \mathbf{v}) \leq \mathcal{L}_S(\mathbf{w}; \mathbf{v}) + \frac{2}{K} \sum_{k=1}^K \text{Rad}_k(\underline{\mathbf{w}}^0, \rho) + \frac{\Psi}{\sqrt{nK}}. \quad (13)$$

Next, we provide an upper bound on the Rademacher complexity.

Theorem 5.4. The Rademacher complexity of client $k \in [K]$ is bounded by

$$\text{Rad}_k(\underline{\mathbf{w}}^0, \rho) \leq \frac{1}{n\sqrt{m}} + \sqrt{\frac{\nu D_\sigma^2 d^2 (\log 4) \log(N_{\theta, \rho} / \delta_1)}{n}},$$

$$\text{where } \nu = (\rho^2 + 3m)/m, N_{\theta, \rho} := 3d^{3/4} \sqrt{\rho D_\sigma n m} \text{ and } \delta_1 := \frac{1}{2mn\sqrt{2D_\sigma d}} \sqrt{\frac{m}{\log 4(\rho^2 + m)}}.$$

Proof: See Appendix K. □

5.1 DISCUSSION

To the best of our knowledge, the above is the first result of its kind for an FL setup. We make the following remarks.

- The generalization error can be made small provided the right-hand side in the Corollary 5.3 is small. The first term, i.e., the empirical loss, depends on the communication rounds and the conditions stated in Theorem 4.5. The latter can be ensured by choosing $\rho = \mathcal{O}(\sqrt{n})$ and $m = \mathcal{O}(n^3)$, as shown in Corollary 4.6. In other words, the radius and the size of the NN scale with n which is not desired in general. However, we believe that this cannot be eliminated unless we make some structural assumptions about the data.
- Note that δ_1 and $N_{\theta, \rho}$ scale with n and m . However, it appears as a logarithmic term, and hence, the Rademacher complexity does not grow linearly with n . The above choices of ρ and m ensure that the Rademacher complexity in Theorem 5.4 goes down as $\mathcal{O}(1/\sqrt{n})$. Also, the choice of ρ cannot scale faster than \sqrt{m} .
- The last term in the generalization result scales down with n as $1/\sqrt{n}$. Based on these observations, it is clear that the generalization error can be made small by choosing large enough communication rounds R and the number of training samples n .
- Here, we present our theoretical insights on the effect of K . From the generalization bound in equation 5.3, it is evident that the last term decreases with K as $1/\sqrt{K}$. However, for larger values of K , the learning rate is impacted by K through $\frac{\zeta_p \rho}{T\sqrt{\Psi_0}}$, which scales as $1/\sqrt{K}$ (see Theorem 3.1). From equation 5, the loss goes down as $\exp\{-\mathcal{O}(R/\sqrt{K})\}$ leading to slower convergence. Thus, the overall effect of increasing K on the generalization is insignificant; this is also demonstrated in our experimental results as well as several existing works.

The above argument shows that the average loss can be made small by choosing sufficiently large m , n , and communication rounds, as shown next.⁴

Corollary 5.5. *With a probability of at least $1 - \delta$, there exists a single hidden layer NN employing the FedAvg algorithm with sufficiently large m , n , and R that achieves a small generalization error. More specifically, the generalization error goes down as $\mathcal{O}(1/\sqrt{n})$.*

6 EXPERIMENTAL RESULTS

In this section, we verify our theoretical findings with experiments performed on an NVIDIA DGX V100 machine. We have used an MNIST image data set LeCun & Cortes (2010) distributed across 5 and 200 clients. We have used the single hidden layer network model with 1000 neurons in the hidden layer and tanh activation function. In both cases, we have maintained around 50 data points at each client, which is less than the dimension of input feature vectors, i.e., around 1200, which satisfies the condition $d \geq n$ and $m \geq nK/d$. We execute FedAvg for $R = 500$ communication rounds along with $T = 5$ round of local updates at each client with i.i.d. data.

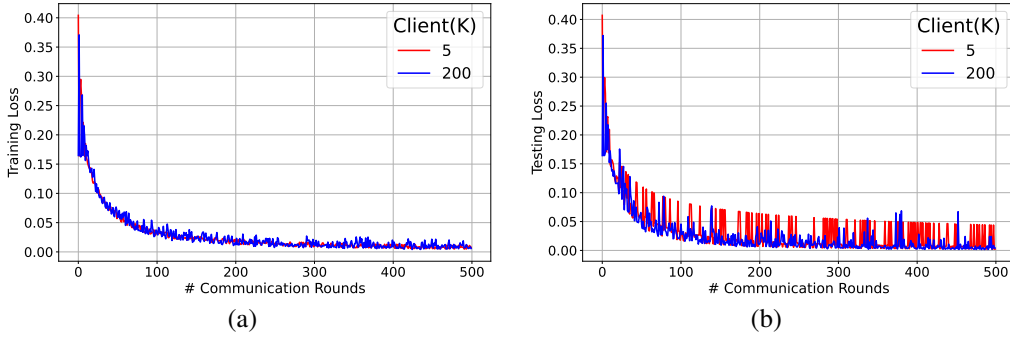


Figure 2: The Figures in (a) and (b) show the effect of the number of clients K on the training and the testing losses, respectively. The experiments are done using MNIST data set.

Figure 2 shows the effect of K on the testing and training errors. As suggested by our theory (see Sec. 5.1), increasing or decreasing K has no effect on the performance (generalization and training loss).

7 CONCLUSIONS

In this work, we addressed the problem of generalization along with convergence guarantees of the widely used FedAvg algorithm for solving Federated Learning (FL) problems. We proved the generalization bound by handling the optimization error and the Rademacher complexity. The optimization error was handled by proposing a novel and new constrained Polyak-Łojasiewicz (PL) type conditions on the (local) loss functions. Under these new conditions, we showed that there exists a global optimum to which the FedAvg converges linearly after $\mathcal{O}(\log(1/\epsilon))$ rounds of communication, where ϵ is the desired optimality gap. Importantly, we demonstrated that a class of single hidden layer NNs satisfy the proposed conditions that are required to establish the linear convergence of FedAvg as long as $m > \frac{nK}{d}$, where m is the number of neurons in the hidden layer, n is the number of samples at each client, K is the number of clients, and d is the feature dimension. Finally, we showed that the generalization error of FedAvg decreases at the rate of $\mathcal{O}(1/\sqrt{n})$ by proving a bound on the Rademacher Complexity using the fact that the neural network parameters are constrained to a neighbourhood around the initialization.

⁴While stating this result, we have ignored log factors.

REFERENCES

- Jing An and Jianfeng Lu. Convergence of stochastic gradient descent under a local łajaszewicz condition for deep neural networks. *arXiv preprint arXiv:2304.09221*, 2023.
- Sanjeev Arora, Simon Du, Wei Hu, Zhiyuan Li, and Ruosong Wang. Fine-grained analysis of optimization and generalization for overparameterized two-layer neural networks. In *International Conference on Machine Learning*, pp. 322–332. PMLR, 2019.
- Idan Attias, Aryeh Kontorovich, and Y. Mansour. Improved generalization bounds for adversarially robust learning. *J. Mach. Learn. Res.*, 23:175:1–175:31, 2018. URL <https://api.semanticscholar.org/CorpusID:250244124>.
- Eugene Bagdasaryan, Andreas Veit, Yiqing Hua, Deborah Estrin, and Vitaly Shmatikov. How to backdoor federated learning. In *International conference on artificial intelligence and statistics*, pp. 2938–2948. PMLR, 2020.
- Keith Bonawitz, Vladimir Ivanov, Ben Kreuter, Antonio Marcedone, H. B. McMahan, Sarvar Patel, Daniel Ramage, Aaron Segal, and Karn Seth. Practical secure aggregation for federated learning on user-held data. *ArXiv*, abs/1611.04482, 2016. URL <https://api.semanticscholar.org/CorpusID:10933707>.
- Olivier Bousquet and André Elisseeff. Stability and generalization. *The Journal of Machine Learning Research*, 2:499–526, 2002.
- Sourav Chatterjee. Convergence of gradient descent for deep neural networks. *arXiv preprint arXiv:2203.16462*, 2022.
- Shuxiao Chen, Qinqing Zheng, Qi Long, and Weijie J Su. A theorem of the alternative for personalized federated learning. *arXiv preprint arXiv:2103.01901*, 2021.
- Melikasadat Emami, Mojtaba Sahraee-Ardakan, Parthe Pandit, Sundeep Rangan, and Alyson Fletcher. Generalization error of generalized linear models in high dimensions. In *International Conference on Machine Learning*, pp. 2892–2901. PMLR, 2020.
- Alireza Fallah, Aryan Mokhtari, and Asuman Ozdaglar. Generalization of model-agnostic meta-learning algorithms: Recurring and unseen tasks. *Advances in Neural Information Processing Systems*, 34:5469–5480, 2021.
- Chen Fan, Christos Thrampoulidis, and Mark Schmidt. Fast convergence of random reshuffling under over-parameterization and the polyak-łojaszewicz condition. In *Joint European Conference on Machine Learning and Knowledge Discovery in Databases*, pp. 301–315. Springer, 2023.
- Robin C. Geyer, Tassilo Klein, and Moin Nabi. Differentially private federated learning: A client level perspective. *ArXiv*, abs/1712.07557, 2017. URL <https://api.semanticscholar.org/CorpusID:3630366>.
- Farzin Haddadpour and Mehrdad Mahdavi. On the convergence of local descent methods in federated learning. *arXiv preprint arXiv:1910.14425*, 2019.
- Farzin Haddadpour, Mohammad Mahdi Kamani, Mehrdad Mahdavi, and Viveck Cadambe. Local sgd with periodic averaging: Tighter analysis and adaptive synchronization. *Advances in Neural Information Processing Systems*, 32, 2019.
- Xiaolin Hu, Shaojie Li, and Yong Liu. Generalization bounds for federated learning: Fast rates, un-participating clients and unbounded losses. In *The Eleventh International Conference on Learning Representations*, 2022.
- Ziwei Ji and Matus Telgarsky. Gradient descent aligns the layers of deep linear networks. *arXiv preprint arXiv:1810.02032*, 2018.
- Peter Kairouz, H Brendan McMahan, Brendan Avent, Aurélien Bellet, Mehdi Bennis, Arjun Nitin Bhagoji, Kallista Bonawitz, Zachary Charles, Graham Cormode, Rachel Cummings, et al. Advances and open problems in federated learning. *Foundations and Trends® in Machine Learning*, 14(1–2):1–210, 2021.

- Hamed Karimi, Julie Nutini, and Mark Schmidt. Linear convergence of gradient and proximal-gradient methods under the polyak-łojasiewicz condition. In *Machine Learning and Knowledge Discovery in Databases: European Conference, ECML PKDD 2016, Riva del Garda, Italy, September 19-23, 2016, Proceedings, Part I 16*, pp. 795–811. Springer, 2016.
- Sai Praneeth Karimireddy, Satyen Kale, Mehryar Mohri, Sashank Reddi, Sebastian Stich, and Ananda Theertha Suresh. Scaffold: Stochastic controlled averaging for federated learning. In *International Conference on Machine Learning*, pp. 5132–5143. PMLR, 2020.
- Ahmed Khaled, Konstantin Mishchenko, and Peter Richtárik. First analysis of local gd on heterogeneous data. *arXiv preprint arXiv:1909.04715*, 2019.
- Jakub Konečný, Brendan McMahan, and Daniel Ramage. Federated optimization: Distributed optimization beyond the datacenter. *arXiv preprint arXiv:1511.03575*, 2015.
- Jakub Konečný, H Brendan McMahan, Felix X Yu, Peter Richtárik, Ananda Theertha Suresh, and Dave Bacon. Federated learning: Strategies for improving communication efficiency. *arXiv preprint arXiv:1610.05492*, 2016.
- Yann LeCun and Corinna Cortes. MNIST handwritten digit database. 2010. URL <http://yann.lecun.com/exdb/mnist/>.
- Tian Li, Anit Kumar Sahu, Ameet Talwalkar, and Virginia Smith. Federated learning: Challenges, methods, and future directions. *IEEE signal processing magazine*, 37(3):50–60, 2020.
- Xiang Li, Kaixuan Huang, Wenhao Yang, Shusen Wang, and Zhihua Zhang. On the convergence of fedavg on non-iid data. *arXiv preprint arXiv:1907.02189*, 2019.
- Xiaoxiao Li, Zhao Song, Runzhou Tao, and Guangyi Zhang. A convergence theory for federated average: Beyond smoothness. In *2022 IEEE International Conference on Big Data (Big Data)*, pp. 1292–1297. IEEE, 2022.
- Chaoyue Liu, Libin Zhu, and Mikhail Belkin. Loss landscapes and optimization in over-parameterized non-linear systems and neural networks. *Applied and Computational Harmonic Analysis*, 59:85–116, 2022.
- Brendan McMahan, Eider Moore, Daniel Ramage, Seth Hampson, and Blaise Agüera y Arcas. Communication-efficient learning of deep networks from decentralized data. In *Artificial intelligence and statistics*, pp. 1273–1282. PMLR, 2017.
- Mehryar Mohri, Afshin Rostamizadeh, and Ameet Talwalkar. *Foundations of machine learning*. 2018.
- Mehryar Mohri, Gary Sivek, and Ananda Theertha Suresh. Agnostic federated learning. In *International Conference on Machine Learning*, pp. 4615–4625. PMLR, 2019.
- Samet Oymak and Mahdi Soltanolkotabi. Overparameterized nonlinear learning: Gradient descent takes the shortest path? In *International Conference on Machine Learning*, pp. 4951–4960. PMLR, 2019.
- Zhaonan Qu, Kaixiang Lin, Jayant Kalagnanam, Zhaojian Li, Jiayu Zhou, and Zhengyuan Zhou. Federated learning’s blessing: Fedavg has linear speedup. *arXiv preprint arXiv:2007.05690*, 2020.
- Zhaonan Qu, Kaixiang Lin, Zhaojian Li, Jiayu Zhou, and Zhengyuan Zhou. Federated learning’s blessing: Fedavg has linear speedup, 2021. URL <https://openreview.net/forum?id=yJHpncwG1B>.
- Shai Shalev-Shwartz and Shai Ben-David. *Understanding Machine Learning - From Theory to Algorithms*. Cambridge University Press, 2014. ISBN 978-1-10-705713-5.
- Shai Shalev-Shwartz, Ohad Shamir, Nathan Srebro, and Karthik Sridharan. Learnability, stability and uniform convergence. *Journal of Machine Learning Research*, 11(90):2635–2670, 2010. URL <http://jmlr.org/papers/v11/shalev-shwartz10a.html>.

- Virginia Smith, Chao-Kai Chiang, Maziar Sanjabi, and Ameet S Talwalkar. Federated multi-task learning. *Advances in neural information processing systems*, 30, 2017.
- Bingqing Song, Prashant Khanduri, Xinwei Zhang, Jinfeng Yi, and Mingyi Hong. Fedavg converges to zero training loss linearly for overparameterized multi-layer neural networks. In *International Conference on Machine Learning*, pp. 32304–32330. PMLR, 2023.
- Sebastian U Stich. Local sgd converges fast and communicates little. *arXiv preprint arXiv:1805.09767*, 2018.
- Zhenyu Sun, Xiaochun Niu, and Ermin Wei. Understanding generalization of federated learning via stability: Heterogeneity matters. *ArXiv*, abs/2306.03824, 2023. URL <https://api.semanticscholar.org/CorpusID:259088815>.
- Matus Telgarsky. Deep learning theory lecture notes. <https://mjt.cs.illinois.edu/dlt/>, 2021. Version: 2021-10-27 v0.0-e7150f2d (alpha).
- Martin J Wainwright. *High-dimensional statistics: A non-asymptotic viewpoint*, volume 48. Cambridge university press, 2019.
- Jianyu Wang and Gauri Joshi. Cooperative sgd: A unified framework for the design and analysis of local-update sgd algorithms. *The Journal of Machine Learning Research*, 22(1):9709–9758, 2021.
- Blake Woodworth, Kumar Kshitij Patel, Sebastian Stich, Zhen Dai, Brian Bullins, Brendan Mcmahon, Ohad Shamir, and Nathan Srebro. Is local sgd better than minibatch sgd? In *International Conference on Machine Learning*, pp. 10334–10343. PMLR, 2020a.
- Blake E Woodworth, Kumar Kshitij Patel, and Nati Srebro. Minibatch vs local sgd for heterogeneous distributed learning. *Advances in Neural Information Processing Systems*, 33:6281–6292, 2020b.
- Bing Xu, Naiyan Wang, Tianqi Chen, and Mu Li. Empirical evaluation of rectified activations in convolutional network. *arXiv preprint arXiv:1505.00853*, 2015.
- Dong Yin, Kannan Ramchandran, and Peter L. Bartlett. Rademacher complexity for adversarially robust generalization. In *International Conference on Machine Learning*, 2018. URL <https://api.semanticscholar.org/CorpusID:53092122>.
- Hao Yu, Rong Jin, and Sen Yang. On the linear speedup analysis of communication efficient momentum sgd for distributed non-convex optimization. In *International Conference on Machine Learning*, pp. 7184–7193. PMLR, 2019.
- Honglin Yuan, Warren Morningstar, Lin Ning, and Karan Singhal. What do we mean by generalization in federated learning? *arXiv preprint arXiv:2110.14216*, 2021a.
- Honglin Yuan, Warren R. Morningstar, Lin Ning, and K. Singhal. What do we mean by generalization in federated learning? *ArXiv*, abs/2110.14216, 2021b. URL <https://api.semanticscholar.org/CorpusID:239998253>.
- Chiyuan Zhang, Samy Bengio, Moritz Hardt, Benjamin Recht, and Oriol Vinyals. Understanding deep learning requires rethinking generalization. *ArXiv*, abs/1611.03530, 2016. URL <https://api.semanticscholar.org/CorpusID:6212000>.
- Jiawei Zhang, Yushun Zhang, Mingyi Hong, Ruoyu Sun, and Zhi-Quan Luo. When expressivity meets trainability: Fewer than n neurons can work. *Advances in Neural Information Processing Systems*, 34:9167–9180, 2021.
- Fan Zhou and Guojing Cong. On the convergence properties of a k -step averaging stochastic gradient descent algorithm for nonconvex optimization. *ArXiv*, abs/1708.01012, 2017a. URL <https://api.semanticscholar.org/CorpusID:3384938>.
- Fan Zhou and Guojing Cong. On the convergence properties of a k -step averaging stochastic gradient descent algorithm for nonconvex optimization. *arXiv preprint arXiv:1708.01012*, 2017b.

A APPENDIX

B RELATED WORK

Convergence of FedAvg Algorithm: McMahan et al. (2017) first introduced federated learning (FL) to learn a global model from distributed data without sharing it with a central server. Majority of the research on FL focuses on communication-efficiency Konečný et al. (2016); McMahan et al. (2017); Li et al. (2020); Smith et al. (2017) and data-privacy Bagdasaryan et al. (2020); Bonawitz et al. (2016); Geyer et al. (2017). FedAvg algorithm has been studied extensively under some key assumptions on loss functions such as convex, strongly convex, and non-convex loss functions with an overparametrized NN setting Stich (2018); Wang & Joshi (2021); Khaled et al. (2019); Yu et al. (2019). Stich (2018) has substantiated that FedAvg (LocalSGD) exhibits a provable achievement of linear speedup with significantly reduced communication requirements, specifically in the realm of strongly convex stochastic optimization. Several notable studies contribute to a comprehensive understanding of FedAvg in different optimization settings. Zhou & Cong (2017a); Wang & Joshi (2021) have delved into the non-convex setting, establishing crucial convergence results. A flurry of research articles has been published so far on the convergence analysis of the FedAvg algorithm Qu et al. (2021); Li et al. (2022); Song et al. (2023). The key assumption, in all the above analyses, is the existence of global minimum and Polyak-Lojasiewicz(PL) condition Karimi et al. (2016), which is assumed to be satisfied by loss function over *whole parameter space* that restrict the class of loss function. Liu et al. (2022) relaxed the above assumption to a small neighbourhood around initialization, and called it modified PL inequality, i.e, loss function satisfies the modified PL inequality within a small neighbourhood around initialization under over-parameterized regime but for centralized data setting. Further, Chatterjee (2022) has derived a sufficient condition that depends on the initialization and radius of the ball for the convergence of the gradient descent (GD) to the global optimum, and the same is extended for SGD by An & Lu (2023) for single client setting. In this paper, we have come up with a sufficient condition on the loss function that guarantees the linear convergence of the FedAvg algorithm to a global minimum where the existence of a global minimum is not an assumption, and it is a part of our conclusion.

Generalization of FedAvg Algorithm: The concept of generalization in centralized learning has been a focal point for researchers for several decades. To assess generalization error, the commonly employed approaches involve utilizing uniform convergence, often tied to VC dimension or Rademacher complexity Shalev-Shwartz et al. (2010); Yin et al. (2018); Attias et al. (2018). However, there is a limitation in cases where uniform convergence yields an excessively loose bound, diminishing its meaningfulness Zhang et al. (2016). This limitation arises because uniform convergence examines the hypothesis class while disregarding the impact of training algorithms responsible for generating these hypotheses. Mohri et al. (2019) establishes a uniform convergence bound of $\mathcal{O}(1/\sqrt{n})$, where n represents the cumulative number of samples collected by all participating clients, for agnostic federated learning problems using Rademacher complexity under binary losses. However, we have not assumed any bound on the assumed loss function. Hence, our result is more general. Recent work by Yuan et al. (2021b); Sun et al. (2023) has explored the meaning of generalization in federated learning, and Sun et al. (2023) has also proved an upper bound on the true and empirical risk of the model obtained by FedAvg, FedProx, and SCAFFOLD. To prove a better generalization of the model, we must achieve low optimization error, preferably linear, which is missing in Sun et al. (2023). In this paper, we have shown that the optimization error of the model obtained by FedAvg converges linearly to 0. Also, we have proved that the Rademacher Complexity tends to 0 for a sufficiently large sample.

C NOTATION

We denote the global round of communication as \underline{w}^r , and $\underline{w}_k^{r,T}$ represents T local rounds of SGD performed by client k after r global rounds of communication. Additionally, we assume that the input space is \mathbb{R}^d , and $\Phi_k : \mathbb{R}^{m(d+1)} \rightarrow \mathbb{R}$ represents the loss function associated with each client $k \in [K]$, where $[K]$ denotes the set $\{1, 2, \dots, N\}$. Let us assume that each client has $n \in \mathbb{N}$ training examples, denoted by $\{\underline{x}_{k,i}, y_{k,i}\}_{i=1}^n$, where k represents the k -th client. Moreover, $\underline{X}_k \in \mathbb{R}^{n \times d}$ represents the input data matrix at client k . We use $\mathbb{B}[\underline{w}, \rho]$ to denote closed Euclidean balls, respectively, centred at \underline{w} with a radius of ρ . We use $\text{vec}[\underline{w}_1, \underline{w}_2, \dots, \underline{w}_n] \in \mathbb{R}^{mn}$ to denote a

vector obtained by stacking each of the vectors $\mathbf{w}_i \in \mathbb{R}^m$, $i \in [n]$. We use \mathbb{E} to denote the expected value and $\mathbb{E}(X|Y)$ or $\mathbb{E}_{|Y}(X)$ to denote the conditional expectation.

D USEFUL LEMMA: CONCENTRATION BOUND

In this subsection, we state the following useful lemmas for use in later proofs.

Lemma D.1. (See Wainwright (2019)) Let $X_i \sim \mathcal{N}(0, 1)$, $i = 1, 2, \dots, M$, then $\Pr \left\{ \sum_{i=1}^M X_i^2 > M(1+t) \right\} \leq e^{-Mt^2/18}$, where $t \in [0, 3]$. Specifically, if $\underline{\mathbf{w}}_0 \sim \mathcal{N}(0, \frac{1}{d} \times I)$, we have

$$\mathbb{P} [\|\underline{\mathbf{w}}^0\|^2 \geq \zeta m] \leq \exp \left\{ -\frac{dm(\zeta - 1)^2}{18} \right\},$$

where $\zeta \in [1, 4]$.

Lemma D.2. Let the local loss function Φ_k for all $k \in [K]$ satisfy Assumption 2.2. Then

$$\mathbb{E} \left[\left\| \widehat{\nabla \Phi_k}(\mathbf{w}_k^{r,\tau}) \right\|^2 \right] := \mathbb{E} \left[\left\| \frac{1}{b} \sum_{i \in \mathcal{B}^{r,\tau}} \nabla \Phi_{k,i}(\mathbf{w}_k^{r,\tau}) \right\|^2 \right] \leq 2l'_k \Phi_k(\mathbf{w}_k^{r,\tau}),$$

where $\mathcal{B}^{r,\tau}$ is random sample-batch of size $|\mathcal{B}^{r,\tau}| = b$ and $l'_k := \max_i l_{k,i}$.

Proof. We can re-write the above expression as

$$\begin{aligned} \mathbb{E} \left[\left\| \frac{1}{b} \sum_{i \in \mathcal{B}^{r,\tau}} \nabla \Phi_{k,i}(\mathbf{w}_k^{r,\tau}) \right\|^2 \right] &= \frac{1}{b^2} \mathbb{E} \left[\sum_{i \in \mathcal{B}^{r,\tau}} \|\nabla \Phi_{k,i}(\mathbf{w}_k^{r,\tau})\|^2 \right] \\ &\quad + \frac{1}{b^2} \sum_{i \neq i'} \mathbb{E} [\langle \nabla \Phi_{k,i}(\mathbf{w}_k^{r,\tau}), \nabla \Phi_{k,i'}(\mathbf{w}_k^{r,\tau}) \rangle] \\ &\leq \frac{1}{b^2} \sum_{i \in \mathcal{B}^{r,\tau}} \mathbb{E} [\|\nabla \Phi_{k,i}(\mathbf{w}_k^{r,\tau})\|^2] + \frac{1}{b^2} \sum_{i \neq i'} \mathbb{E} [\|\nabla \Phi_{k,i}(\mathbf{w}_k^{r,\tau})\| \|\nabla \Phi_{k,i'}(\mathbf{w}_k^{r,\tau})\|] \\ &\leq \frac{2}{b^2} \sum_{i \in \mathcal{B}^{r,\tau}} l_{k,i} \mathbb{E} [\Phi_{k,i}(\mathbf{w}_k^{r,\tau})] + \frac{2}{b^2} \sum_{i \neq i'} l_{k,i} \mathbb{E} [\Phi_{k,i}(\mathbf{w}_k^{r,\tau})], \end{aligned}$$

where the above inequality follows from the Assumption 2.2 and Jensen's inequality. Using the above, we have

$$\begin{aligned} \mathbb{E} \left[\left\| \frac{1}{b} \sum_{i \in \mathcal{B}^{r,\tau}} \nabla \Phi_{k,i}(\mathbf{w}_k^{r,\tau}) \right\|^2 \right] &\leq \left(\frac{2}{b} + \frac{2(b-1)}{b} \right) l'_k \Phi_k(\mathbf{w}_k^{r,\tau}) \\ &= 2l'_k \Phi_k(\mathbf{w}_k^{r,\tau}), \end{aligned}$$

where we used the unbiased estimate of the local loss function, i.e., $\mathbb{E} [\Phi_{k,i}(\mathbf{w}_k^{r,\tau})] = \Phi_k(\mathbf{w}_k^{r,\tau})$ and $l'_k := \max_i l_{k,i}$. \square

E NEW CONDITION FOR LINEAR CONVERGENCE: SINGLE CLIENT SETTING

In this subsection, we show that the local SGD iterates for any client $k \in [K]$ stays within the ball $\mathbb{B}[\underline{\mathbf{w}}^0, \rho]$. This result will be used to prove new conditions for the linear convergence of FedAvg.

Theorem E.1. Let $\Phi_k : \mathbb{R}^{d'} \rightarrow [0, \infty)$ be a non-negative L_k -smooth function. Take any $\underline{\mathbf{w}}^0 \in \mathbb{R}^{d'}$ and $\rho > 0$. Let $\mathbb{B}[\underline{\mathbf{w}}^0, \rho]$ denote the closed Euclidean ball of radius ρ , centered at $\underline{\mathbf{w}}^0$. Assume that for $\epsilon \in (0, 1)$ and $\eta \leq \min \left\{ \frac{\alpha_k}{2L_k l'_k}, \sqrt{\frac{3}{L_k l'_k}}, \frac{\rho}{G}, \frac{2}{\alpha_k} \right\}$, the following holds:

$$32\Phi_k(\underline{\mathbf{w}}^0) \leq \rho^2 \alpha_k, \quad (14)$$

where α_k is defined in 2 but with $f(\cdot)$ replaced by $\Phi_k(\cdot)$. Consider the SGD update rule defined in Algorithm 1. Then $\mathbf{w}_k^{0,T} \in \mathbb{B}[\underline{\mathbf{w}}^0, \rho]$ for all $T \geq 0$.⁵

Proof. If $\Phi_k(\underline{\mathbf{w}}^0) = 0$, then $\nabla \Phi_k(\underline{\mathbf{w}}^0) = 0$, for all $k \geq 0$. So, let $\Phi_k(\underline{\mathbf{w}}^0) > 0$. The proof is by the method of induction on $T \in \mathbb{N}$. By definition of euclidean ball, we have $\underline{\mathbf{w}}^0 = \mathbf{w}_k^{0,0} \in \mathbb{B}[\underline{\mathbf{w}}^0, \rho]$. Let the induction hypothesis be $\mathbf{w}_k^{0,1}, \mathbf{w}_k^{0,2}, \dots, \mathbf{w}_k^{0,T-1} \in \mathbb{B}[\underline{\mathbf{w}}^0, \rho]$ for some $T \geq 1$. We need to show that $\mathbf{w}_k^{0,T} \in \mathbb{B}[\underline{\mathbf{w}}^0, \rho]$ using induction hypothesis. Toward this, we need the following lemmas:

Lemma E.2. For each $\tau > 0$, define

$$\Lambda_\tau := \Phi_k(\mathbf{w}_k^{0,\tau+1}) - \Phi_k(\mathbf{w}_k^{0,\tau}) + \eta \|\nabla \Phi_k(\mathbf{w}_k^{0,\tau})\|^2$$

Then for all $0 \leq \tau \leq T-1$, we have

$$|\mathbb{E}[\Lambda_\tau]| \leq \eta^2 L_{\max} l'_{\max} \Phi_k(\mathbf{w}_k^{0,\tau}),$$

where $L_{\max} := \max_k L_k$ and $l'_{\max} := \max_k l'_k$.

Proof. By L_k -smoothness Assumption 2.1 of local loss function Φ_k and local iterates of SGD, for $1 \leq \tau \leq T-1$, we have

$$\begin{aligned} \Lambda_\tau &= \Phi_k(\mathbf{w}_k^{0,\tau} - \eta \widehat{\nabla \Phi_k}(\mathbf{w}_k^{0,\tau})) - \Phi_k(\mathbf{w}_k^{0,\tau}) + \eta \|\nabla \Phi_k(\mathbf{w}_k^{0,\tau})\|^2 \\ &= -\eta \langle \widehat{\nabla \Phi_k}(\mathbf{w}_k^{0,\tau}), \nabla \Phi_k(\mathbf{w}_k^{0,\tau}) \rangle + \frac{\eta^2 L_k}{2} \|\widehat{\nabla \Phi_k}(\mathbf{w}_k^{0,\tau})\|^2 + \eta \|\nabla \Phi_k(\mathbf{w}_k^{0,\tau})\|^2, \end{aligned}$$

where $\widehat{\nabla \Phi_k}$ is an unbiased estimate of true gradient over mini-batch $\mathcal{B}^{0,\tau}$ of samples. Next, we take the expectation over randomness in mini-batch of samples. We get

$$\begin{aligned} \mathbb{E}[\Lambda_\tau] &= -\eta \langle \mathbb{E}[\widehat{\nabla \Phi_k}(\mathbf{w}_k^{0,\tau})], \nabla \Phi_k(\mathbf{w}_k^{0,\tau}) \rangle + \frac{\eta^2 L_k}{2} \mathbb{E}[\|\widehat{\nabla \Phi_k}(\mathbf{w}_k^{0,\tau})\|^2] + \eta \|\nabla \Phi_k(\mathbf{w}_k^{0,\tau})\|^2, \\ &\leq \frac{\eta^2 L_k}{2} \mathbb{E}[\|\widehat{\nabla \Phi_k}(\mathbf{w}_k^{0,\tau})\|^2]. \end{aligned}$$

Further, we reduce the above inequality as follows

$$\begin{aligned} |\mathbb{E}[\Lambda_\tau]| &\leq \frac{\eta^2 L_k}{2} \mathbb{E}[2l'_k \widehat{\Phi_k}(\mathbf{w}_k^{0,\tau})] \\ &= \eta^2 L_{\max} l'_{\max} \Phi_k(\mathbf{w}_k^{0,\tau}), \end{aligned}$$

where the first inequality follows from the Lemma D.2 and the last equality is due to true estimate assumption, i.e., $\mathbb{E}[\widehat{\Phi_k}(\mathbf{w}_k^{0,\tau})] = \Phi_k(\mathbf{w}_k^{0,\tau})$ and, $L_k \leq L_{\max}$ and $l'_k \leq l'_{\max}$ for all $k \in [K]$. \square

Lemma E.3. By choosing $\eta \leq \min \left\{ \frac{2}{\alpha_{\min}}, \frac{\alpha_{\min}}{2L_{\max}l'_{\max}} \right\}$ and for any $0 \leq \tau \leq T-1$, we have

$$\Phi_k(\mathbf{w}_k^{0,\tau}) \leq \left(1 - \frac{\alpha_{\min}\eta}{2}\right)^\tau \Phi_k(\underline{\mathbf{w}}^0). \quad (15)$$

Further, we have $\frac{\alpha_{\min}\eta}{2} \leq 1$ where $\alpha_{\min} := \min_k \{\alpha_1, \alpha_2, \dots, \alpha_k\}$

Proof. Since $\mathbf{w}_k^{0,\tau} \in \mathbb{B}[\underline{\mathbf{w}}^0, \rho]$ for $\tau \leq T-1$. By rearranging the terms in the definition of Λ_τ from Lemma E.2 under expectation,

$$\begin{aligned} \Phi_k(\mathbf{w}_k^{0,\tau+1}) &= \Phi_k(\mathbf{w}_k^{0,\tau}) - \eta \|\nabla \Phi_k(\mathbf{w}_k^{0,\tau})\|^2 + \mathbb{E}[\Lambda_\tau] \\ &\stackrel{(a)}{\leq} \Phi_k(\mathbf{w}_k^{0,\tau}) - \eta \alpha_k \Phi_k(\mathbf{w}_k^{0,\tau}) + \eta^2 L_{\max} l'_{\max} \Phi_k(\mathbf{w}_k^{0,\tau}) \\ &= (1 - \eta \alpha_k + \eta^2 L_{\max} l'_{\max}) \Phi_k(\mathbf{w}_k^{0,\tau}) \\ &= (1 - \eta(\alpha_{\min} - \eta L_{\max} l'_{\max}))^{\tau+1} \Phi_k(\underline{\mathbf{w}}^0), \end{aligned}$$

⁵Note that we are not assuming that the gradient is bounded throughout but due to smoothness, it is indeed bounded over a closed ball. This only impacts the learning rate η , which we are free to choose.

where (a) follows from the definition of α_k equation 2 and Lemma E.2. The last inequality follows from iteration over τ and $\alpha_{\min} \leq \alpha_k$ for all $k \in [K]$. By choosing $\eta < \frac{\alpha_{\min}}{2L_{\max}l'_{\max}}$, we have

$$\Phi_k(\mathbf{w}_k^{0,\tau+1}) \leq \left(1 - \frac{\alpha_{\min}\eta}{2}\right)^{\tau+1} \Phi_k(\mathbf{w}^0). \quad (16)$$

Now, since $\Phi_k(\mathbf{w}^0) > 0$ and $\Phi_k(\mathbf{w}_k^{0,1}) \geq 0$. Choose $\tau = 0$ and dividing equation 16 by $\Phi_k(\mathbf{w}^0)$, we get

$$\begin{aligned} 0 \leq \frac{\Phi_k(\mathbf{w}_k^{0,1})}{\Phi_k(\mathbf{w}^0)} &\leq 1 - \frac{\alpha_{\min}\eta}{2} \\ \frac{\alpha_{\min}\eta}{2} &\leq 1 \end{aligned}$$

which completes the proof. \square

Lemma E.4. For each $\tau \leq T - 1$,

$$\eta \|\nabla \Phi_k(\mathbf{w}_k^{0,\tau})\|^2 \leq (1 + \eta^2 L_{\max} l'_{\max}) \Phi_k(\mathbf{w}_k^{0,\tau}) - \Phi_k(\mathbf{w}_k^{0,\tau+1}).$$

Further, we have

$$\Phi_k(\mathbf{w}_k^{0,\tau+1}) \leq (1 + \eta^2 L_{\max} l'_{\max}) \Phi_k(\mathbf{w}_k^{0,\tau}).$$

Proof. By rearranging the terms in the definition of Λ_τ under expectation and using Lemma E.2, for $\tau \leq T - 1$, we have

$$\begin{aligned} \Phi_k(\mathbf{w}_k^{0,\tau+1}) &= \Phi_k(\mathbf{w}_k^{0,\tau}) - \eta \|\nabla \Phi_k(\mathbf{w}_k^{0,\tau})\|^2 + \mathbb{E}[\Lambda_\tau] \\ &\leq \Phi_k(\mathbf{w}_k^{0,\tau}) - \eta \|\nabla \Phi_k(\mathbf{w}_k^{0,\tau})\|^2 + \eta^2 L_{\max} l'_{\max} \Phi_k(\mathbf{w}_k^{0,\tau}). \end{aligned}$$

Next, we again rearrange the above inequality to upper-bound the norm of gradient square, i.e., $\eta \|\nabla \Phi_k(\mathbf{w}_k^{0,\tau})\|^2$,

$$\begin{aligned} \eta \|\nabla \Phi_k(\mathbf{w}_k^{0,\tau})\|^2 &\leq \Phi_k(\mathbf{w}_k^{0,\tau}) - \Phi_k(\mathbf{w}_k^{0,\tau+1}) + \eta^2 L_{\max} l'_{\max} \Phi_k(\mathbf{w}_k^{0,\tau}) \\ &\leq (1 + \eta^2 L_{\max} l'_{\max}) \Phi_k(\mathbf{w}_k^{0,\tau}) - \Phi_k(\mathbf{w}_k^{0,\tau+1}). \end{aligned}$$

This completes the proof of the first part.

To prove the second part, we observe that $\eta \|\nabla \Phi_k(\mathbf{w}_k^{0,\tau})\|^2 \geq 0$. Using this in the first part of this lemma proves the second part, i.e.,

$$\Phi_k(\mathbf{w}_k^{0,\tau+1}) \leq (1 + \eta^2 L_{\max} l'_{\max}) \Phi_k(\mathbf{w}_k^{0,\tau}).$$

\square

Lemma E.5. For all $\tau \leq T - 1$, we have

$$\sum_{l=\tau}^{T-1} \eta \left\| \nabla \Phi_k(\mathbf{w}_k^{0,l}) \right\| \leq (1 - (\alpha_{\min} - \epsilon)\eta)^{\tau/2} \frac{\sqrt{32\Phi_k(\mathbf{w}^0)\eta}}{(\alpha_{\min} - \epsilon)\eta}, \quad (17)$$

provided $\eta < \min \left\{ \sqrt{\frac{3}{L_{\max} l'_{\max}}}, \frac{\alpha_{\min}}{2L_{\max} l'_{\max}} \right\}$.

Proof. From Lemma E.4, we have

$$\begin{aligned} \eta \left\| \nabla \Phi_k(\mathbf{w}_k^{0,\tau}) \right\| &= \sqrt{\eta^2 \|\nabla \Phi_k(\mathbf{w}_k^{0,\tau})\|^2} \\ &\leq \left[\eta (1 + \eta^2 L_{\max} l'_{\max}) \Phi_k(\mathbf{w}_k^{0,\tau}) - \eta \Phi_k(\mathbf{w}_k^{0,\tau+1}) \right]^{1/2}. \end{aligned}$$

Let us denote $1 + \eta^2 L_{\max} l'_{\max}$ by γ , i.e., $\gamma = 1 + \eta^2 L_{\max} l'_{\max}$, in the above, we get

$$\eta \left\| \nabla \Phi_k(\mathbf{w}_k^{0,\tau}) \right\| \leq \sqrt{\eta} \left[\gamma \Phi_k(\mathbf{w}_k^{0,\tau}) - \Phi_k(\mathbf{w}_k^{0,\tau+1}) \right]^{1/2}.$$

We can rewrite the above inequality as

$$\eta \|\nabla \Phi_k(\mathbf{w}_k^{0,\tau})\| \leq \sqrt{\eta} \left[\left(\sqrt{\gamma \Phi_k(\mathbf{w}_k^{0,\tau})} + \sqrt{\Phi_k(\mathbf{w}_k^{0,\tau+1})} \right) \left(\sqrt{\gamma \Phi_k(\mathbf{w}_k^{0,\tau})} - \sqrt{\Phi_k(\mathbf{w}_k^{0,\tau+1})} \right) \right]^{1/2}.$$

Taking summation on both sides of the above inequality from $l = \tau$ to $T - 1$. Then, the above inequality is reduced to

$$\begin{aligned} \sum_{l=\tau}^{T-1} \eta \|\nabla \Phi_k(\mathbf{w}_k^{0,l})\| &\leq \sqrt{\eta} \sum_{l=\tau}^{T-1} \left[\left(\sqrt{\gamma \Phi_k(\mathbf{w}_k^{0,l})} + \sqrt{\Phi_k(\mathbf{w}_k^{0,l+1})} \right) \left(\sqrt{\gamma \Phi_k(\mathbf{w}_k^{0,l})} - \sqrt{\Phi_k(\mathbf{w}_k^{0,l+1})} \right) \right]^{1/2} \\ &\stackrel{(a)}{\leq} \sqrt{\eta} \left[\sum_{l=\tau}^{T-1} \left(\sqrt{\gamma \Phi_k(\mathbf{w}_k^{0,l})} + \sqrt{\Phi_k(\mathbf{w}_k^{0,l+1})} \right) \sum_{l=\tau}^{T-1} \left(\sqrt{\gamma \Phi_k(\mathbf{w}_k^{0,l})} - \sqrt{\Phi_k(\mathbf{w}_k^{0,l+1})} \right) \right]^{1/2} \\ &\stackrel{(b)}{\leq} \sqrt{\eta} \left[\sum_{l=\tau}^{T-1} \left(2\sqrt{\gamma \Phi_k(\mathbf{w}_k^{0,l})} \right) \sum_{l=\tau}^{T-1} \left(\sqrt{(1 + \eta^2 L_{\max} l'_{\max}) \Phi_k(\mathbf{w}_k^{0,l})} - \sqrt{\Phi_k(\mathbf{w}_k^{0,l+1})} \right) \right]^{1/2}, \end{aligned}$$

where (a) follows from Cauchy-Schwarz inequality and (b) follows from the second part of Lemma E.4. In the last inequality above, we have substituted the $\gamma = 1 + \eta^2 L_{\max} l'_{\max}$. Further, we can write the above inequality as

$$\begin{aligned} \sum_{l=\tau}^{T-1} \eta \|\nabla \Phi_k(\mathbf{w}_k^{0,l})\| &\leq \sqrt{\eta} \left(\sum_{l=\tau}^{T-1} \left(2\sqrt{\gamma \Phi_k(\mathbf{w}_k^{0,l})} \right) \right) \left[\sum_{l=\tau}^{T-1} \left(\sqrt{\Phi_k(\mathbf{w}_k^{0,l})} - \sqrt{\Phi_k(\mathbf{w}_k^{0,l+1})} \right) \right. \\ &\quad \left. + \sum_{l=\tau}^{T-1} \frac{\eta^2 L_{\max} l'_{\max}}{2} \sqrt{\Phi_k(\mathbf{w}_k^{0,l})} \right]^{1/2} \\ &\leq \sqrt{\eta} \left[\sum_{l=\tau}^{T-1} \left(2\sqrt{\gamma \Phi_k(\mathbf{w}_k^{0,l})} \right) \right]^{1/2} \left[\sqrt{\Phi_k(\mathbf{w}_k^{0,\tau})} + \frac{\eta^2 L_{\max} l'_{\max}}{2} \sum_{l=\tau}^{T-1} \sqrt{\Phi_k(\mathbf{w}_k^{0,l})} \right]^{1/2}, \end{aligned}$$

where the first inequality follows from the result $(1+x)^{\frac{1}{2}} \leq 1 + \frac{x}{2}$ whereas in the second inequality,

we have reduced the telescopic sum and ignored the positive quantity $\sqrt{\Phi_k(\mathbf{w}_k^{0,T})}$. We can further reduce the above inequality using distributive law as follows

$$\begin{aligned} \sum_{l=\tau}^{T-1} \eta \|\nabla \Phi_k(\mathbf{w}_k^{0,l})\| &\leq \sqrt{\eta} \left[2\sqrt{\gamma \Phi_k(\mathbf{w}_k^{0,\tau})} \sum_{l=\tau}^{T-1} \sqrt{\Phi_k(\mathbf{w}_k^{0,l})} + 2\sqrt{\gamma} \times \frac{\eta^2 L_{\max} l'_{\max}}{2} \left(\sum_{l=\tau}^{T-1} \sqrt{\Phi_k(\mathbf{w}_k^{0,l})} \right)^2 \right]^{1/2} \\ &\leq \sqrt{\eta} \left[2\sqrt{\gamma} \Phi_k(\mathbf{w}_0) \left(1 - \frac{\alpha_{\min} \eta}{2} \right)^{\frac{\tau}{2}} \sum_{l=\tau}^{T-1} \left(1 - \frac{\alpha_{\min} \eta}{2} \right)^{\frac{l}{2}} + \sqrt{\gamma} \eta^2 L_{\max} l'_{\max} \Phi_k(\mathbf{w}_0) \right. \\ &\quad \left. \left(\sum_{l=\tau}^{T-1} \left(1 - \frac{\alpha_{\min} \eta}{2} \right)^{\frac{l}{2}} \right)^2 \right]^{1/2}, \end{aligned} \tag{18}$$

where the above inequality follows from the Lemma E.3. Next, for $x \in [0, 1]$, we have the inequality $1 - x \leq (1 - \frac{x}{2})^2$. We use this inequality to upper bound the last term of the above inequality as follows

$$\begin{aligned} \sum_{l=\tau}^{T-1} \left(1 - \frac{\alpha_{\min} \eta}{2} \right)^{\frac{l}{2}} &\leq \left(1 - \frac{\alpha_{\min} \eta}{2} \right)^{\frac{\tau}{2}} \sum_{q=0}^{\infty} \left(1 - \frac{\alpha_{\min} \eta}{2} \right)^{\frac{q}{2}} \\ &\leq \left(1 - \frac{\alpha_{\min} \eta}{2} \right)^{\frac{\tau}{2}} \sum_{q=0}^{\infty} \left(1 - \frac{\alpha_{\min} \eta}{4} \right)^q \\ &= \left(1 - \frac{\alpha_{\min} \eta}{2} \right)^{\frac{\tau}{2}} \times \frac{4}{\alpha_{\min} \eta}. \end{aligned}$$

Next, we use the above upper-bound to equation 18, which will reduce the upper bound to

$$\begin{aligned}
\sum_{l=\tau}^{T-1} \eta \|\nabla \Phi_k(\mathbf{w}_k^{0,l})\| &\leq \sqrt{\eta} \left[2\sqrt{\gamma} \Phi_k(\underline{\mathbf{w}}^0) \left(1 - \frac{\alpha_{\min} \eta}{2}\right)^\tau \times \frac{4}{\alpha_{\min} \eta} + \sqrt{\gamma} \eta^2 L_{\max} l'_{\max} \Phi_k(\underline{\mathbf{w}}^0) \right. \\
&\quad \left. \left(\left(1 - \frac{\alpha_{\min} \eta}{2}\right)^{\frac{\tau}{2}} \times \frac{4}{\alpha_{\min} \eta} \right)^2 \right]^{1/2} \\
&= \left[2\eta \sqrt{\gamma} \Phi_k(\underline{\mathbf{w}}^0) \left(1 - \frac{\alpha_{\min} \eta}{2}\right)^\tau \times \frac{4}{\alpha_{\min} \eta} + \sqrt{\gamma} \eta^3 L_{\max} l'_{\max} \Phi_k(\underline{\mathbf{w}}^0) \right. \\
&\quad \left. \left(\left(1 - \frac{\alpha_{\min} \eta}{2}\right)^{\frac{\tau}{2}} \times \frac{4}{\alpha_{\min} \eta} \right)^2 \right]^{1/2} \\
&= \left[\frac{8\sqrt{1 + \eta^2 L_{\max} l'_{\max}} \Phi_k(\underline{\mathbf{w}}^0)}{\alpha_{\min}} + \frac{16\eta L_{\max} l'_{\max} \sqrt{1 + \eta^2 L_{\max} l'_{\max}} \Phi_k(\underline{\mathbf{w}}^0)}{\alpha_{\min}^2} \right]^{\frac{1}{2}} \\
&\quad \times \left(1 - \frac{\alpha_{\min} \eta}{2}\right)^{\frac{\tau}{2}}.
\end{aligned}$$

If we choose $\eta < \sqrt{\frac{3}{L_{\max} l'_{\max}}}$, the above inequality reduces to

$$\sum_{l=\tau}^{T-1} \eta \|\nabla \Phi_k(\mathbf{w}_k^{0,l})\| \leq \left(\frac{16}{\alpha_{\min}} + \frac{32\eta L_{\max} l'_{\max}}{\alpha_{\min}^2} \right)^{\frac{1}{2}} \sqrt{\Phi_k(\underline{\mathbf{w}}_0)} \times \left(1 - \frac{\alpha_{\min} \eta}{2}\right)^{\frac{\tau}{2}}.$$

Again choosing $\eta < \frac{\alpha_{\min}}{2L_{\max} l'_{\max}}$, we get

$$\begin{aligned}
\sum_{l=\tau}^{T-1} \eta \|\nabla \Phi_k(\mathbf{w}_k^{0,l})\| &\leq \left(\frac{16}{\alpha_{\min}} + \frac{16}{\alpha_{\min}} \right)^{\frac{1}{2}} \sqrt{\Phi_k(\underline{\mathbf{w}}_0)} \times \left(1 - \frac{\alpha_{\min} \eta}{2}\right)^{\frac{\tau}{2}} \\
&\leq \sqrt{\frac{32\Phi_k(\underline{\mathbf{w}}_0)}{\alpha_{\min}}} \left(1 - \frac{\alpha_{\min} \eta}{2}\right)^{\frac{\tau}{2}}.
\end{aligned}$$

□

Now, we are ready to complete the proof of Theorem E.1.

Proof. (Proof of Theorem E.1): Applying the Lemma E.5 with $\tau = 0$ and assumed condition, we have

$$\begin{aligned}
\|\mathbf{w}_k^{0,T} - \mathbf{w}_k^{0,0}\| &\leq \sum_{\tau=0}^{T-1} \eta \|\nabla \Phi_k(\mathbf{w}_k^{0,\tau})\| \\
&\leq \sqrt{\frac{32\Phi_k(\underline{\mathbf{w}}_0)}{\alpha_{\min}}} < \rho.
\end{aligned}$$

The last inequality follows from equation 14. This shows that $\mathbf{w}_k^{0,T} \in \mathbb{B}(\underline{\mathbf{w}}_0, \rho)$, which by induction is true for all $T \in \mathbb{N}$. This completes the proof. □

□

F PROOF OF THEOREM 3.1

To prove the Theorem 3.1, the foremost requirement is to show that the local and the global updates of FedAvg Algorithm 1 at each local round T and the global rounds R stay within a closed ball of radius $\rho > 0$ and centre at initialization $\underline{\mathbf{w}}^0$, i.e., $\mathbf{w}_k^{R,T} \in \mathbb{B}[\underline{\mathbf{w}}^0, \rho]$ for all $R \geq 0$ and $T \geq 0$. We use

the method of induction for two variables to prove this. By our initial hypothesis, the initialization $\mathbf{w}_k^{0,0} = \underline{\mathbf{w}}^0 \in \mathbb{B}[\underline{\mathbf{w}}^0, \rho]$, and hence the hypothesis is true for $R = 0$ and $T = 0$. First we show that the sequence $\{\mathbf{w}_k^{0,T}\}_{T \geq 0}$ stays within the ball for all T and $k \in [K]$ under the Assumption 2.1 and Assumption 2.4 [see Theorem E.1]. Towards this, assume the induction hypothesis that $\mathbf{w}_k^{0,1}, \mathbf{w}_k^{0,2}, \dots, \mathbf{w}_k^{0,T-1} \in \mathbb{B}[\underline{\mathbf{w}}^0, \rho]$ for all $k \in [K]$, then we need to prove that $\mathbf{w}_k^{0,T} \in \mathbb{B}[\underline{\mathbf{w}}^0, \rho]$. This is shown in Theorem E.1. Next, we assume that $\mathbf{w}_k^{1,T}, \mathbf{w}_k^{2,T}, \dots, \mathbf{w}_k^{R-1,T} \in \mathbb{B}[\underline{\mathbf{w}}^0, \rho]$ as the induction hypothesis on the first variable, and prove that $\mathbf{w}_k^{R,T} \in \mathbb{B}[\underline{\mathbf{w}}^0, \rho]$. This is equivalent to saying $\|\mathbf{w}_k^{R,T} - \mathbf{w}_k^{0,0}\| \leq \rho$ for all $k \in [K]$. This requires the following set of lemmas.

Lemma F.1. For $0 \leq r \leq R-1$ and $\tau \leq T$, we have for all $k \in [K]$,

$$\Phi_k(\mathbf{w}_k^{r,\tau}) \leq \left(1 - \frac{\eta \alpha_{\min}}{2}\right)^\tau \Phi_k(\mathbf{w}_k^{r,0}),$$

where $\alpha_{\min} := \min\{\alpha_1(\underline{\mathbf{w}}^0, \rho), \alpha_2(\underline{\mathbf{w}}^0, \rho), \dots, \alpha_K(\underline{\mathbf{w}}^0, \rho)\}$ provided $\eta < \frac{\alpha_{\min}}{4L_{\max}l'_{\max}}$ where $L_{\max} := \max_k L_k$ and $l'_{\max} := \max_k l'_k$.

Proof. For brevity, in this proof, we will use α_k in place of $\alpha_k(\underline{\mathbf{w}}^0, \rho)$. Since, $r \leq R-1$ and $\tau \leq T$, we have $\mathbf{w}_k^{r,\tau} \in \mathbb{B}[\underline{\mathbf{w}}^0, \rho]$ by induction hypothesis and Assumption 2.1 on Φ_k , we have

$$\begin{aligned} \Phi_k(\mathbf{w}_k^{r,\tau}) &\leq \Phi_k(\mathbf{w}_k^{r,\tau-1}) + \langle \mathbf{w}_k^{r,\tau} - \mathbf{w}_k^{r,\tau-1}, \nabla \Phi_k(\mathbf{w}_k^{r,\tau-1}) \rangle + \frac{\eta^2 L_k}{2} \|\mathbf{w}_k^{r,\tau} - \mathbf{w}_k^{r,\tau-1}\|^2 \\ &\leq \Phi_k(\mathbf{w}_k^{r,\tau-1}) - \eta \langle \widehat{\nabla} \Phi_k(\mathbf{w}_k^{r,\tau-1}), \nabla \Phi_k(\mathbf{w}_k^{r,\tau-1}) \rangle + \frac{\eta^2 L_k}{2} \|\widehat{\nabla} \Phi_k(\mathbf{w}_k^{r,\tau-1})\|^2. \end{aligned}$$

Next, we take expectation over mini-batch of samples at client $k \in [K]$. The above inequality reduces to

$$\begin{aligned} \Phi_k(\mathbf{w}_k^{r,\tau}) &\leq \Phi_k(\mathbf{w}_k^{r,\tau-1}) - \eta \langle \nabla \Phi_k(\mathbf{w}_k^{r,\tau-1}), \nabla \Phi_k(\mathbf{w}_k^{r,\tau-1}) \rangle + \frac{\eta^2 L_k}{2} \mathbb{E} \|\widehat{\nabla} \Phi_k(\mathbf{w}_k^{r,\tau-1})\|^2 \\ &\stackrel{(a)}{\leq} \Phi_k(\mathbf{w}_k^{r,\tau-1}) - \eta \alpha_k \Phi_k(\mathbf{w}_k^{r,\tau-1}) + 2\eta^2 L_k l'_k \Phi_k(\mathbf{w}_k^{r,\tau-1}) \\ &\leq (1 - \eta \alpha_{\min} + 2\eta^2 L_k l'_k) \Phi_k(\mathbf{w}_k^{r,\tau-1}) \\ &\leq (1 - \eta \alpha_{\min} + 2\eta^2 L_{\max} l'_{\max}) \Phi_k(\mathbf{w}_k^{r,\tau-1}), \end{aligned}$$

where (a) follows from the Definition 2.3 and Lemma D.2. Last inequality follows as $L_k \leq L_{\max}$ and $l'_k \leq l'_{\max}$ for all $k \in [K]$. Using $\eta < \frac{\alpha_{\min}}{4L_{\max}l'_{\max}}$, $\underline{\mathbf{w}}^r = \mathbf{w}_k^{r,0}$, and recursively iterating over τ , we have the desired upper bound

$$\Phi_k(\mathbf{w}_k^{r,\tau}) \leq \left(1 - \frac{\eta \alpha_{\min}}{2}\right)^\tau \Phi_k(\underline{\mathbf{w}}^r).$$

□

The above lemma gives an upper bound on the loss function at each client $k \in [K]$ for which $\mathbf{w}_k^{r,\tau} \in \mathbb{B}[\underline{\mathbf{w}}^0, \rho]$ for $r \leq R-1$ and $\tau \leq T$. In the next Lemma, we prove an upper bound for the same when $r > R-1$.

Lemma F.2. For $r \geq R$, we have for all $k \in [K]$

$$\Phi_k(\mathbf{w}_k^{r,\tau}) \leq (1 + 3\eta L_k)^\tau \Phi_k(\mathbf{w}_k^{r,0}),$$

provided $\eta < \frac{1}{2l'_{\max}}$.

Proof. For $r \geq R$, the induction hypothesis is not valid. Therefore, by the Assumption 2.1 i.e., Φ_k is L_k -Smooth, we can write

$$\begin{aligned} \Phi_k(\mathbf{w}_k^{r,\tau}) &\leq \Phi_k(\mathbf{w}_k^{r,\tau-1}) + \langle \mathbf{w}_k^{r,\tau} - \mathbf{w}_k^{r,\tau-1}, \nabla \Phi_k(\mathbf{w}_k^{r,\tau-1}) \rangle + \frac{\eta^2 L_k}{2} \|\mathbf{w}_k^{r,\tau} - \mathbf{w}_k^{r,\tau-1}\|^2 \\ &\leq \Phi_k(\mathbf{w}_k^{r,\tau-1}) - \eta \langle \widehat{\nabla} \Phi_k(\mathbf{w}_k^{r,\tau-1}), \nabla \Phi_k(\mathbf{w}_k^{r,\tau-1}) \rangle + \frac{\eta^2 L_k}{2} \|\widehat{\nabla} \Phi_k(\mathbf{w}_k^{r,\tau-1})\|^2. \end{aligned}$$

Taking expectation over mini-batch of samples reduces the above inequality to

$$\begin{aligned}\Phi_k(\mathbf{w}_k^{r,\tau}) &\leq \Phi_k(\mathbf{w}_k^{r,\tau-1}) - \eta \|\nabla \Phi_k(\mathbf{w}_k^{r,\tau-1})\|^2 + \frac{\eta^2 L_k}{2} \mathbb{E} \|\widehat{\nabla \Phi_k}(\mathbf{w}_k^{r,\tau-1})\|^2 \\ &\leq (1 + 2\eta L_k + 2\eta^2 L_k l'_k) \Phi_k(\mathbf{w}_k^{r,\tau-1}) \\ &\leq (1 + 2\eta L_k + 2\eta^2 L_k l'_{\max}) \Phi_k(\mathbf{w}_k^{r,\tau-1}),\end{aligned}$$

where we have used the unbiased gradient assumption and the fact that $-\eta \|\nabla \Phi_k(\mathbf{w}_k^{r,\tau-1})\|^2 \leq 2\eta L_k \Phi_k(\mathbf{w}_k^{r,\tau-1})$; which can be easily derived from smoothness assumption 2.1. Also, the last term is upper-bounded by using Lemma D.2 Using $\eta < \frac{1}{2l'_{\max}}$ in the above results and iterating over τ results in

$$\Phi_k(\mathbf{w}_k^{r,\tau}) \leq (1 + 3\eta L_k)^\tau \Phi_k(\mathbf{w}_k^{r,0}).$$

□

Next, we define the error of linear approximation of global loss function as Γ_r and show that it is bounded above.

Lemma F.3. Define the error as $\Gamma_r = \Phi(\underline{\mathbf{w}}^{r+1}) - \Phi(\underline{\mathbf{w}}^r) + \eta \|\nabla \Phi(\underline{\mathbf{w}}^r)\|^2$. Then for $r \leq R-1$, $\tau \leq T$ and $T \geq 3$, we have

$$\Gamma_r \leq \eta \left(\alpha_g - \frac{\alpha_g T}{4} \right) \Phi(\underline{\mathbf{w}}^r)$$

provided $\eta < \min \left\{ \frac{1}{3L_{\max}T}, \frac{\alpha_g \alpha_{\min}}{4T(4L_{\max}^2 l'_{\max} + L l'_{\max} \alpha_{\min})} \right\}$, where $L_{\max} := \max_{k \in [K]} L_k$, $l'_{\max} := \max_{k \in [K]} l'_k$ and $\alpha_{\min} := \min_{k \in [K]} \alpha_k$.

Proof. Since, for $r \leq R-1$ and $\tau \leq T$, we have $\mathbf{w}_k^{r,\tau} \in \mathbb{B}[\underline{\mathbf{w}}^0, \rho]$. By using FedAvg update rule in Algorithm 1, we can write the linear approximation error Γ_r as

$$\begin{aligned}\Gamma_r &\leq \Phi \left(\underline{\mathbf{w}}^r - \frac{\eta}{K} \sum_{k=1}^K \sum_{\tau=0}^{T-1} \widehat{\nabla \Phi_k}(\mathbf{w}_k^{r,\tau}) \right) - \Phi(\underline{\mathbf{w}}^r) + \eta \|\nabla \Phi(\underline{\mathbf{w}}^r)\|^2 \\ &\leq -\eta \sum_{\tau=0}^{T-1} \left\langle \frac{1}{K} \sum_{k=1}^K \widehat{\nabla \Phi_k}(\mathbf{w}_k^{r,\tau}), \nabla \Phi(\underline{\mathbf{w}}^r) \right\rangle + \frac{L\eta^2}{2K^2} \left\| \sum_{k=1}^K \sum_{\tau=0}^{T-1} \widehat{\nabla \Phi_k}(\mathbf{w}_k^{r,\tau}) \right\|^2 + \eta \|\nabla \Phi(\underline{\mathbf{w}}^r)\|^2,\end{aligned}$$

where the above inequality follows from L -smoothness Assumption 2.1 for the global loss function Φ . Further, we can take the expectation over a mini-batch of samples and write the above inequality as

$$\begin{aligned}\Gamma_r &\leq -\eta \sum_{\tau=0}^{T-1} \left\langle \frac{1}{K} \sum_{k=1}^K \nabla \Phi_k(\mathbf{w}_k^{r,\tau}), \nabla \Phi(\underline{\mathbf{w}}^r) \right\rangle + \frac{L\eta^2}{2K^2} \mathbb{E} \left[\left\| \sum_{k=1}^K \sum_{\tau=0}^{T-1} \widehat{\nabla \Phi_k}(\mathbf{w}_k^{r,\tau}) \right\|^2 \right] + \eta \|\nabla \Phi(\underline{\mathbf{w}}^r)\|^2 \\ &\leq -\frac{\eta}{2} \sum_{\tau=0}^{T-1} \left\| \frac{1}{K} \sum_{k=1}^K \nabla \Phi_k(\mathbf{w}_k^{r,\tau}) \right\|^2 - \frac{\eta}{2} \sum_{\tau=0}^{T-1} \|\nabla \Phi(\underline{\mathbf{w}}^r)\|^2 + \eta \|\nabla \Phi(\underline{\mathbf{w}}^r)\|^2 \\ &\quad + \frac{L\eta^2 T}{2K} \sum_{k=1}^K \sum_{\tau=0}^{T-1} \mathbb{E} [\|\widehat{\nabla \Phi_k}(\mathbf{w}_k^{r,\tau})\|^2] + \frac{\eta}{2} \sum_{\tau=0}^{T-1} \left\| \frac{1}{K} \sum_{k=1}^K \nabla \Phi_k(\mathbf{w}_k^{r,\tau}) - \nabla \Phi(\underline{\mathbf{w}}^r) \right\|^2,\end{aligned}$$

where the above inequality follows from the identity $\|\mathbf{a} - \mathbf{b}\|^2 = \|\mathbf{a}\|^2 + \|\mathbf{b}\|^2 - 2\langle \mathbf{a}, \mathbf{b} \rangle$ and inequality $\|\sum_{i=1}^n \mathbf{a}_i\|^2 \leq n \sum_{i=1}^n \|\mathbf{a}_i\|^2$. The above inequality can be further simplified with the

help of smoothness as

$$\begin{aligned}
\Gamma_r &\stackrel{(a)}{\leq} -\eta \left(\frac{T}{2} - 1 \right) \|\nabla \Phi(\underline{\mathbf{w}}^r)\|^2 + \frac{\eta}{2K} \sum_{k=1}^K \sum_{\tau=0}^{T-1} \|\nabla \Phi_k(\mathbf{w}_k^{r,\tau}) - \nabla \Phi_k(\underline{\mathbf{w}}^r)\|^2 \\
&\quad + \frac{2Ll'_{\max}\eta^2 T}{K} \sum_{k=1}^K \sum_{\tau=0}^{T-1} \Phi_k(\mathbf{w}_k^{r,\tau}) \\
&\stackrel{(b)}{\leq} -\eta\alpha_g \left(\frac{T}{2} - 1 \right) \Phi(\underline{\mathbf{w}}^r) + \frac{\eta L_{\max}^2}{2K} \sum_{k=1}^K \sum_{\tau=0}^{T-1} \|\mathbf{w}_k^{r,\tau} - \underline{\mathbf{w}}^r\|^2 \\
&\quad + \frac{2Ll'_{\max}\eta^2 T}{K} \sum_{k=1}^K \sum_{\tau=0}^{T-1} (1 + 3\eta L_k)^\tau \Phi_k(\mathbf{w}_k^{r,0}) \\
&\leq -\eta\alpha_g \left(\frac{T}{2} - 1 \right) \Phi(\underline{\mathbf{w}}^r) + \frac{\eta L_{\max}^2}{2K} \sum_{k=1}^K \sum_{\tau=0}^{T-1} \|\mathbf{w}_k^{r,\tau} - \underline{\mathbf{w}}^r\|^2 \\
&\quad + \frac{2Ll'_{\max}\eta^2 T}{K} \sum_{k=1}^K \sum_{\tau=0}^{T-1} (1 + 3\eta L_{\max})^\tau \Phi_k(\mathbf{w}_k^{r,0}),
\end{aligned}$$

where the last term of inequality (a) follows from sample-wise smoothness Assumption 2.2 and unbiased estimate assumption of loss function Φ_k , here $l'_{\max} := \max_k l'_k$. The first term of the above inequality (b) is a direct consequence of the Assumption 2.4 as $T \geq 3$. The second term follows from the L_k -smoothness Assumption 2.1 for loss function Φ_k whereas the third term follows from Lemma F.2. Here, $L_{\max} := \max_k L_k$. Note that we have ignored the first term as it is negative. We can further reduce the above inequality Γ_r as follows,

$$\begin{aligned}
\Gamma_r &\stackrel{(a)}{\leq} -\eta\alpha_g \left(\frac{T}{2} - 1 \right) \Phi(\underline{\mathbf{w}}^r) + \frac{\eta L_{\max}^2}{2K} \sum_{k=1}^K \sum_{\tau=0}^{T-1} \left\| -\eta \sum_{t=0}^{\tau-1} \widehat{\nabla} \Phi_k(\mathbf{w}_k^{r,t}) \right\|^2 + Ll'_{\max} \eta^2 T^2 e \Phi(\underline{\mathbf{w}}^r) \\
&\stackrel{(b)}{\leq} -\eta\alpha_g \left(\frac{T}{2} - 1 \right) \Phi(\underline{\mathbf{w}}^r) + \frac{2\eta^3 L_{\max}^2 l'_{\max}}{K} \sum_{k=1}^K \sum_{\tau=0}^{T-1} \tau \sum_{t=0}^{\tau-1} \Phi_k(\mathbf{w}_k^{r,t}) + Ll'_{\max} \eta^2 T^2 e \Phi(\underline{\mathbf{w}}^r),
\end{aligned}$$

where the second term of inequality (a) follows from the telescopic sum for the SGD update rule, and in the last term, we have assumed the $\eta < \frac{1}{3L_{\max}T}$. The second term of inequality (b) follows from sample-wise smoothness Assumption 2.2. Next, we use Lemma F.1 in the above inequality to get the second term in the global loss function Φ .

$$\begin{aligned}
\Gamma_r &\leq -\eta\alpha_g \left(\frac{T}{2} - 1 \right) \Phi(\underline{\mathbf{w}}^r) + \frac{2\eta^3 L_{\max}^2 l'_{\max}}{K} \sum_{k=1}^K \sum_{\tau=0}^{T-1} \tau \sum_{t=0}^{\tau-1} \left(1 - \frac{\eta\alpha_{\min}}{2} \right)^t \Phi_k(\mathbf{w}_k^{r,0}) \\
&\quad + Ll'_{\max} \eta^2 T^2 e \Phi(\underline{\mathbf{w}}^r) \\
&\leq -\eta\alpha_g \left(\frac{T}{2} - 1 \right) \Phi(\underline{\mathbf{w}}^r) + \frac{4\eta^2 T^2 L_{\max}^2 l'_{\max}}{\alpha_{\min}} \Phi(\underline{\mathbf{w}}^r) + Ll'_{\max} \eta^2 T^2 e \Phi(\underline{\mathbf{w}}^r) \\
&\leq \left(\eta\alpha_g - \frac{\eta\alpha_g T}{2} + \eta^2 T^2 e \left(\frac{4L_{\max}^2 l'_{\max} + Ll'_{\max} \alpha_{\min}}{\alpha_{\min}} \right) \right) \Phi(\underline{\mathbf{w}}^r),
\end{aligned}$$

where the second inequality follows from the infinite geometric sum and the sum of the first T natural numbers. Using $\eta < \frac{\alpha_g \alpha_{\min}}{4T(4L_{\max}^2 l'_{\max} + Ll'_{\max} \alpha_{\min})}$, the above inequality boils down to desired result,

$$\Gamma_r \leq \eta \left(\alpha_g - \frac{\alpha_g T}{4} \right) \Phi(\underline{\mathbf{w}}^r).$$

□

In the next lemma, we will use the above error bound to show the linear convergence relation of Φ .

Lemma F.4. For each $r \leq R$ i.e, for $\underline{\mathbf{w}}^r \in \mathbb{B}[\underline{\mathbf{w}}^0, \rho]$, we have

$$\Phi(\underline{\mathbf{w}}^r) \leq \left(1 - \frac{\eta\alpha_g T}{4}\right)^R \Phi(\underline{\mathbf{w}}^0).$$

Proof. From the definition of Γ_r for $r \leq R$, we have

$$\begin{aligned} \Phi(\underline{\mathbf{w}}^{r+1}) &= \Phi(\underline{\mathbf{w}}^r) - \eta \|\nabla \Phi(\underline{\mathbf{w}}^r)\|^2 + \Gamma_r \\ &\leq \Phi(\underline{\mathbf{w}}^r) - \eta\alpha_g \Phi(\underline{\mathbf{w}}^r) + \eta \left(\alpha_g - \frac{\alpha_g T}{4}\right) \Phi(\underline{\mathbf{w}}^r) \\ &\leq \left(1 - \frac{\eta\alpha_g T}{4}\right) \Phi(\underline{\mathbf{w}}^r), \end{aligned}$$

where the second inequality follows from the definition 2.3 for global loss function Φ and Lemma F.3. Next, we iterate over r and this proves the claim of the Lemma F.4 i.e,

$$\Phi(\underline{\mathbf{w}}^r) \leq \left(1 - \frac{\eta\alpha_g T}{4}\right)^r \Phi(\underline{\mathbf{w}}^0).$$

□

The above Lemma F.4 is valid as long as $\underline{\mathbf{w}}^r \in \mathbb{B}[\underline{\mathbf{w}}^0, \rho]$ for $r \leq R$. Therefore, we need to show that $\underline{\mathbf{w}}^r \in \mathbb{B}[\underline{\mathbf{w}}^0, \rho]$ for all r . Towards this, we need the following drift to be bounded. Especially, for the case, when $r = R$ as $\underline{\mathbf{w}}_k^{R,T}$ may lie outside of the ball $\mathbb{B}[\underline{\mathbf{w}}^0, \rho]$.

Lemma F.5. For $r \leq R$ and $\eta < \min \left\{ \frac{1}{3L_{\max}T}, \frac{1}{T\sqrt{2\epsilon l'_{\max} K \Phi(\underline{\mathbf{w}}^0)}} \right\}$, we have

$$\|\underline{\mathbf{w}}_k^{r,T} - \underline{\mathbf{w}}_k^{r,0}\| \leq \zeta_\rho \rho$$

for some $\zeta_\rho \in (0, 1)$.

Proof. Writing the drift term as telescoping sum, we have

$$\begin{aligned} \|\underline{\mathbf{w}}_k^{r,T} - \underline{\mathbf{w}}_k^{r,0}\| &= \left\| \sum_{\tau=0}^{T-1} (\underline{\mathbf{w}}_k^{r,\tau+1} - \underline{\mathbf{w}}_k^{r,\tau}) \right\| \\ &\leq \eta \sqrt{2l'_k} \sum_{\tau=0}^{T-1} \sqrt{\Phi_k(\underline{\mathbf{w}}_k^{r,\tau})}, \end{aligned}$$

where the above inequality follows from the local SGD update rule and Lemma D.2 which is a consequence of sample-wise smoothness Assumption 2.1 of Φ_k for all $k \in [K]$. Next, we use Lemma F.2 as $\underline{\mathbf{w}}_k^{R,T}$ may lie outside of the ball $\mathbb{B}[\underline{\mathbf{w}}^0, \rho]$, which reduces the above inequality to,

$$\begin{aligned} \|\underline{\mathbf{w}}_k^{r,T} - \underline{\mathbf{w}}_k^{r,0}\| &\leq \eta \sqrt{2l'_k} \sum_{t=0}^{T-1} (1 + 3\eta L_k)^{t/2} \sqrt{\Phi_k(\underline{\mathbf{w}}_k^{r,0})} \\ &\leq \eta T \sqrt{2l'_k K} \left(1 + \frac{3\eta L_k}{2}\right)^T \sqrt{\Phi(\underline{\mathbf{w}}^r)} \\ &\leq \eta T \sqrt{2l'_k K \Phi(\underline{\mathbf{w}}^0)} \left(1 + \frac{3\eta L_{\max}}{2}\right)^T \left(1 - \frac{\eta\alpha_g T}{8}\right)^r \\ &\leq \eta T \sqrt{2l'_k K \Phi(\underline{\mathbf{w}}^0)} \exp\left(\frac{3\eta L_{\max} T}{2}\right) \left(1 - \frac{\eta\alpha_g T}{8}\right)^r, \end{aligned}$$

where the second inequality follows from the result $\sqrt{1+x} \leq (1+\frac{x}{2})$ for $x \in (0, 1)$ and $\Phi_k(\underline{\mathbf{w}}^r) \leq \sum_{k=1}^K \Phi_k(\underline{\mathbf{w}}^r)$. Since, $\underline{\mathbf{w}}^r \in \mathbb{B}[\underline{\mathbf{w}}^0, \rho]$ for all $r \leq R$ by the induction hypothesis. Therefore, by

Lemma F.1, we get the third inequality. Last inequality follows from $(1+x) \leq \exp(x)$ for all $x \geq 0$. Note that $L_k \leq L_{\max}$ for all $k \in [K]$. Further,

$$\begin{aligned} \|\mathbf{w}_k^{r,T} - \mathbf{w}_k^{r,0}\| &\leq \eta T \sqrt{2el'_k K \Phi(\underline{\mathbf{w}}^0)} \exp\left(-\frac{\eta r \alpha_g T}{8}\right) \\ &\leq \eta T \sqrt{2el'_{\max} K \Phi(\underline{\mathbf{w}}^0)}, \end{aligned} \quad (19)$$

where we have used the inequality $(1+x) \leq e^x$ for all $x \in \mathbb{R}$ and $\eta < \frac{1}{3L_{\max}T}$ in the first inequality. The last inequality follows from the fact $\exp(-x) \leq 1$ for all $x \geq 0$ and $l'_k \leq l'_{\max}$. Next, by choosing $\eta < \frac{\zeta_\rho \rho}{T \sqrt{2el'_{\max} K \Phi(\underline{\mathbf{w}}^0)}}$, we have

$$\|\mathbf{w}_k^{r,T} - \mathbf{w}_k^{r,0}\| < \zeta_\rho \rho.$$

□

Next, we upper bound the global drift for $r \leq R$, especially for $r = R$, to get the desired result of Theorem 3.1.

Lemma F.6. *For $r \leq R$, we have*

$$\|\underline{\mathbf{w}}^r - \underline{\mathbf{w}}^0\| \leq \frac{\sqrt{128el'_{\max} K \Phi(\underline{\mathbf{w}}^0)}}{\alpha_g}, \quad (20)$$

provided $\eta < \frac{1}{3L_{\max}T}$, where $L_{\max} := \max_{k \in [K]} L_k$ and $l'_{\max} := \max_k l'_k$.

Proof. Consider the following

$$\begin{aligned} \|\underline{\mathbf{w}}^r - \underline{\mathbf{w}}^0\| &= \left\| \sum_{s=0}^{r-1} (\underline{\mathbf{w}}^{s+1} - \underline{\mathbf{w}}^s) \right\| \\ &\leq \sum_{s=0}^{r-1} \left\| \underline{\mathbf{w}}^s - \frac{\eta}{K} \sum_{k=1}^K \sum_{\tau=0}^{T-1} \widehat{\nabla} \Phi_k(\mathbf{w}_k^{s,\tau}) - \underline{\mathbf{w}}^s \right\| \\ &\leq \frac{\eta}{K} \sum_{k=1}^K \sum_{s=0}^{r-1} \sum_{\tau=0}^{T-1} \left\| \widehat{\nabla} \Phi_k(\mathbf{w}_k^{s,\tau}) \right\|, \end{aligned}$$

where the first equality follows from the telescoping sum whereas the second inequality follows from triangle inequality and FedAvg update rule in Algorithm 1. The last inequality follows from the triangle inequality. The above can be further bounded as follows,

$$\|\underline{\mathbf{w}}^r - \underline{\mathbf{w}}^0\| \leq \frac{\eta \sqrt{2l'_{\max}}}{K} \sum_{k=1}^K \sum_{s=0}^{r-1} \sum_{\tau=0}^{T-1} \sqrt{\widehat{\Phi}_k(\mathbf{w}_k^{s,\tau})},$$

where we have used the samples-smoothness Assumption 2.2. Next, we take expectation over a mini-batch of samples and apply Jensen's inequality. We get

$$\begin{aligned} \|\underline{\mathbf{w}}^r - \underline{\mathbf{w}}^0\| &\leq \frac{\eta \sqrt{2l'_{\max}}}{K} \sum_{k=1}^K \sum_{s=0}^{r-1} \sum_{\tau=0}^{T-1} \sqrt{\Phi_k(\mathbf{w}_k^{s,\tau})} \\ &\leq \frac{\eta \sqrt{2l'_{\max}}}{K} \sum_{k=1}^K \sum_{s=0}^{r-1} \sum_{\tau=0}^{T-1} (1 + 3\eta L_k)^{\tau/2} \sqrt{\Phi_k(\underline{\mathbf{w}}^s)} \\ &\leq \frac{\eta \sqrt{2l'_{\max}}}{K} \sum_{k=1}^K \sum_{s=0}^{r-1} \sum_{\tau=0}^{T-1} (1 + 3\eta L_{\max})^{T/2} \sqrt{\Phi_k(\underline{\mathbf{w}}^s)} \\ &\leq \frac{\eta \sqrt{2l'_{\max}}}{K} \sum_{k=1}^K \sum_{s=0}^{r-1} \sum_{\tau=0}^{T-1} \exp\left(\frac{3\eta L_{\max} T}{2}\right) \sqrt{\Phi_k(\underline{\mathbf{w}}^s)}, \end{aligned}$$

where the second inequality follows from Lemma F.2 and $\mathbf{w}_k^{s,0} = \underline{\mathbf{w}}^s$. Also, $L_k \leq L_{\max}$. Next, we choose $\eta < \frac{1}{3L_{\max}T}$ and use the fact that $\sqrt{\Phi_k(\underline{\mathbf{w}}^s)} \leq \sqrt{\sum_{k=1}^K \Phi_k(\underline{\mathbf{w}}^s)}$ and the Lemma F.1 in the above inequality.

$$\begin{aligned} \|\underline{\mathbf{w}}^r - \underline{\mathbf{w}}^0\| &\leq \frac{\eta T \sqrt{2el'_{\max}}}{\sqrt{K}} \sum_{k=1}^K \sum_{s=0}^{r-1} \sqrt{\Phi(\underline{\mathbf{w}}^s)} \\ &\leq \frac{\eta T \sqrt{2el'_{\max}} K}{K} \sum_{k=1}^K \sum_{s=0}^{r-1} \left(1 - \frac{\eta \alpha_g T}{4}\right)^{s/2} \sqrt{\Phi(\underline{\mathbf{w}}^0)} \\ &\leq \frac{\sqrt{128el'_{\max}} K \Phi(\underline{\mathbf{w}}^0)}{\alpha_g}, \end{aligned}$$

where the final inequality is due to the result $(1-x)^{r/2} \leq (1-\frac{x}{2})^r$ for $x \in (0,1)$ and the infinite geometric sum. \square

F.1 PROOF OF THEOREM 3.1

Proof. From the Lemma F.5 and Lemma F.6, for $r = R$, we have

$$\begin{aligned} \|\mathbf{w}_k^{R,T} - \mathbf{w}_k^{0,0}\| &\leq \|\mathbf{w}_k^{R,T} - \mathbf{w}_k^{R,0}\| + \|\mathbf{w}_k^{R,0} - \mathbf{w}_k^{0,0}\| \\ &\leq \zeta \rho + \frac{\sqrt{128el'_{\max}} K \Phi(\underline{\mathbf{w}}^0)}{\alpha_g} \\ &\leq \rho, \end{aligned}$$

where the last inequality follows from the assumed new condition 2.4(b) on α_g assumed in Theorem 3.1. Therefore, by induction $\mathbf{w}_k^{R,T} \in \mathbb{B}[\underline{\mathbf{w}}^0, \rho]$ for $R \geq 0$ and $T \geq 0$. Consequently, $\underline{\mathbf{w}}^R \in \mathbb{B}[\underline{\mathbf{w}}^0, \rho]$ for all $R \geq 0$ and lemma F.4 holds for $R \geq 0$. Taking $R \rightarrow \infty$, results in $\Phi(\mathbf{w}^*) = 0$, where $\mathbf{w}^* \in \mathbb{B}(\underline{\mathbf{w}}^0)$ is the optimal point. \square

G ANALYTICAL AND EXPERIMENTAL JUSTIFICATION OF ASSUMPTION 4.4

By definition of $\lambda_{k,\min}(m)$ defined in equation 9, we have

$$\begin{aligned} \frac{\mathbb{E}_{\underline{\mathbf{w}}^0} \mathbf{e}_k^\top \mathbf{H}_k(\underline{\mathbf{w}}^0) \mathbf{H}_k(\underline{\mathbf{w}}^0)^\top \mathbf{e}_k}{\|\mathbf{e}_k\|^2} &= \sum_{i,j,p} u_{k,i} u_{k,j} v_p^2 G_p(k,i,j) \mathbf{x}_{k,i}^\top \mathbf{x}_{k,j} \\ &= \sum_{i,j} u_{k,i} u_{k,j} \sum_p G_p(k,i,j) \mathbf{x}_{k,i}^\top \mathbf{x}_{k,j} \\ &= m \sum_{i,j} \mathbb{E}_{\mathbf{w}_p^0} [\sigma'(\mathbf{w}_p^{0\top} \mathbf{x}_{k,i}) \sigma'(\mathbf{w}_p^{0\top} \mathbf{x}_{k,j})] u_{k,i} u_{k,j} \mathbf{x}_{k,i}^\top \mathbf{x}_{k,j}, \end{aligned}$$

where $G_p(k,i,j) := \mathbb{E}_{\mathbf{w}_p^0} [\sigma'(\mathbf{w}_p^{0\top} \mathbf{x}_{k,i}) \sigma'(\mathbf{w}_p^{0\top} \mathbf{x}_{k,j})]$ is independent of p , $u_{k,i} := \frac{e_{k,i}}{\|\mathbf{e}_k\|}$ and \mathbf{w}_p^0 is the p^{th} vector of size d in the $\underline{\mathbf{w}}^0$ vector. In the above, $e_{k,i}$ is the i^{th} entry of \mathbf{e}_k , and the last inequality follows from the fact that $G_p(i,j,k)$ is independent of p since the $\underline{\mathbf{w}}^0$ is sampled in an i.i.d. fashion from a Gaussian distribution (see the discussion above Lemma 4.3 or Algorithm 2). Based on the above, it seems reasonable to expect that both $\lambda_{k,\rho}^-(m)$ and $\lambda_{\rho}^-(m)$ scale with m provided the matrices $\mathbf{H}(\underline{\mathbf{w}}^0) \mathbf{H}(\underline{\mathbf{w}}^0)^\top$ and $\mathbf{H}(\underline{\mathbf{w}}^0) \mathbf{H}(\underline{\mathbf{w}}^0)^\top$ are of full ranks respectively, which is shown in Lemma 4.3.

H BOUND ON THE OUTPUT OF THE NN2

In this section, we state and prove a bound on the output of the NN. The following lemma will come in handy while proving the generalization result for unbounded loss.

Lemma H.1. *There exists a $NN_{\beta=1}$ with some weights $\mathbf{v} \in \{-1, 1\}^m$ in the output layer such that $|f_{\mathbf{w}, \mathbf{v}}(\mathbf{x})| \leq \Delta := \sqrt{2}D_\sigma d \sqrt{\frac{\rho^2 + 3m}{m} \log(4)}$ for any random data point $\mathbf{x} \in \mathcal{X}_k := \{\mathbf{x}_{k,i} : i \in [n]\}$, $k \in [K]$, sampled i.i.d. from $p_k(\mathbf{x})$, and $\mathbf{w} \in \mathbb{B}[\underline{\mathbf{w}}^0, \rho]$.*

Proof: Recall that the weights in the final layer of the $NN_{\beta=1}$, i.e., \mathbf{v} is sampled randomly. We will show that in the ensemble of $NN_{\beta=1}$, there exists at least one NN whose output is bounded by Δ regardless of the input, i.e., $\Pr \{\exists \mathbf{v} \in \{-1, 1\}^m : |f_{\mathbf{w}, \mathbf{v}}(\mathbf{x})| \leq \Delta\} > \frac{1}{2}^6$. This probability can be written as

$$\begin{aligned} \Pr \{\exists \mathbf{v} \in \{-1, 1\}^m : |f_{\mathbf{w}, \mathbf{v}}(\mathbf{x})| \leq \Delta\} &= 1 - \Pr \left\{ \bigcap_{\mathbf{v} \in \{-1, 1\}^m} |f_{\mathbf{w}, \mathbf{v}}(\mathbf{x})| \leq \Delta \right\} \\ &> 1 - \Pr \{|f_{\mathbf{w}, \mathbf{v}}(\mathbf{x})| \leq \Delta\}. \end{aligned} \quad (21)$$

Thus, we need to upper bound $\Pr \{|f_{\mathbf{w}, \mathbf{v}}(\mathbf{x})| \leq \Delta\}$. Let us start with the complement of the event for any $\mathbf{x} \sim p_k(\mathbf{x})$, i.e.,

$$\mathbb{P} \left[\left| \frac{1}{\sqrt{m}} \sum_{l=1}^m v_l \sigma(\mathbf{w}_l^\top \mathbf{x}) \right| > \Delta \right] = \mathbb{E}_{\mathbf{x} \sim p_k(\mathbf{x})} \mathbb{P} \left[\left| \frac{1}{\sqrt{m}} \sum_{l=1}^m v_l \sigma(\mathbf{w}_l^\top \mathbf{x}) \right| > \Delta \mid \mathbf{x} \right]. \quad (22)$$

Since $\sigma(\cdot)$ is a smooth function, it follows from the remainder form of Taylor's expansion around $\mathbf{w} = \mathbf{0}$ that⁷

$$\begin{aligned} \sigma(\mathbf{w}_l^\top \mathbf{x}) &= \sigma(\mathbf{0}) + \mathbf{w}_l^\top \nabla \sigma(\mathbf{w}_{l*}^\top \mathbf{x}) \\ &= \mathbf{w}_l^\top \nabla \sigma(\mathbf{w}_{l*}^\top \mathbf{x}), \end{aligned}$$

where \mathbf{w}_{l*} lies between the line joining the points $\mathbf{0} \in \mathbb{R}^d$, $\mathbf{w}_l \in \mathbb{R}^d$, and $\mathbf{x} \in \mathcal{X}_k$. Using the above approximation, we get the following upper bound over the ball $\mathbb{B}[\underline{\mathbf{w}}^0, \rho]$ for all $l \in [m]$

$$\left| \frac{v_l \sigma(\mathbf{w}_l^\top \mathbf{x})}{\sqrt{m}} \right| \leq \frac{\|\mathbf{w}_l\| d D_\sigma}{\sqrt{m}},$$

where we have used the assumed upper-bound on the activation function, i.e., $\sigma'(\mathbf{w}_l^\top \mathbf{x}) \leq D_\sigma$. As the random variable of interest is bounded, we get the following upper bound using Hoeffding's inequality in equation 22

$$\begin{aligned} \mathbb{P} \left[\left| \frac{1}{\sqrt{m}} \sum_{l=1}^m v_l \sigma(\mathbf{w}_l^\top \mathbf{x}) \right| > \Delta \right] &\leq \mathbb{E}_{p_k(\mathbf{x})} \left[\exp \left\{ -\frac{2\Delta^2}{\sum_{l=1}^m \left[2 \frac{\|\mathbf{w}_l\| d D_\sigma}{\sqrt{m}} \right]^2} \right\} \mid \mathbf{x} \right] \\ &\leq \mathbb{E} \left[\exp \left\{ -\frac{m\Delta^2}{2d^2 D_\sigma^2 (\|\mathbf{w} - \underline{\mathbf{w}}^0\|^2 + \|\underline{\mathbf{w}}^0\|^2)} \right\} \right] \\ &\leq \mathbb{E} \left[\exp \left\{ -\frac{m\Delta^2}{2d^2 \rho^2 D_\sigma^2 + 2d^2 D_\sigma^2 \|\underline{\mathbf{w}}^0\|^2} \right\} \right], \end{aligned}$$

where the expectation is with respect to the random initialization. Next, we upper-bound the right-hand side of the above inequality using the total expectation law as follows,

$$\begin{aligned} \mathbb{E} \left[\exp \left\{ -\frac{m\Delta^2}{2d^2 D_\sigma^2 (\rho^2 + \|\underline{\mathbf{w}}^0\|^2)} \right\} \right] &\leq \mathbb{E} \left[\exp \left\{ -\frac{m\Delta^2}{2d^2 D_\sigma^2 (\rho^2 + \|\underline{\mathbf{w}}^0\|^2)} \right\} \mid \|\underline{\mathbf{w}}^0\|^2 < \zeta \right] \\ &\quad + \mathbb{E} \left[\exp \left\{ -\frac{m\Delta^2}{2d^2 D_\sigma^2 (\rho^2 + \|\underline{\mathbf{w}}^0\|^2)} \right\} \mid \|\underline{\mathbf{w}}^0\|^2 > \zeta \right] \\ &\quad \times \mathbb{P}[\|\underline{\mathbf{w}}^0\|^2 > \zeta] \\ &\leq \exp \left\{ -\frac{m\Delta^2}{2d^2 D_\sigma^2 (\rho^2 + \zeta)} \right\} + \mathbb{P}[\|\underline{\mathbf{w}}^0\|^2 > \zeta] \end{aligned} \quad (23)$$

⁶We just need to show that this quantity is non-zero. However, without loss of optimality, we choose $1/2$.

⁷Note that the expansion around the origin does not contradict the previous argument that the ball should not contain the origin.

where the last inequality, we have used $\mathbb{E} \left[\exp \left\{ -\frac{m\Delta^2}{2d^2 D_\sigma^2 (\rho^2 + \|\underline{\mathbf{w}}^0\|^2)} \right\} \middle| \|\underline{\mathbf{w}}^0\|^2 > \zeta \right] \leq 1$. Using the Lemma in Appendix D and $\zeta = 4m$, equation 23 reduces to

$$\begin{aligned} \mathbb{E} \left[\exp \left\{ -\frac{m\Delta^2}{2d^2 D_\sigma^2 (\rho^2 + \|\underline{\mathbf{w}}^0\|^2)} \right\} \right] &\leq \exp \left\{ -\frac{m\Delta^2}{2d^2 D_\sigma^2 (\rho^2 + 4m)} \right\} + \exp \left\{ -\frac{dm}{2} \right\} \\ &\leq 2 \exp \left\{ -\frac{m}{2} \min \left(\frac{\Delta^2}{d^2 D_\sigma^2 (\rho^2 + 4m)}, \frac{d}{2} \right) \right\} \end{aligned} \quad (24)$$

Note that by choosing $\Delta < \frac{d^3 D_\sigma^2 (\rho^2 + m)}{2}$ ensures that the first term, i.e., $\frac{\Delta^2}{d^2 D_\sigma^2 (\rho^2 + 4m)}$ is the minimum. Further, by choosing $\Delta > \sqrt{2} D_\sigma d \sqrt{\frac{\rho^2 + 3m}{m} \log(4)}$ ensures that the above is less than $1/2$, as desired. This completes the proof. \square

Lemma H.2. *With a probability of at least $1 - \delta/2$, the loss function satisfies*

$$\Phi(\underline{\mathbf{w}}^0) \leq n \times b := \left[\frac{2D_\sigma^2 \rho^2 d \log(2n/\delta)}{m} + 2y_{\max}^2 \right]. \quad (25)$$

Proof. Consider the squared loss for the 2-layer NN defined in equation 6

$$\begin{aligned} \Phi_k(\mathbf{w}) &= \sum_{i=1}^n [f_{\mathbf{w},v}(\mathbf{x}_{k,i}) - \mathbf{y}_{k,i}]^2 \\ &\leq 2 \sum_{i=1}^n [(f_{\mathbf{w},v}(\mathbf{x}_{k,i}))^2 + (\mathbf{y}_{k,i})^2] \\ &\leq n \times \left[\frac{2D_\sigma^2 \rho^2 d \log(2n/\delta)}{m} + 2y_{\max}^2 \right], \end{aligned}$$

where the last inequality follows from Lemma H.1 with $\delta_1 = \frac{\delta}{2}$ and $y_{\max}^2 := \max_{i \in [n]} y_{k,i}^2$, and hence satisfies with a probability of at least $(1 - \frac{\delta}{2})$. \square

I PROOF OF THEOREM 4.5

We consider the following average loss function:

$$\Phi(\mathbf{w}) := \frac{1}{K} \sum_{k=1}^K \Phi_k(\mathbf{w}), \quad (26)$$

where $\Phi_k : \mathbb{R}^{md} \rightarrow \mathbb{R}$ is squared loss function for each client $k \in [K]$ and is defined as

$$\Phi_k(\mathbf{w}) = \sum_{i=1}^n [f_{\mathbf{w},v}(\mathbf{x}_{k,i}) - \mathbf{y}_{k,i}]^2.$$

In the following, we show that the above NN with initialization provided in Algorithm 2 satisfies the conditions provided in Assumption 2.4 (1) and (2) provided the number of neurons in the first layer scales as $\mathcal{O}(\frac{nK}{d})$. Note that Algorithm 2 uses the following matrix

$$\mathbf{J}_k(\mathbf{w}, \mathbf{v}) := \frac{1}{\sqrt{m}} \times \mathbf{H}_k(\mathbf{w}, \mathbf{v}) = \frac{1}{\sqrt{m}} \times \begin{bmatrix} v_1 \sigma'(\mathbf{w}_1^\top \mathbf{x}_{k,1}) \mathbf{x}_{k,1}^\top & \dots & v_m \sigma'(\mathbf{w}_m^\top \mathbf{x}_{k,1}) \mathbf{x}_{k,1}^\top \\ \vdots & \ddots & \vdots \\ v_1 \sigma'(\mathbf{w}_1^\top \mathbf{x}_{k,n}) \mathbf{x}_{k,n}^\top & \dots & v_m \sigma'(\mathbf{w}_m^\top \mathbf{x}_{k,n}) \mathbf{x}_{k,n}^\top \end{bmatrix},$$

where $k = 1, 2, \dots, K$. The size of the above matrix is $\mathbb{R}^{n \times md}$. Note that the neural coefficients v_i 's for the second layer are sampled from a uniform distribution as shown in Algorithm 2, and are the same across clients.

Consider a single hidden layer NN defined in Sec. 4. For clarity, with a slight abuse of notation, we use \mathbf{w} to represent the weights of the first layer while we use \mathbf{v} to denote the weights of the second

layer. Note that the conditions in Assumption 2.4 amount to finding a lower bound on the norm squared of the gradient in terms of the average loss function inside a ball of radius ρ centred at \underline{w}^0 . Therefore, we consider the gradient of the average loss function in 26

$$\begin{aligned}\nabla_w \Phi(w) &= \frac{1}{K} \sum_{k=1}^K \sum_{i=1}^n \nabla_w (f_w(\mathbf{x}_{k,i}) - y_{k,i})^2 \\ &= \frac{2}{K} \sum_{k=1}^K \sum_{i=1}^n \nabla_w (f_w(\mathbf{x}_{k,i}))^\top (f_w(\mathbf{x}_{k,i}) - y_{k,i}) \\ &= \frac{2}{K} \sum_{k=1}^K \mathbf{J}_k^\top(w) (\mathbf{f}_w(\mathbf{X}_k) - \mathbf{y}_k),\end{aligned}$$

where $\mathbf{J}_k(w) \in \mathbb{R}^{n \times md}$ is the Jacobian matrix with the weight of the output layer scaled by β , defined in Algorithm 2, of the loss function at each client $k \in [K]$, and is given by

$$\mathbf{J}_k(w) := \begin{bmatrix} \nabla_w f_w(\mathbf{x}_{k,1})^\top \\ \nabla_w f_w(\mathbf{x}_{k,2})^\top \\ \vdots \\ \nabla_w f_w(\mathbf{x}_{k,n})^\top \end{bmatrix} = \frac{1}{\sqrt{m}} \times \mathbf{H}_k(w)$$

and is explicitly defined in Sec. 4, and $[\mathbf{f}_w(\mathbf{X}_k) - \mathbf{y}_k] \in \mathbb{R}^n$ be the column vector of error on each of n input data available at client $k \in [K]$. For simplicity, we denote error vector for client $k \in [K]$ by \mathbf{e}_k , and is given by

$$\mathbf{e}_k := [\mathbf{f}_w(\mathbf{X}_k) - \mathbf{y}_k] = \begin{bmatrix} f_w(\mathbf{x}_{k,1}) - y_{k,1} \\ f_w(\mathbf{x}_{k,2}) - y_{k,2} \\ \vdots \\ f_w(\mathbf{x}_{k,n}) - y_{k,n} \end{bmatrix} \in \mathbb{R}^{n \times 1}.$$

Using the above notations, we can re-write the gradient of the average loss in a more compact matrix form as

$$\begin{aligned}\nabla \Phi(w) &= \frac{2}{K} \sum_{k=1}^K \mathbf{J}_k^\top(w) \mathbf{e}_k \\ &= \frac{2}{K} \mathbf{J}(w)^\top \mathbf{e},\end{aligned}\tag{27}$$

where $\mathbf{J}(w) := [\mathbf{J}_1(w), \mathbf{J}_2(w), \dots, \mathbf{J}_K(w)] \in \mathbb{R}^{nK \times md}$, and the column vector formed by concatenating the error vectors \mathbf{e}_k and is given by $\mathbf{e} = [\mathbf{e}_1, \mathbf{e}_2, \dots, \mathbf{e}_K]^\top \in \mathbb{R}^{nK}$. Now, using equation 27, the norm squared of the gradient can be written as

$$\begin{aligned}\|\nabla \Phi(w)\|^2 &= \frac{4}{K^2} (\mathbf{e}^\top [\mathbf{J}(w) \mathbf{J}(w)^\top] \mathbf{e}) \\ &= \frac{4}{K^2} \|\mathbf{J}(w)^\top \mathbf{e}\|^2.\end{aligned}\tag{28}$$

Next, we approximate the matrix $\mathbf{H}_k(w)$ around the origin \underline{w}^0 . Note that $(i, j)^{th}$ entry of the matrix $\mathbf{H}_k(w)$ is given by $[\mathbf{H}_k(w)]_{i,j} = v_j \sigma'(\mathbf{w}_j^\top \mathbf{x}_{k,i}) \mathbf{x}_{k,i}^\top$. Therefore,

$$\begin{aligned}\|(\mathbf{H}_k^\top(w) - \mathbf{H}_k^\top(\underline{w}^0)) \mathbf{e}_k\|_2^2 &= \sum_{i=1}^n \sum_{j=1}^m |v_j \sigma'(\mathbf{w}_j^\top \mathbf{x}_{k,i}) \mathbf{x}_{k,i}^\top - v_j \sigma'((\underline{w}^0)^\top \mathbf{x}_{k,i}) \mathbf{x}_{k,i}^\top|^2 \\ &= \sum_{i=1}^n \sum_{j=1}^m |\sigma'(\mathbf{w}_j^\top \mathbf{x}_{k,i}) - \sigma'((\underline{w}^0)^\top \mathbf{x}_{k,i})|^2 |e_{k,i}|^2 \|\mathbf{x}_{k,i}^\top\|^2 \\ &= \sum_{j=1}^m |\sigma'(\mathbf{w}_j^\top \mathbf{x}_{k,i}) - \sigma'((\underline{w}^0)^\top \mathbf{x}_{k,i})|^2 \sum_{i=1}^n |e_{k,i}|^2,\end{aligned}\tag{29}$$

where the above inequality follows from the fact that $\|\mathbf{x}_{k,i}^\top\|^2 = 1$, and $e_{k,i}$ is the i -th component of the vector \mathbf{e}_k . By Taylor's expansion around $\underline{\mathbf{w}}^0$, we can write

$$\sigma'(\mathbf{w}_j^\top \mathbf{x}_{k,i}) = \sigma'(\underline{\mathbf{w}}^0) + (\mathbf{w}_j - \underline{\mathbf{w}}_j^0)^\top \nabla \sigma'(\mathbf{w}_{j*}^\top \mathbf{x}_{k,i}) \quad (31)$$

for some \mathbf{w}_{j*} in the line joining \mathbf{w}_j and $\underline{\mathbf{w}}^0$. Now, using the above equation, we have

$$\begin{aligned} \|(\mathbf{H}_k^\top(\mathbf{w}) - \mathbf{H}_k^\top(\underline{\mathbf{w}}^0))\mathbf{e}_k\|_2^2 &= \sum_{i=1}^n \sum_{j=1}^m |\sigma'(\underline{\mathbf{w}}^0) + (\mathbf{w}_j - \underline{\mathbf{w}}_j^0)^\top \nabla \sigma'(\mathbf{w}_{j*}^\top \mathbf{x}_{k,i}) - \sigma'(\underline{\mathbf{w}}^0)|^2 |e_{k,i}|^2 \\ &= \sum_{i=1}^n \sum_{j=1}^m |(\mathbf{w}_j - \underline{\mathbf{w}}_j^0)^\top \nabla \sigma'(\mathbf{w}_{j*}^\top \mathbf{x}_{k,i})|^2 |e_{k,i}|^2 \\ &\leq \sum_{i=1}^n \sum_{j=1}^m \|\mathbf{w}_j - \underline{\mathbf{w}}_j^0\|^2 \|\nabla \sigma'(\mathbf{w}_{j*}^\top \mathbf{x}_{k,i})\|^2 |e_{k,i}|^2 \\ &\leq d\Delta_\sigma^2 \sum_{i=1}^n |e_{k,i}|^2 \left(\sum_{j=1}^m \|\mathbf{w}_j - \underline{\mathbf{w}}_j^0\|^2 \right) \\ &\leq d\Delta_\sigma^2 \rho^2 \|\mathbf{e}_k\|^2, \end{aligned}$$

where we have used the assumption (see Assumption 4.1) $\|\nabla \sigma'(\mathbf{w}_{j*}^\top \mathbf{x}_{k,i})\| \leq d\Delta_\sigma$. Summing the above for all $k \in [K]$, the following holds

$$\|(\mathbf{H}(\mathbf{w}) - \mathbf{H}(\underline{\mathbf{w}}^0))\mathbf{e}\|_F^2 \leq d\Delta_\sigma^2 \rho^2 \|\mathbf{e}\|^2, \quad (32)$$

where we have used the definition of $\|\mathbf{e}\|^2 = \sum_{k=1}^K \|\mathbf{e}_k\|^2$. Next, we can write

$$\|\mathbf{H}(\underline{\mathbf{w}}^0)\mathbf{e} - \mathbf{H}(\mathbf{w})\mathbf{e} + \mathbf{H}(\mathbf{w})\mathbf{e}\|^2 \leq 2\|\mathbf{H}(\underline{\mathbf{w}}^0)\mathbf{e} - \mathbf{H}(\mathbf{w})\mathbf{e}\|^2 + 2\|\mathbf{H}(\mathbf{w})\mathbf{e}\|^2.$$

The above equation can be re-written as

$$\|\mathbf{H}(\mathbf{w})\mathbf{e}\|^2 \geq \frac{1}{2}\|\mathbf{H}(\underline{\mathbf{w}}^0)\mathbf{e}\|^2 - \|(\mathbf{H}(\underline{\mathbf{w}}^0) - \mathbf{H}(\mathbf{w}))\mathbf{e}\|_F^2. \quad (33)$$

Using the equation equation 33 in equation 28,

$$\begin{aligned} \|\nabla \Phi(\mathbf{w})\|^2 &\geq \frac{4}{mK^2} \left[\frac{1}{2}\|\mathbf{H}(\underline{\mathbf{w}}^0)\mathbf{e}\|^2 - \|(\mathbf{H}(\underline{\mathbf{w}}^0) - \mathbf{H}(\mathbf{w}))\mathbf{e}\|_F^2 \right] \\ &\geq \frac{4}{mK^2} \left[\frac{1}{2}\|\mathbf{H}(\underline{\mathbf{w}}^0)\mathbf{e}\|^2 - d\Delta_\sigma^2 \rho^2 \|\mathbf{e}\|^2 \right] \\ &= \frac{4}{mK} \left[\frac{\|\mathbf{H}(\underline{\mathbf{w}}^0)\mathbf{e}\|^2}{2\|\mathbf{e}\|^2} - d\Delta_\sigma^2 \rho^2 \right] \frac{\|\mathbf{e}\|^2}{K}, \end{aligned}$$

where in the second inequality, we have used the upper bound from equation equation 32. Now, we take infimum over the ball $\mathbb{B}(\underline{\mathbf{w}}^0, \rho)$ on both side of the above equation

$$\inf_{\mathbf{w} \in \mathbb{B}(\underline{\mathbf{w}}^0, \rho)} \frac{\|\nabla \Phi(\mathbf{w})\|^2}{\Phi(\mathbf{w})} \geq \frac{2}{mK} \inf_{\mathbf{w} \in \mathbb{B}(\underline{\mathbf{w}}^0, \rho)} \frac{\mathbf{e}^\top \mathbf{H}(\underline{\mathbf{w}}^0) \mathbf{H}(\underline{\mathbf{w}}^0)^\top \mathbf{e}}{\|\mathbf{e}\|^2} - \frac{4d\Delta_\sigma^2 \rho^2}{mK}.$$

The above equation can be written as

$$\alpha_g(\underline{\mathbf{w}}^0, \rho) \geq \frac{2}{mK} \lambda_\rho^-(m) - \frac{4d\Delta_\sigma^2 \rho^2}{mK}, \quad (34)$$

where $\lambda_\rho^-(m) := \inf_{\mathbf{w} \in \mathbb{B}(\underline{\mathbf{w}}^0, \rho)} \frac{\mathbf{e}^\top \mathbf{H}(\underline{\mathbf{w}}^0) \mathbf{H}(\underline{\mathbf{w}}^0)^\top \mathbf{e}}{\|\mathbf{e}\|^2}$. Similarly, for the single-client setting, we have

$$\alpha_k(\underline{\mathbf{w}}^0, \rho) \geq 2 \frac{\lambda_{k,\rho}^-(m)}{m} - \frac{4d\Delta_\sigma^2 \rho^2}{m} \quad (35)$$

for all $k \in [K]$. Now, we need a bound on l'_{max} . However, we will find a bound on L instead, which may be of independent interest. Later, we find a bound on l'_{max} as a special case. Consider

$$\begin{aligned}\|\nabla\Phi(\mathbf{w})\|^2 &= \frac{4}{mK^2}\|\mathbf{H}(\mathbf{w})\mathbf{e}\|^2 \\ &\leq \frac{4}{mK^2}\{2\|(\mathbf{H}(\mathbf{w}) - \mathbf{H}(\underline{\mathbf{w}}^0))\mathbf{e}\|_F^2 + 2\|\mathbf{H}(\underline{\mathbf{w}}^0)\mathbf{e}\|^2\} \\ &\leq \frac{4}{mK^2}\{d\Delta_\sigma^2\rho^2\|\mathbf{e}\|^2 + 2\|\mathbf{H}(\underline{\mathbf{w}}^0)\mathbf{e}\|^2\},\end{aligned}$$

where the second inequality from the result $\|a+b\|^2 \leq 2\|a\|^2 + 2\|b\|^2$ and the last inequality follows from equation 32. Further we divide both side by $\Phi(\mathbf{w}) = \frac{\|\mathbf{e}\|^2}{K}$, we get

$$\frac{\|\nabla\Phi(\mathbf{w})\|^2}{\Phi(\mathbf{w})} \leq \frac{8}{mK} \sup_{\mathbf{w} \in \mathbb{B}[\underline{\mathbf{w}}^0, \rho]} \frac{\|\mathbf{H}(\underline{\mathbf{w}}^0)\mathbf{e}\|^2}{\|\mathbf{e}\|^2} + \frac{d\Delta_\sigma^2\rho^2}{mK} \quad (36)$$

By comparing the above inequality with $\|\nabla\Phi(\mathbf{w})\|^2 \leq 2L\Phi(\mathbf{w})$, we get

$$L = \frac{4}{mK}\lambda_\rho^+(m) + \frac{d\Delta_\sigma^2\rho^2}{2mK}, \quad (37)$$

where $\lambda_\rho^+(m) := \sup_{\mathbf{w} \in \mathbb{B}[\underline{\mathbf{w}}^0, \rho]} \frac{\|\mathbf{H}(\underline{\mathbf{w}}^0)\mathbf{e}\|^2}{\|\mathbf{e}\|^2}$.

Bound on l'_{max} (Special Case: $K = 1, n = 1$): For single client ($K = 1$) and single data point ($n = 1$), from equation 36 we have,

$$\frac{\|\nabla\Phi_{k,i}(\mathbf{w})\|^2}{\Phi_{k,i}(\mathbf{w})} \leq \frac{8}{m} \sup_{\mathbf{w} \in \mathbb{B}[\underline{\mathbf{w}}^0, \rho]} \frac{\|\mathbf{H}_k(\underline{\mathbf{w}}^0)e_{k,i}\|^2}{e_{k,i}^2} + \frac{d\Delta_\sigma^2\rho^2}{m} \quad (38)$$

By comparing the above inequality with $\|\nabla\Phi_k(\mathbf{w})\|^2 \leq 2l'_{max}\Phi(\mathbf{w})$, where $l'_{max} := \max_{k,i} l_{k,i}$, we get

$$l'_{max} = \frac{4}{m}\lambda_{\rho,max}^+(m) + \frac{d\Delta_\sigma^2\rho^2}{2m}, \quad (39)$$

where $\lambda_{\rho,max}^+(m) := \max_{i \in [n], k \in [K]} \left[\sup_{\mathbf{w} \in \mathbb{B}[\underline{\mathbf{w}}^0, \rho]} \frac{\|\mathbf{H}(\underline{\mathbf{w}}^0)e_{k,i}\|^2}{e_{k,i}^2} \right]$.

To satisfy the condition mentioned in Assumption equation 2.4, the NN should be designed such that the following two properties are satisfied:

- **Condition (1) of Assumption 2.4:** The first condition is $\alpha_k > \frac{32\Phi_k(\underline{\mathbf{w}}^0)}{\rho^2}$ for $k \in [K]$. Using equation 35, the condition is satisfied if

$$\frac{\lambda_{k,\rho}^-(m)}{m} > 2 \times \left[\frac{\Delta_\sigma^2 d \rho^2}{m} + \frac{8bn}{\rho^2} \right],$$

where we have used the fact that with a probability of at least $1 - \delta/2$, $\Phi_k(\underline{\mathbf{w}}^0) \leq n \times b$ (see Lemma H.2).

- **Condition (2) of Assumption 2.4:** The second condition is $\alpha_g(\underline{\mathbf{w}}^0, \rho) \geq \frac{\sqrt{128e l'_{max} K \Phi(\underline{\mathbf{w}}^0)}}{(1-\zeta_\rho)\rho}$. Using equation 34 and substituting for L equation 39, and rearranging, we get the following condition

$$\frac{\lambda_\rho^-(m)}{m} > \frac{4K}{(1-\zeta_\rho)\rho} \sqrt{ebnK \left(\frac{\lambda_{\rho,max}^+(m)}{m} + \frac{4\Delta_\sigma^2 d \rho^2}{m} \right)} + \frac{2\Delta_\sigma^2 d \rho}{m},$$

where again we have used the fact that $\Phi_k(\underline{\mathbf{w}}^0) \leq n \times b$ with a probability of at least $1 - \delta/2$ from Lemma H.2. This completes the proof. \square

J GENERALIZATION BOUND: PROOF OF THEOREM 5.2

Overview of the proof: Note that the weights of the last layer of the neural network is sampled i.i.d. from $\{-1, 1\}$. Hence, there are 2^m possible neural networks (equal probability). In order to prove generalization bound, we invoke McDiarmid's inequality, which requires the output to be bounded. First, we prove that there exists NNs, i.e., $\mathbf{v} \in \{-1, 1\}^m$ such that the output is bounded; see Lemma H.1 in Appendix H. We denote such realizations of \mathbf{v} by the set \mathcal{G}_v . Once the existence is proved, we select the NN for which the output is bounded, and prove generalization for the chosen neural network. This is done by conditioning on the event $\mathbf{v} \in \mathcal{G}_v$. One can notice that such conditioning appears in the Rademacher complexity expression as well. In the following, we provide the details.

First, consider the following empirical loss function conditioned on $\mathbf{v} \in \mathcal{G}_v$ ⁸

$$\Phi(\mathbf{w}, \mathbf{v}) := \frac{1}{K} \sum_{k=1}^K \sum_{i=1}^n \left(\frac{1}{\sqrt{m}} \sum_{j=1}^m v_j \sigma(\mathbf{w}_j^\top \mathbf{x}_{k,i}) - y_{k,i} \right)^2. \quad (40)$$

Proving a typical PAC-style result requires one to apply McDiarmid's inequality. Towards this, consider two sample data points $S = (S_1, S_2, \dots, S_K)$ and $S' = (S'_1, S'_2, \dots, S'_K)$ differing only by points $x_{k,i}$ in S_k and $x'_{k,i}$ in S'_k and define

$$\Psi(S_1, S_2, \dots, S_K) = \sup_{\mathbf{w} \in \mathbb{B}[\underline{\mathbf{w}}^0, \rho]} (\Phi(\mathbf{w}, \mathbf{v}) - \Phi_S(\mathbf{w}, \mathbf{v})).$$

Since the difference of the suprema is upper bounded by the supremum of the difference over the same set, we have

$$\begin{aligned} \Psi(S') - \Psi(S) &= \sup_{\mathbf{w} \in \mathbb{B}[\underline{\mathbf{w}}^0, \rho]} (\Phi(\mathbf{w}, \mathbf{v}) - \Phi_{S'}(\mathbf{w}, \mathbf{v})) - \sup_{\mathbf{w} \in \mathbb{B}[\underline{\mathbf{w}}^0, \rho]} (\Phi(\mathbf{w}, \mathbf{v}) - \Phi_S(\mathbf{w}, \mathbf{v})) \\ &\leq \sup_{\mathbf{w} \in \mathbb{B}[\underline{\mathbf{w}}^0, \rho]} (\Phi_S(\mathbf{w}, \mathbf{v}) - \Phi_{S'}(\mathbf{w}, \mathbf{v})) \\ &= \frac{1}{K} \sup_{\mathbf{w} \in \mathbb{B}[\underline{\mathbf{w}}^0, \rho]} \{ (f_{\mathbf{w}, \mathbf{v}}(\mathbf{x}_{k,i}) - y_{k,i})^2 - (f_{\mathbf{w}, \mathbf{v}}(\mathbf{x}'_{k,i}) - y'_{k,i})^2 \}, \\ &\leq \frac{2}{K} \left[\sup_{\mathbf{w} \in \mathbb{B}[\underline{\mathbf{w}}^0, \rho]} f_{\mathbf{w}, \mathbf{v}}^2(\mathbf{x}_{k,i}) + y_{max}^2 \right], \end{aligned} \quad (41)$$

where in the first inequality above, all the terms cancel each other except for different data points $x_{k,i}$ and $x'_{k,i}$. Now, we use Lemma H.1 in Appendix H to show that the output is bounded conditioned on $\mathbf{v} \in \mathcal{G}_v$; this allows us to apply McDiarmid's inequality. It is important to note that the data and \mathbf{v} are independent of each other. The proof of McDiarmid's inequality depends on the randomness in the data and hence, does not change even after conditioning on $\mathbf{v} \in \mathcal{G}_v$, which is equivalent to the output being bounded (see Lemma H.1). Invoking Lemma H.1, the following bound holds

$$\sup_{\mathbf{w} \in \mathbb{B}[\underline{\mathbf{w}}^0, \rho]} f_{\mathbf{w}, \mathbf{v}}^2(\mathbf{x}_{k,i}) \leq (\rho^2 + 3m) \frac{2D_\sigma^2 d^2 \log 4}{m}.$$

Using the above in equation 41, we get the following bound

$$\Psi(S') - \Psi(S) \leq \frac{2}{K} \left((\rho^2 + 3m) \frac{2D_\sigma^2 d^2 \log 4}{m} + y_{max}^2 \right).$$

Now, applying McDiarmid's inequality to $\Psi(S) - \mathbb{E}_{|\mathcal{G}_v}[\Psi(S)]$, the following inequality holds with a probability of at least $1 - \delta$ for any $\mathbf{w} \in \mathbb{B}[\underline{\mathbf{w}}^0, \rho]$ conditioned on \mathcal{G}_v

$$\Psi(S) \leq \mathbb{E}_{|\mathcal{G}_v}[\Psi(S)] + \left((\rho^2 + 3m) \frac{2D_\sigma^2 d^2 \log 4}{m} + y_{max}^2 \right) \sqrt{\frac{2n \log(\frac{1}{\delta})}{K}}.$$

The expectation term appearing in the above inequality can be upper bounded in terms of the Rademacher Complexity using a standard approach as in Mohri et al. (2019) to get

$$\Phi(\mathbf{w}, \mathbf{v}) \leq \Phi_S(\mathbf{w}, \mathbf{v}) + \frac{2n}{K} \sum_{k=1}^K \text{Rad}_k(\underline{\mathbf{w}}^0, \rho) + \left(\frac{2D_\sigma^2 d^2 (\log 4)(\rho^2 + 3m)}{m} + y_{max}^2 \right) \sqrt{\frac{2n \log(\frac{1}{\delta})}{K}},$$

where the Rademacher complexity is as in Definition 12. This completes the proof. \square

⁸For the ease of exposition, we do not show the conditioning explicitly.

K PROOF OF THEOREM 5.4 (BOUND ON THE RADEMACHER COMPLEXITY)

By the definition of Rademacher complexity of client $k \in [K]$, we have

$$\begin{aligned} n \cdot \text{Rad}_k(\underline{\mathbf{w}}^0, \rho) &:= \mathbb{E}_{\mathbf{v} \in \mathcal{G}_v} \left[\sup_{\mathbf{w} \in \mathbb{B}(\underline{\mathbf{w}}^0, \rho)} \sum_{i=1}^n \zeta_i f_{\mathbf{w}, \mathbf{v}}(\mathbf{x}_{k,i}) \right] \\ &= \mathbb{E}_{\mathbf{v} \in \mathcal{G}_v} \left[\sup_{\mathbf{w} \in \mathbb{B}(\underline{\mathbf{w}}^0, \rho)} \frac{1}{\sqrt{m}} \sum_{i=1}^n \zeta_i \sum_{l=1}^m v_l \sigma(\mathbf{w}_l^\top \mathbf{x}_{k,i}) \right]. \end{aligned} \quad (42)$$

Consider $\mathbb{P} \left[\sup_{\mathbf{w} \in \mathbb{B}(\underline{\mathbf{w}}^0, \rho)} \sum_{i=1}^n \zeta_i f_{\mathbf{w}}(\mathbf{x}) > \epsilon \mid \mathbf{v} \in \mathcal{G}_v \right]$. Now we use the standard procedure of using covering to reduce the above supremum to countable union. In particular, it turns out that there exist balls of radius θ centred at points \mathbf{w}_{l_j} , $j = 1, 2, \dots, N_{\theta, \rho}$ which covers $\mathbb{B}(\underline{\mathbf{w}}^0, \rho)$. Here, $N_{\theta, \rho} = \left(\frac{3\rho\sqrt{d}}{\theta} \right)^d$ (see Shalev-Shwartz & Ben-David (2014)). Using this, we get

$$\begin{aligned} \mathbb{P} \left[\sup_{\mathbf{w} \in \mathbb{B}(\mathbf{0}, \rho)} \sum_{i=1}^n \zeta_i f_{\mathbf{w}, \mathbf{v}}(\mathbf{x}) > \epsilon \mid G \right] &= \mathbb{P} \left[\bigcup_{j=1}^{N_{\theta, \rho}} \left\{ \sup_{\mathbf{w} \in \mathbb{B}(\mathbf{w}_{l_j}, \theta)} \sum_{i=1}^n \zeta_i f_{\mathbf{w}}(\mathbf{x}) > \epsilon \right\} \mid G \right] \\ &\leq \sum_{j=1}^{N_{\theta, \rho}} \mathbb{P} \left[\left\{ \sup_{\mathbf{w} \in \mathbb{B}(\mathbf{w}_{l_j}, \theta)} \sum_{i=1}^n \zeta_i f_{\mathbf{w}}(\mathbf{x}) > \epsilon \right\} \mid G \right], \end{aligned} \quad (43)$$

where the above inequality follows from the union bound. As a part of the covering principle, we approximate the NN function around the centre of the covering balls $\mathbb{B}(\mathbf{w}_{l_j}, \theta)$ using the remainder form of Taylor series:

$$\begin{aligned} f_{\mathbf{w}, \mathbf{v}}(\mathbf{x}) &= \frac{1}{\sqrt{m}} \sum_{l=1}^m v_l \sigma(\mathbf{w}_l^\top \mathbf{x}) \\ &= \frac{1}{\sqrt{m}} \sum_{l=1}^m v_l \left[\sigma(\mathbf{w}_{l_j}^\top \mathbf{x}) + (\mathbf{w}_{l_j} - \mathbf{w}_l)^\top \nabla \sigma(\mathbf{w}_{l_j}^\top \mathbf{x}) \right] \\ &\leq f_{\mathbf{w}_{l_j}, \mathbf{v}}(\mathbf{x}) + \frac{1}{\sqrt{m}} \sum_{l=1}^m \|\mathbf{w}_{l_j} - \mathbf{w}_l\| \cdot \|\nabla \sigma(\mathbf{w}_{l_j}^\top \mathbf{x})\| \\ &\leq f_{\mathbf{w}_{l_j}, \mathbf{v}}(\mathbf{x}) + \theta d D_\sigma \sqrt{m}, \end{aligned}$$

where the first inequality is due to the well-known result of the Cauchy-Schwarz Inequality, and the last follows from the assumed upper bound on the activation function $\sigma(\cdot)$. Using the above approximation around the centre of the covering ball, we have

$$\sup_{\mathbf{w} \in \mathbb{B}(\mathbf{w}_{l_j}, \theta)} \sum_{i=1}^n \zeta_i f_{\mathbf{w}, \mathbf{v}}(\mathbf{x}) \leq \sum_{i=1}^n \zeta_i f_{\mathbf{w}_{l_j}}(\mathbf{x}) + n\theta D_\sigma d \sqrt{m}.$$

Note that the event due to the above inequality is bigger than that of in equation 43. Therefore, by the monotonic property of probability measure, we have

$$\begin{aligned} \sum_{j=1}^{N_{\theta, \rho}} \mathbb{P} \left[\sup_{\mathbf{w} \in \mathbb{B}(\mathbf{w}_{l_j}, \theta)} \sum_{i=1}^n \zeta_i f_{\mathbf{w}, \mathbf{v}}(\mathbf{x}) > \epsilon \mid \mathbf{v} \in \mathcal{G}_v \right] &\leq \sum_{j=1}^{N_{\theta, \rho}} \mathbb{P} \left[\sum_{i=1}^n \zeta_i f_{\mathbf{w}_{l_j}, \mathbf{v}}(\mathbf{x}) > \epsilon - n\theta D_\sigma d \sqrt{m} \mid \mathbf{v} \in \mathcal{G}_v \right] \\ &\leq \sum_{j=1}^{N_{\theta, \rho}} \exp \left\{ -\frac{m(\epsilon - n\theta D_\sigma d \sqrt{m})^2}{(\rho^2 + 3m) D_\sigma^2 d^2 \log \left(\frac{2}{\delta_2} \right)} \right\} \\ &\leq N_{\theta, \rho} \exp \left\{ -\frac{m(\epsilon - n\theta D_\sigma d \sqrt{m})^2}{(\rho^2 + 3m) D_\sigma^2 d^2 \log 4} \right\}, \end{aligned}$$

where the second inequality follows from the well-known Hoeffding's Inequality. Note that the above is less than δ_1 when $\epsilon := n\theta D_\sigma d \sqrt{m} + \sqrt{\frac{(\rho^2 + 3m) D_\sigma^2 d^2 \log 4 \log(N_{\theta, \rho}/\delta_1)}{m}}$. Now, choosing

$\theta = \frac{1}{2nD_\sigma dm}$ leads to the following

$$\mathbb{P} \left\{ \sup_{\mathbf{w} \in \mathbb{B}[\underline{\mathbf{w}}^0, \rho]} \sum_{i=1}^n \zeta_i f_{\mathbf{w}, \mathbf{v}}(\mathbf{x}) \leq \frac{1}{2\sqrt{m}} + \sqrt{\frac{(\rho^2 + 3m)D_\sigma^2 d^2 \log 4 \log(N_{\theta, \rho}/\delta_1)}{m}} \middle| \mathbf{v} \in \mathcal{G}_v \right\} \geq 1 - \delta_1.$$

Now, we are ready to prove a bound on the Rademacher Complexity using the above inequalities.

Let $\widehat{\text{Rad}}_k(\underline{\mathbf{w}}^0, \rho) := \sup_{\mathbf{w} \in \mathbb{B}[\underline{\mathbf{w}}^0, \rho]} \sum_{i=1}^n \zeta_i f_{\mathbf{w}, \mathbf{v}}(\mathbf{x})$. By the total expectation law, we have

$$\begin{aligned} \mathbb{E}_{\zeta_i \in \{-1, +1\}} \left[\widehat{\text{Rad}}_k(\underline{\mathbf{w}}^0, \rho) | \mathbf{v} \in \mathcal{G}_v \right] &\leq \mathbb{E}[\widehat{\text{Rad}}_k(\underline{\mathbf{w}}^0, \rho) | \widehat{\text{Rad}}_k(\underline{\mathbf{w}}^0, \rho) > \epsilon, \mathbf{v} \in \mathcal{G}_v] \\ &\quad \times \mathbb{P} \left[\widehat{\text{Rad}}_k(\underline{\mathbf{w}}^0, \rho) > \epsilon | \mathbf{v} \in \mathcal{G}_v \right] \\ &\quad + \mathbb{E}[\widehat{\text{Rad}}_k(\underline{\mathbf{w}}^0, \rho) | \widehat{\text{Rad}}_k(\underline{\mathbf{w}}^0, \rho) < \epsilon | \mathbf{v} \in \mathcal{G}_v]. \end{aligned}$$

We bound the first term of the above inequality using the upper bound on the output of the NN (see Lemma H.1) function as follows

$$\mathbb{E}_{\zeta_i \in \{-1, +1\}} \left[n \cdot \widehat{\text{Rad}}_k(\underline{\mathbf{w}}^0, \rho) | \mathbf{v} \in \mathcal{G}_v \right] \leq n\sqrt{2}D_\sigma d \sqrt{\frac{\rho^2 + m}{m} (\log 4) \delta_1} + \epsilon.$$

We are free to choose δ_1 as it appears in the log term. Choosing $\delta_1 = \frac{1}{2mn\sqrt{2}D_\sigma d} \sqrt{\frac{m}{\log 4(\rho^2 + m)}}$ and dividing by n leads to

$$\text{Rad}_k(\underline{\mathbf{w}}^0, \rho) \leq \frac{1}{n\sqrt{m}} + \sqrt{\frac{(\rho^2 + 3m)D_\sigma^2 d^2 (\log 4) \log(N_{\theta, \rho}/\delta_1)}{mn}}.$$

Note that the first term above converges fast. However, the overall complexity is of the order $1/\sqrt{n}$ for any $\rho \leq \mathcal{O}(\sqrt{m})$. This completes the proof. \square