

# On the Unreasonable Effectiveness of Federated Averaging with Heterogeneous Data

**Jianyu Wang\***

*Carnegie Mellon University*

*jianyu.wang.thu@gmail.com*

**Rudrajit Das**

*University of Texas at Austin*

*rdas@utexas.edu*

**Gauri Joshi**

*Carnegie Mellon University*

*gaurij@andrew.cmu.edu*

**Satyen Kale**

*Google Research*

*satyen@satyenkale.com*

**Zheng Xu**

*Google Research*

*xuzheng@google.com*

**Tong Zhang**

*University of Illinois Urbana-Champaign*

*tongzhang@tongzhang-ml.org*

**Reviewed on OpenReview:** <https://openreview.net/forum?id=zF76Ga4EPs>

## Abstract

Existing theoretical results (such as (Woodworth et al., 2020a)) predict that the performance of federated averaging (FedAvg) is exacerbated by high data heterogeneity. However, in practice, FedAvg converges pretty well on several naturally heterogeneous datasets. In order to explain this seemingly unreasonable effectiveness of FedAvg that contradicts previous theoretical predictions, this paper introduces the *client consensus hypothesis*: the average of local models updates on clients starting from the optimum is close to zero. We prove that under this hypothesis, data heterogeneity does not exacerbate the convergence of FedAvg. Moreover, we show that this hypothesis holds for a linear regression problem and some naturally heterogeneous datasets such as FEMNIST and StackOverflow. Therefore, we believe that this hypothesis can better explain the performance of FedAvg in practice.

## 1 Introduction

Federated learning (FL) is an emerging distributed training paradigm (Kairouz et al., 2019; Wang et al., 2021), which enables a large number of clients to collaboratively train a powerful machine learning model without the need of uploading any raw training data. One of the most popular FL algorithms is Federated Averaging (FEDAVG), proposed by McMahan et al. (2017). In each round of FEDAVG, a small subset of clients are randomly selected to perform local model training. Then, the local model changes from clients are aggregated at the server to update the global model. This general local-update framework only requires infrequent communication between server and clients, and thus, is especially suitable for FL settings where the communication cost is a major bottleneck.

Due to its simplicity and empirical effectiveness, FEDAVG has become the basis of almost all subsequent FL optimization algorithms. Nonetheless, its convergence behavior – especially when the clients have

---

\*Currently at Apple.

heterogeneous data – has not been fully understood yet. Existing theoretical results such as Woodworth et al. (2020a); Glasgow et al. (2021) predict that FEDAVG’s convergence is greatly affected by the data heterogeneity. In particular, when the local gradients on clients become more different from each other (i.e., more data heterogeneity), FEDAVG may require much more communication rounds to converge. These theoretical predictions match with the observations on some datasets with artificially partitioned or synthetic non-IID data (Hsu et al., 2019; Li et al., 2020a; Zhao et al., 2018). However, on many real-world FL training tasks, FEDAVG actually performs extremely well (McMahan et al., 2017; Charles et al., 2021), which appears to be unreasonable based on the existing theory. In fact, many advanced methods aimed at mitigating the negative effects of data heterogeneity performs similar to FEDAVG in real-world FL training tasks. For example, SCAFFOLD (Karimireddy et al., 2020b) needs much fewer communication rounds to converge than FEDAVG in theory and when run on a synthetic dataset. However, SCAFFOLD and FEDAVG are reported to have roughly identical empirical performance on many realistic federated datasets, see the results in (Reddi et al., 2021; Wu et al., 2023; Li et al., 2021).

The above gap between theoretical results and practical observations motivates us to think about whether the current theoretical analyses are too pessimistic about FEDAVG, since most of them only focus on the worst cases. It remains as an open problem whether the data heterogeneity modeled by theory and simulated via artificially constructed non-IID datasets matches the heterogeneity in real-world applications. Is it possible that the data heterogeneity in practice has some special properties that allow FEDAVG to enjoy good convergence? The answers to the above questions will serve as important guidelines for the design and evaluation of future federated algorithms.

**Main Contributions.** In this paper, we take the first step towards explaining the practical effectiveness of FEDAVG under a special but realistic hypothesis. In particular, our contributions are as follows.

- By performing experiments on naturally heterogeneous federated datasets, we show that previous theoretical predictions do not align well with practice. FEDAVG can have nearly identical performance on both IID and non-IID versions of these datasets. Thus, previous worst-case analyses may be too pessimistic for such datasets. This is likely because the heterogeneity metric used by existing analyses may be too loose for naturally heterogeneous datasets.
- In order to explain the practical effectiveness of FEDAVG, we propose a *client consensus hypothesis*: the average of local model updates starting from the optimum is close to zero. For smooth and strongly-convex functions, we prove that under this hypothesis, there is no negative impact on the convergence of FEDAVG even with unbounded gradient dissimilarity.
- We further validate that the client consensus hypothesis can hold in many scenarios. We firstly consider a linear regression problem where all the clients have the same conditional probability of label given data, and show that the client consensus hypothesis holds. Besides, we empirically show that natural federated datasets such as FEMNIST and StackOverflow satisfy this hypothesis. Indeed, data heterogeneity has very limited impact on these datasets.

We would like to clarify that we are *not* trying to provide any practical approach/criterion to predict which datasets FedAvg will converge on. We are just diving deep into the notion of heterogeneity and trying to provide some insights on which notion seems to be more aligned with the behavior of FedAvg in practice.

## 2 Preliminaries and Related Work

**Problem Formulation.** We consider a FL setting with total  $M$  clients, where each client  $c$  has a local objective function  $F_c(\mathbf{w})$  defined on its local dataset  $\mathcal{D}_c$ . The goal of FL training is to minimize a global objective function, defined as a weighted average over all clients:

$$F(\mathbf{w}) = \sum_{c=1}^M p_c F_c(\mathbf{w}) = \mathbb{E}_c[F_c(\mathbf{w})], \tag{1}$$

where  $p_c$  is the relative weight for client  $c$ . For the ease of writing, in the rest of this paper, we will use  $\mathbb{E}_c[\mathbf{a}_c] = \sum_{c=1}^M p_c \mathbf{a}_c$  to represent the weighted average over clients.

**Update Rule of FedAvg.** FEDAVG (McMahan et al., 2017) is a popular algorithm to minimize (1) without the need of uploading raw training data. In round  $t$  of FEDAVG, client  $c$  performs  $H$  steps of SGD from the current global model  $\mathbf{w}^{(t)}$  to a local model  $\mathbf{w}_c^{(t,H)}$  with a local learning rate  $\eta$ . Then, at the server, the local model changes are aggregated to update the global model as follows:

$$\mathbf{w}^{(t+1)} = \mathbf{w}^{(t)} - \alpha \mathbb{E}_c[\mathbf{w}^{(t)} - \mathbf{w}_c^{(t,H)}]. \quad (2)$$

Here  $\alpha$  denotes the server learning rate. Our results in this paper are for the full-device participation case, i.e., when all the clients participate in each round. We discuss how our results can be extended to the partial-device participation case at the end of Section 4.2.

**Theoretical Analysis of FedAvg.** When clients have *homogeneous* data, many works have provided error upper bounds to guarantee the convergence of FEDAVG (also called Local SGD) (Stich, 2019; Yu et al., 2019b; Wang & Joshi, 2018; Zhou & Cong, 2018; Khaled et al., 2020; Li et al., 2019). In these papers, FEDAVG was treated as a method to reduce the communication cost in distributed training. It has been shown that in the stochastic setting, using a proper  $H > 1$  in FEDAVG will not negatively influence the dominant convergence rate. Hence FEDAVG can save communication rounds compared to the algorithm with  $H = 1$ . Later in Woodworth et al. (2020b), the authors compared FEDAVG with the mini-batch SGD baseline, and showed that in certain regimes, FEDAVG provably improves over mini-batch SGD. These upper bounds on FEDAVG was later proved by Glasgow et al. (2021) to be tight and not improvable for general convex functions.

When clients have *heterogeneous* data, in order to analyze the convergence of FEDAVG, it is common to make the following assumption to bound the difference among local gradients.

**Assumption 1** (Bounded Gradient Dissimilarity). *There exists a positive constant  $\zeta$  such that  $\forall \mathbf{w} \in \mathbb{R}^d$ , the difference between local and global gradients are uniformly bounded:*

$$\mathbb{E}_c \|\nabla F_c(\mathbf{w}) - \nabla F(\mathbf{w})\|^2 \leq \zeta^2. \quad (3)$$

This assumption first appeared in decentralized optimization literature (Lian et al., 2017; Yuan et al., 2016; Assran et al., 2018; Koloskova et al., 2020; Wang et al., 2022), and has been subsequently used in the analysis of FEDAVG and related algorithms (Yu et al., 2019a; Khaled et al., 2020; Karimireddy et al., 2020b; Reddi et al., 2020; Wang et al., 2020a;b; Haddadpour & Mahdavi, 2019; Karimireddy et al., 2020a; Das et al., 2022; Zindari et al., 2023; Crawshaw et al., 2024). Under the bounded gradient dissimilarity assumption, FEDAVG cannot outperform the simple mini-batch SGD baseline unless  $\zeta$  is extremely small ( $\zeta < 1/T$  where  $T$  is the total communication rounds) (Woodworth et al., 2020a); the deterministic version of FEDAVG (i.e., Local GD) has even slower convergence rate than vanilla GD (Khaled et al., 2020). Again, these upper bounds match a lower bound constructed in Glasgow et al. (2021) for general convex functions, suggesting that they are tight in the worst case. In this paper, we do not aim to improve these bounds, which are already tight. Instead, we argue that since the existing analyses only consider the worst case, they may be too pessimistic for practical applications. The data heterogeneity induced by the bounded gradient dissimilarity assumption may be different from the real-world heterogeneity.

Finally, we note that there is another line of works, namely, Malinovskiy et al. (2020); Charles & Konečný (2021); Charles & Rush (2022), using a different analysis technique from the above literature. They showed that FEDAVG (with full client participation) is equivalent to performing gradient descent on a surrogate loss function. However, so far this technique still has some limitations. It can only be applied to deterministic settings with quadratic (or a very special class of) loss functions. Additionally, Wang et al. (2023a) propose a heterogeneity-driven Lipschitz assumption to better capture the effect of local steps, and derive a convergence result for FEDAVG with this assumption. Gu et al. (2023) discuss the reasons why FEDAVG can have better generalization than mini-batch SGD.

**Comparisons of Data Heterogeneity Assumptions.** In Table 1, we summarize some commonly used data heterogeneity assumptions in literature. It is worth highlighting that previous literature, no matter

which heterogeneity measures they used, suggested that only when all local functions are the same or share the same optimum, there are no additional error terms (i.e.,  $\zeta = 0$ ) caused by data heterogeneity. However, in our paper, we show that even if local functions are heterogeneous and have different local optima, data heterogeneity can have no negative impact in some regimes (which likely happen in practice). This new result helps us gain a deeper understanding of the great performance of FedAvg observed in practice and cannot be obtained from previous works.

Among all these heterogeneity measures, gradient dissimilarity (at optimum) is considered as the most general one and widely used in literature (Wang et al., 2021). But as we will discuss in the paper, it can be pessimistic in practice. Besides, the gradient diversity assumption used in (Li et al., 2020a; Haddadpour & Mahdavi, 2019) implicitly forces all local functions share the same optimum. So it is much more restrictive than the gradient dissimilarity.

Table 1: Summary of data heterogeneity measures/assumptions used in literature. The above table is adapted from Table 6 in the survey (Kairouz et al., 2019). In the above works, their error upper bounds have the same dependency on  $\zeta$ , though  $\zeta$ 's definition is different.

Name	Definition	Example usage
grad. dissimilarity	$\mathbb{E}_c \ \nabla F_c(\mathbf{w}) - \nabla F(\mathbf{w})\  \leq \zeta^2$	Woodworth et al. (2020a)
grad. dissimilarity at opt.	$\mathbb{E}_c \ \nabla F_c(\mathbf{w}^*)\ ^2 \leq \zeta^2$	Khaled et al. (2020); Koloskova et al. (2020)
grad. diversity	$\mathbb{E}_c \ \nabla F_c(\mathbf{w})\ ^2 \leq \lambda^2 \ \nabla F(\mathbf{w})\ ^2$	Li et al. (2020a); Haddadpour & Mahdavi (2019)
general grad. diversity	$\mathbb{E}_c \ \nabla F_c(\mathbf{w})\ ^2 \leq \lambda^2 \ \nabla F(\mathbf{w})\ ^2 + \zeta^2$	Wang et al. (2020a); Karimireddy et al. (2020b)
grad. norm	$\ \nabla F_c(\mathbf{w})\ ^2 \leq \zeta^2$	Yu et al. (2019b); Li et al. (2020b)
opt. diff	$\ \mathbb{E}_c \mathbf{w}_c^* - \mathbf{w}^*\ ^2 \leq \zeta^2$	Wang et al. (2023b)

### 3 Mismatch Between Theory and Practice on FedAvg

In this section, we will introduce the existing convergence analysis in detail and compare the theory with experimental results. We shall show that there is a gap between the theory and practice of FEDAVG.

#### 3.1 Existing Theoretical Analysis of FedAvg

Besides the bounded gradient dissimilarity assumption (3), the analysis of FEDAVG relies on the following common assumptions.

**Assumption 2 (Lipschitz Smoothness).** *There exists a constant  $L$  such that,  $\forall \mathbf{w}, \mathbf{u} \in \mathbb{R}^d, \forall c$ ,*

$$\|\nabla F_c(\mathbf{w}) - \nabla F_c(\mathbf{u})\|^2 \leq L \|\mathbf{w} - \mathbf{u}\|^2 \quad (4)$$

**Assumption 3 (Local Unbiased Gradient).** *Local stochastic gradient on each client  $c$  is unbiased with expectation  $\nabla F_c(\mathbf{w})$  and bounded variance  $\sigma^2$ .*

Using the above assumptions, previous works derived an upper bound for the optimization error with non-convex objective functions. We take the following theorem from (Jhunjunwala et al., 2022) as an example.

**Theorem 1.** *Under Assumptions 1 to 3, if FEDAVG's learning rates satisfy  $\eta \leq 1/8LH, \alpha \leq 1/24LH$ , then the global gradient norm  $\min_t \mathbb{E} \|\nabla F(\mathbf{w}^{(t)})\|^2$  can be upper bounded by*

$$\mathcal{O} \left( \frac{F(\mathbf{w}^{(0)}) - F^*}{\alpha \eta H T} \right) + \mathcal{O} \left( \frac{\alpha \eta L \sigma^2}{M} + \eta^2 L^2 H \sigma^2 \right) + \mathcal{O}(\eta^2 L^2 H^2 \zeta^2). \quad (5)$$

Theorem 1 shows that when the learning rates  $\alpha, \eta$  are fixed and the communication round  $T$  is limited, data heterogeneity always introduces an addition term  $\mathcal{O}(\zeta^2)$  to the optimization error bound. It is worth noting that, with an optimized learning rate, the above upper bound matches a lower bound constructed in Glasgow et al. (2021) for general convex functions, suggesting that it is tight. In the worst case, the convergence of

FEDAVG will become worse compared to the homogeneous setting. Although some papers argued that the data heterogeneity only influences the higher order terms when using optimized learning rates (Yu et al., 2019a; Khaled et al., 2020), this conclusion only holds asymptotically w.r.t.  $T$ . Also, in some special cases, FEDAVG’s convergence rate with data heterogeneity can be substantially slower. For instance, when there is no stochastic noise  $\sigma = 0$ , FEDAVG in homogeneous setting can achieve a rate of  $T^{-1}$  but with heterogeneous data, the rate degrades to  $T^{-2/3}$  (Woodworth et al., 2020a; Wang et al., 2021; Jhunjhunwala et al., 2022).

### 3.2 Empirical Observations on FedAvg

**Results on Artificial Non-IID Datasets.** In order to corroborate the above theoretical results, most of previous works constructed datasets with artificially high data heterogeneity. Given a common benchmark classification dataset (such as MNIST, CIFAR-10/100), researchers simulated the data heterogeneity by assigning different class distributions to clients. For example, each client may only hold one or very few classes of data (Zhao et al., 2018), or has data for all classes but the amount of each class is randomly drawn from a Dirichlet distribution (Hsu et al., 2019). On these datasets, the empirical convergence of FEDAVG is greatly impacted and its final accuracy is much lower than the one in the centralized training setting. As shown in Zhao et al. (2018), FEDAVG’s final test accuracy on non-IID versions of CIFAR-10 can be 10 – 40% lower than its on the standard IID version.

**Results on Natural Non-IID Datasets.** While the negative results on artificial non-IID datasets are widely cited to claim FEDAVG suffers due to data heterogeneity, we doubt whether the results are representative and general enough to cover all practical applications. Is it possible that there are some scenarios where the data heterogeneity has benign effects and has different characteristics from the heterogeneity simulated through these originally IID datasets?

We conduct some experiments on StackOverflow, a naturally non-IID split dataset for next-word prediction. Each client in the dataset corresponds to a unique user on the Stack Overflow site. The data of each client consists of questions and answers written by the corresponding user. More details about the experimental setup, such as optimizer and learning rate choices, can be found in the Appendix. From the naturally heterogeneous StackOverflow dataset, we create its IID version by aggregating and shuffling the data from all clients, and then re-assigning the IID data back to clients. Surprisingly, the results in Figure 1 show that the convergence of FEDAVG is nearly identical (with limited communication rounds) on the new IID dataset and its original non-IID version. This observation contradicts the conventional wisdom that data heterogeneity always adds an additional error to the optimization procedure, as well as the empirical results on *artificial non-IID* datasets. There is indeed a gap between the theory and practice of FEDAVG.

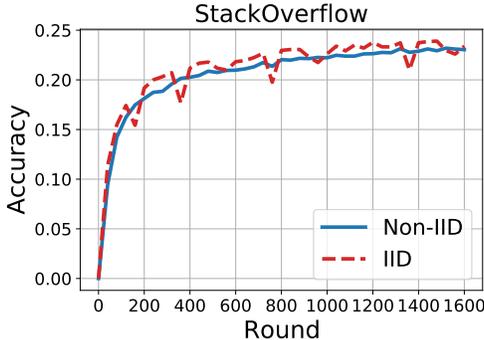


Figure 1: **Results on the naturally non-IID dataset StackOverflow.** The IID version of StackOverflow is created via aggregating and shuffling all the data from clients. We observe that FEDAVG achieves roughly the same convergence in both IID and non-IID settings.

Experimental results from some previous papers serve as additional evidence to support our claim. For instance, FEDAVG converges significantly faster than FEDSGD on Shakespeare and StackOverflow datasets

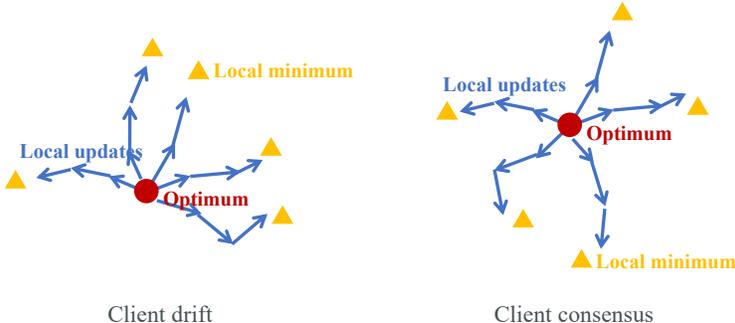


Figure 2: Comparison between client drift and client consensus. (Left) Previous works focus on the worst case where the average of clients’ local updates (or drifts) drives the global model away from the optimum; (Right) Our client consensus hypothesis states that the average local drifts at the optimum is very small or close to zero. This structured heterogeneity does not impact the convergence of FEDAVG.

even with strong heterogeneity (McMahan et al., 2017; Charles et al., 2021). Even though a more sophisticated method like SCAFFOLD does better on artificial datasets like EMNIST, it does not outperform FEDAVG on many naturally heterogeneous datasets like FEMNIST, StackOverflow, etc., as shown in (Reddi et al., 2021; Wu et al., 2023; Li et al., 2021). Explanations to these observations still remain mysteries.

## 4 Client Consensus Hypothesis

In previous sections, we saw that there is a gap between the theory and practice of FEDAVG. While theory predicts that FEDAVG performs worse in the presence of data heterogeneity, we did not observe this on naturally non-IID datasets. Also, many advanced methods designed to tackle data heterogeneity problem do not outperform FEDAVG. In order to explain the above phenomenon, in this section, we propose the client consensus hypothesis. Many realistic federated datasets may satisfy this special property such that data heterogeneity has very limited impacts on the convergence of FEDAVG.

### 4.1 Key Insights

FEDAVG is commonly viewed as a perturbed version of vanilla SGD. To see this, we can consider the accumulated local updates on each client as a pseudo-gradient, which is an approximation of the batch client gradient. When clients have heterogeneous data, the average of the pseudo-gradients across all clients can be very different from the original global gradient. Especially, when the global model approaches the optimum (or a stationary point) with a constant learning rate, since the average of pseudo-gradients is not zero, the global model cannot stay and may drift away towards a different point. This is referred to as client drift problem in literature (Karimireddy et al., 2020b). An illustration is provided in Figure 2.

The above insight is true in the worst case. But it is possible that real-world FL datasets are far away from this worst case. In this paper, we make a hypothesis about a special class of data heterogeneity. In particular, clients can still have drastically different data distributions and local minima. But, at the global optimum (or stationary point), clients’ local updates (i.e., pseudo-gradients) cancel out with each other and hence, the average pseudo-gradient at the optimum is or close to zero. As a consequence, when the global model reach the global optimum, it can stay there and will not drift away. That is, clients reach some kind of consensus at the optimum. For this reason, we name the hypothesis as *client consensus hypothesis*. We will show it in Section 4.2 that this special class of non-IIDness has no negative impacts on the convergence of FEDAVG. Therefore, the hypothesis can serve as a possible explanation of the effectiveness of FEDAVG.

More formally, we define the deterministic pseudo-gradient for client  $c$  as follows:

$$\mathcal{G}_c(\mathbf{w}) \triangleq \frac{1}{\eta H}(\mathbf{w} - \mathbf{w}_{c,\text{GD}}^{(H)}) \tag{6}$$

where  $\mathbf{w}_{c,\text{GD}}^{(H)}$  denotes the locally trained model at client  $c$  after performing  $H$  steps of GD from  $\mathbf{w}$  using a local learning rate  $\eta$ . Given this definition, the client consensus hypothesis is formally stated below.

**Assumption 4 (Client Consensus Hypothesis).** *On real-world federated datasets, for the values of  $\eta, H$  used in FEDAVG, the average pseudo-gradient at the optimum (i.e., average client model drift at the optimum)*

$$\rho \triangleq \left\| \mathbb{E}_c[\mathcal{G}_c(\mathbf{w}^*)] \right\|$$

*is very small or close to zero, where  $\mathbf{w}^*$  is the global optimum or stationary point of the global objective (1).*

Note that the Client Consensus Hypothesis is satisfied whenever FEDAVG converges with a constant learning rate. Thus, in a sense, it is the minimal assumption under which one can expect to prove improved convergence rates for FEDAVG.  $\rho$  can be interpreted as the *average drift at optimum*; this quantity is akin to a *heterogeneity metric* and our theoretical results will be in terms of it.

**Connections to Existing Analysis with Bounded Gradient Dissimilarity.** Here, we discuss the connections between client consensus hypothesis and previous convergence analysis with bounded gradient dissimilarity assumption. As we mentioned before, FEDAVG can be viewed as a perturbed version of vanilla SGD. This alternative view is also critical in the convergence analysis. A key step in previous FEDAVG analysis is to bound the perturbations, i.e., the difference between the average pseudo-gradient  $\mathbb{E}_c[\mathcal{G}_c]$  and the batch global gradient  $\nabla F(\mathbf{w}) = \mathbb{E}_c[\nabla F_c(\mathbf{w})]$ . Previous works always upper bound this difference using Jensen’s inequality as follows:

$$\left\| \mathbb{E}_c[\mathcal{G}_c(\mathbf{w}) - \nabla F_c(\mathbf{w})] \right\|^2 \leq \mathbb{E}_c \left\| \mathcal{G}_c(\mathbf{w}) - \nabla F_c(\mathbf{w}) \right\|^2. \quad (7)$$

Then, the right-hand-side (RHS) can be further bounded using the gradient dissimilarity assumption. However, we note that the above upper bound omits the correlations among different clients. While the RHS (i.e., average of squared  $\ell_2$  differences) can be large or unbounded, the LHS (i.e., squared  $\ell_2$  norm of the average difference) can be small or just zero, especially at the optimum  $\mathbf{w}^*$ . Therefore, using (7) in the analysis may result in a pessimistic estimate.

Instead, in this paper, the client consensus hypothesis states that the LHS of (7) at the optimum is close to, if not, zero. Using the hypothesis and our new analysis, we no longer need to provide an upper bound of the RHS of (7) for all points. As a consequence, the new analysis does not require the gradient dissimilarity to be bounded.

## 4.2 Convergence of FedAvg under Client Consensus Hypothesis

In this subsection, we will provide a convergence analysis for FEDAVG under the client consensus hypothesis (Assumption 4) for strongly-convex loss functions. In particular, we show that when  $\rho = 0$  in Assumption 4, *data heterogeneity has no negative impact* on the convergence of FEDAVG.

Besides the pseudo-gradient defined in (6), we need to define its stochastic version:

$$\hat{\mathcal{G}}_c(\mathbf{w}) \triangleq \frac{1}{\eta H} (\mathbf{w} - \mathbf{w}_c^{(H)}) \quad (8)$$

where  $\mathbf{w}_c^{(H)}$  denotes the locally trained model at client  $c$  using  $H$  steps of SGD with learning rate  $\eta$ , starting from  $\mathbf{w}$ . Let us define  $\mathcal{G} := \mathbb{E}_c[\mathcal{G}_c]$  and  $\hat{\mathcal{G}} := \mathbb{E}_c[\hat{\mathcal{G}}_c]$ . Then, we have the following result.

**Theorem 2.** *Under Assumptions 2 to 4, when each local objective function is  $\mu$ -strongly convex and the learning rates satisfy  $\alpha \leq 1/4, \eta \leq \min\{1/L, 1/\mu H\}$ , after  $T$  rounds of FEDAVG, we have*

$$\begin{aligned} \mathbb{E} \left\| \mathbf{w}^{(T)} - \mathbf{w}^* \right\|^2 &\leq \left(1 - \frac{1}{2}\alpha\eta H\mu\right)^T \left\| \mathbf{w}^{(0)} - \mathbf{w}^* \right\|^2 + \frac{2\alpha\eta H}{\mu} \max_{\mathbf{w}} \text{Var}[\hat{\mathcal{G}}(\mathbf{w})] \\ &\quad + \frac{20}{\mu^2} \max_{\mathbf{w}} \mathbb{E}_c \left\| \delta_c(\mathbf{w}) \right\|^2 + \frac{20\rho^2}{\mu^2} \end{aligned} \quad (9)$$

where  $\mathbb{E}[\cdot], \text{Var}[\cdot]$  are taken with respect to random noise in stochastic local updates, and  $\delta_c(\mathbf{w}) = (\mathbf{w}_{c,\text{GD}}^{(H)} - \mathbb{E}[\mathbf{w}_c^{(H)}]) / (\eta H)$  denotes the iterate bias between local GD and local SGD on client  $c$ .

From Theorem 2, we observe that the stochastic noise during local updates influences the second and the third terms on the right-hand-side of (9). The upper bounds of these two terms only depend on the dynamics of SGD, which has been well understood in literature. Specifically, in Khaled et al. (2020), the authors show that  $\mathbb{E}\|\mathbf{w}^{(H)} - \mathbb{E}[\mathbf{w}^{(H)}]\|^2 \leq 2\eta^2 H\sigma^2$ . As a consequence, we directly obtain that

$$\text{Var}[\hat{\mathcal{G}}(\mathbf{w})] = \frac{\mathbb{E}\|\mathbf{w}_c^{(H)} - \mathbb{E}[\mathbf{w}_c^{(H)}]\|^2}{\eta^2 H^2 M} \leq \frac{2\sigma^2}{MH}. \quad (10)$$

As for the iterate bias, one can obtain

$$\mathbb{E}_c \|\delta_c(\mathbf{w})\|^2 \leq \eta^2 L^2 \sigma^2 (H - 1), \quad (11)$$

the proof of which is provided in the Appendix. After substituting (10) and (11) into (9) and optimizing the learning rates, we can obtain the following convergence rate for FEDAVG.

**Corollary 1 (Convergence Rate for Strongly Convex Functions).** *In the same setting as Theorem 2, when  $\alpha = 1/4, \eta = \mathcal{O}(1/\mu HT)$ , the convergence rate of FEDAVG is<sup>1</sup>*

$$\mathbb{E}\|\mathbf{w}^{(T)} - \mathbf{w}^*\|^2 = \tilde{\mathcal{O}}\left(\frac{\sigma^2}{MHT} + \frac{\sigma^2}{HT^2} + \rho^2\right). \quad (12)$$

If clients perform local GD instead of local SGD, then when  $\eta = \min\{1/L, 1/(\mu H)\}$ , we have

$$\|\mathbf{w}^{(T)} - \mathbf{w}^*\|^2 = \mathcal{O}\left(\exp\left(-\frac{T}{16\kappa} \min\{\kappa, H\}\right) + \rho^2\right) \quad (13)$$

where  $\kappa = L/\mu$  denotes the condition number.

**Effects of Data Heterogeneity.** In the special regime of  $\rho \approx 0$ , i.e., when the client consensus hypothesis holds, Theorem 2 and Corollary 1 state that data heterogeneity does not have any negative impact on the convergence of FEDAVG. However, in previous works based on gradient dissimilarity, even if  $\rho = 0$ , there is an additional error caused by the positive gradient dissimilarity. Compared to the centralized setting where  $M = 1$ , (12) suggests that FEDAVG can provide linear speedup due to the usage of multiple workers.

**Extensions.** The above analysis can be extended to various settings. Below we provide several examples.

(1) *Client Sampling:* If we consider client sampling in FEDAVG, then only the variance term in (9) will change while the other terms will be unaffected. One can obtain new convergence guarantees by analyzing the variance of different sampling schemes and then simply substituting them into (9). Standard techniques to analyze client sampling (Yang et al., 2020) can be directly applied.

(2) *Third-order Smoothness:* When the local objective functions satisfy third-order smoothness ( $\|\nabla^2 F_c(\mathbf{w}) - \nabla^2 F_c(\mathbf{u})\| \leq Q \|\mathbf{w} - \mathbf{u}\|$ ), the bound of the iterate bias  $\delta(\mathbf{w})$  can be further improved while all other terms remain unchanged. According to Glasgow et al. (2021), one can obtain  $\|\delta(\mathbf{w})\| \leq \frac{1}{4}\eta^2 H Q \sigma^2$ . In the special case when local objective functions are quadratic, we have  $Q = 0$ . That is, there is no iterate bias. As a consequence, the convergence rate of FEDAVG can be significantly improved.

Moreover, a result for the general convex case is provided in the Appendix. Analysis for non-convex functions requires a different framework and is non-trivial. So we leave it for future work.

**Comparison with Previous Works.** In Table 2, we summarize the convergence rates of FEDAVG in different papers. All previous results depend on the gradient dissimilarity upper bound  $\zeta$ , which is generally large for heterogeneous data settings. However, in our result, under client consensus hypothesis, we show that the effects of data heterogeneity can be measured by average drift at optimum  $\rho$ , which can be close to zero even in presence of strong data heterogeneity, as we showed in the quadratic example and experiments on FEMNIST and StackOverflow. When the average drift at optimum is zero, we can get an improved convergence rate compared to existing results.

<sup>1</sup> $\tilde{\mathcal{O}}(\cdot)$  hides log factors.

Table 2: Comparison with existing results for strongly convex objectives functions with deterministic local updates. In the table, the error is measured by the distance to the optimum  $\|\mathbf{w} - \mathbf{w}^*\|^2$ , and  $\kappa = L/\mu$  is the condition number. Also, we omit logarithmic factors. Compared to previous results, we show that in the considered setting: (i) FEDAVG enjoys linear convergence to the global optimum, and (ii) the multiple local steps of FEDAVG mitigate the impact of ill-conditioning (high condition number).

Algorithm	Worst-case error	Comm. rounds to attain $\epsilon$ error when $\rho = 0$
GD	$\exp(-T/\kappa)$	$\mathcal{O}(\kappa \log(1/\epsilon))$
FEDAVG (Koloskova et al., 2020)	$\zeta^2/T^2$	$\mathcal{O}(1/\epsilon^2)$
FEDAVG (Woodworth et al., 2020a)	$1/(HT + H^2T^2) + \zeta^2/T^2$	$\mathcal{O}(1/\epsilon^2)$
FEDAVG (Ours)	$\exp(-T \min\{1, H/\kappa\}) + \rho^2$	$\mathcal{O}(\max\{1, \kappa/H\} \log(1/\epsilon))$

## 5 Validating the Client Consensus Hypothesis

Here we will present some evidence to show that the client consensus hypothesis is realistic and practical; we do so (a) theoretically, in a linear regression setting, and (b) empirically, on some naturally split datasets such as FEMNIST and StackOverflow.

### 5.1 Linear Regression Example

**Intuition on When Client Consensus Hypothesis Holds.** One of the key insights of the client consensus hypothesis is that clients do reach some consensus and a single global model can work reasonably well for all clients though they have heterogeneous data. Inspired by this, we assume that all clients have the same conditional probability of the label given data, i.e.,  $p_c(y|\mathbf{x}) = p(y|\mathbf{x}), \forall c$ , where  $\mathbf{x}, y$  denote the input data and its label, respectively. In this case, clients still have heterogeneous data distributions, as they may have drastically different  $p_c(\mathbf{x})$  and  $p_c(\mathbf{x}, y)$ . However, the client’s updates should not conflict with each other, as the learning algorithms tend to learn the same  $p(y|\mathbf{x})$  on all clients.

Let us study a concrete *linear regression* setting satisfying the above property. Specifically, we assume that the label of the  $i^{\text{th}}$  data sample on client  $c$  is generated as follows:

$$y_{c,i} = \langle \mathbf{w}^*, \mathbf{x}_{c,i} \rangle + \epsilon_{c,i}, \quad (14)$$

where  $\mathbf{w}^* \in \mathbb{R}^d$  denotes the optimal model, and  $\epsilon_{c,i} \sim \mathcal{P}_\epsilon$  is a zero-mean random noise and independent from  $\mathbf{x}_{c,i}$  (this is a common assumption in statistical learning). We also assume that all  $\|\mathbf{x}_{c,i}\|$  have bounded variance. Note that both  $\mathbf{w}^*$  and  $\mathcal{P}_\epsilon$  are the same across all the clients, i.e., all clients have the same label generation process and hence, the same conditional probability  $p(y|\mathbf{x})$ . Our goal is to find the optimal model  $\mathbf{w}^*$  given a large amount of clients with *finite* data samples (which is a common cross-device FL setting (Kairouz et al., 2019)). We use the squared loss function which makes our objective function *quadratic*; specifically, it is:

$$\begin{aligned} F_c(\mathbf{w}) &= \frac{1}{n} \sum_{i=1}^n \frac{1}{2} (y_{c,i} - \mathbf{w}^\top \mathbf{x}_{c,i})^2 \\ &= \frac{1}{2} (\mathbf{w} - \mathbf{w}^*)^\top \mathbf{A}_c (\mathbf{w} - \mathbf{w}^*) - \mathbf{b}_c^\top (\mathbf{w} - \mathbf{w}^*) + \text{const.}, \end{aligned} \quad (15)$$

where  $\mathbf{A}_c = \sum_{i=1}^n \mathbf{x}_{c,i} \mathbf{x}_{c,i}^\top / n$ ,  $\mathbf{b}_c = \sum_{i=1}^n \epsilon_{c,i} \mathbf{x}_{c,i} / n$ . The minimizer of local objective  $F_c(\mathbf{w})$  is  $\mathbf{w}_c^* = \mathbf{w}^* + \mathbf{A}_c^{-1} \mathbf{b}_c$ , which is different from the global minimizer  $\mathbf{w}^*$  as  $\mathbf{b}_c \neq 0$ .

**Client Consensus Hypothesis.** In this problem, we can show that client consensus hypothesis holds. In particular, one can show that the average pseudo-gradient at the optimum  $\rho$  is

$$\mathbb{E}_c[\mathcal{G}_c(\mathbf{w}^*)] = \mathbb{E}_c \left[ \frac{1}{H} \sum_{h=0}^{H-1} [\mathbf{I} - (\mathbf{I} - \eta \nabla^2 \mathbf{A}_c)^h] \mathbf{b}_c \right]. \quad (16)$$

Due to the independence of  $\epsilon_{c,i}$  and  $\mathbf{x}_{c,i}$ , for any choice of  $H$ ,  $\rho = \|\mathbb{E}_c[\mathcal{G}_c(\mathbf{w}^*)]\| \rightarrow 0$  almost surely when the number of clients  $M \rightarrow \infty$ .

**Unbounded Gradient Dissimilarity.** In contrast, if we check the gradient dissimilarity, we have:

$$\mathbb{E}_c \|\nabla F_c(\mathbf{w}^*) - \nabla F(\mathbf{w}^*)\|^2 = \mathbb{E}_c \|\mathbf{b}_c\|^2. \quad (17)$$

Observe that  $\epsilon$  can have extremely large variance such that the gradient dissimilarity bound  $\zeta$  is arbitrarily large. As a result, existing analyses, which rely on the bounded gradient dissimilarity, may predict that FEDAVG is much worse than its non-local counterparts. However, by simple manipulations on the update rule of FEDAVG, one can prove the following theorem.

**Theorem 3.** *Suppose that the weighting of the clients is uniform, and each client has a small finite amount of data. For the above linear regression setting, the iterates of Local GD (i.e., deterministic version of FEDAVG) satisfy the following equation almost surely as the number of clients goes to infinity:*

$$\mathbf{w}^{(T)} - \mathbf{w}^* = \left[ \mathbb{E}_c \left[ (\mathbf{I} - \eta \mathbf{A}_c)^H \right] \right]^T (\mathbf{w}^{(0)} - \mathbf{w}^*). \quad (18)$$

The proof is relegated to the Appendix. From Theorem 3, it is clear that if the learning rate  $\eta$  is properly set such that  $(\mathbf{I} - \eta \mathbf{A}_c)$  is positive definite, then performing more local updates (larger  $H$ ) will lead to faster linear convergence rate  $\mathcal{O}(\exp(-T))$  to the global optimum  $\mathbf{w}^*$ . That is, local GD is strictly better than vanilla GD. However, previous papers based on gradient dissimilarity will get a substantially slower rate of  $\mathcal{O}(1/T^2)$ . In this example, while gradient dissimilarity can be arbitrarily large, data heterogeneity does not have any negative impacts on FEDAVG as the client consensus hypothesis holds.

## 5.2 Experiments on Naturally Non-IID Datasets

In this subsection, we empirically show that the proposed client consensus hypothesis approximately holds across multiple practical training tasks.

In Figure 3, we first run mini-batch SGD on Federated EMNIST (FEMNIST) (McMahan et al., 2017) and StackOverflow Next Word Prediction datasets (Reddi et al., 2019) to obtain an approximation for the optimal model  $\mathbf{w}^*$ . Then we evaluate the average drift at optimum  $\rho = \|\mathbb{E}_c B_c(\mathbf{w}^*)\|$  and its upper bound as given in (7) on these datasets. We summarize the important observations below.

- As shown in Figures 3a and 3b, the average drift (or pseudo-gradient) at the optimum, i.e.  $\rho$ , is very close to zero on naturally non-IID datasets FEMNIST and StackOverflow. The proposed client consensus hypothesis is true on these two realistic datasets.
- From Figures 3a and 3b, we also observe that there is a large gap between  $\rho$  and its upper bound as given in (7). This suggests that previous analyses using this upper bound can be loose.
- We run the same set of experiments on a non-IID CIFAR-100 dataset. Figure 3c shows that on this artificial dataset, the client consensus hypothesis no longer holds. The average drift at optimum  $\rho$  is pretty close to its upper bound and far from zero. This is the scenario where FEDAVG fails with heterogeneous data and faces the client drift problem.

The drastically different observations on FEMNIST, StackOverflow and artificial CIFAR-100 highlight that there are various kinds of data heterogeneity. For the heterogeneity satisfying client consensus hypothesis, it may have very limited impacts on the convergence of FEDAVG.

Furthermore, we run FEDAVG on FEMNIST dataset following the same setup as (Reddi et al., 2021) and check the difference between the average of pseudo-gradient and the true batch gradient at several intermediate points on the optimization trajectory. For each point, we let clients perform local GD with the same local learning rate for multiple steps. As shown in Figure 4, we observe a significant gap between the norm of average differences (red lines) and its upper bound (7) (i.e., the average of  $\ell_2$  difference, blue lines). Especially, at the 50<sup>th</sup> and 100<sup>th</sup> rounds, the upper bound is about 10 times larger. These observations suggest that the upper bounds based on the gradient dissimilarity are loose in practice.

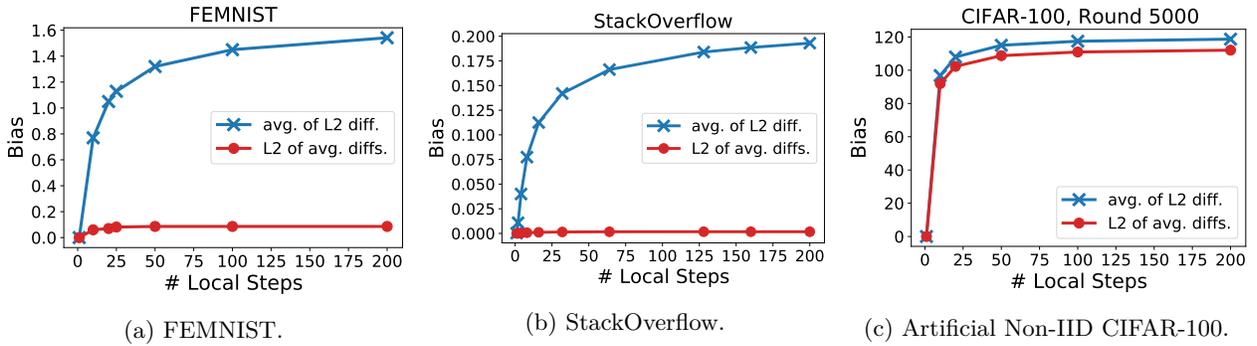


Figure 3: Difference between the average pseudo-gradient and the global gradient at the optimal point  $w^*$  on three different datasets. We observe that the norm of the average gradient differences at  $w^*$  (red line) nearly remain zero on all natural non-IID datasets but its upper bound used in previous analyses (7) (blue lines) slowly become larger when  $H$  increases. On a artificial non-IID CIFAR-100 dataset, these two values are pretty close to each other.

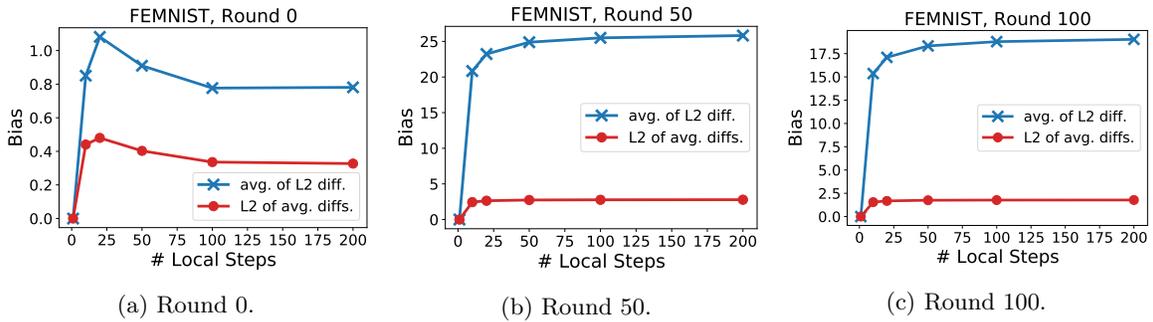


Figure 4: Difference between the average pseudo-gradient and the global gradient on a FEDAVG’s optimization trajectory. There is a significant gap between the average gradient differences (red line) and its upper bound (blue lines). Both of them increase and then saturate when increasing  $H$ .

## 6 Conclusions

In this paper, we aim to bridge the gap between theory and practice of the popular FEDAVG algorithm. We found that previous analyses based on the bounded gradient dissimilarity assumption can be too pessimistic for practical applications. On some natural federated datasets, FEDAVG may have identical performance on both IID and non-IID settings. In order to explain this phenomenon, we introduced the client consensus hypothesis and formally proved that under this hypothesis, data heterogeneity does not exacerbate the convergence of FEDAVG. More importantly, we show that this hypothesis holds for a linear regression problem and many practical federated datasets, including FEMNIST and StackOverflow. So the proposed hypothesis is realistic and can better explain the empirical success of FEDAVG.

Given the above contributions, several future directions are ripe for exploration. Our client consensus hypothesis is expressed in terms of a quantity ( $\rho$ ) that is akin to a heterogeneity metric; it would be interesting to come up with a formal metric that applies to all settings and can replace the existing loose gradient dissimilarity metric. Perhaps such metrics could also be helpful in guiding the design principles for FL algorithms. Besides, it may be also worthwhile to design a more practical criterion to check whether a dataset satisfies the client consensus hypothesis.

## References

- Mahmoud Assran, Nicolas Loizou, Nicolas Ballas, and Michael Rabbat. Stochastic gradient push for distributed deep learning. *arXiv preprint arXiv:1811.10792*, 2018.
- Zachary Charles and Jakub Konečný. Convergence and accuracy trade-offs in federated learning and meta-learning. In *International Conference on Artificial Intelligence and Statistics*, pp. 2575–2583. PMLR, 2021.
- Zachary Charles and Keith Rush. Iterated vector fields and conservatism, with applications to federated learning. In *International Conference on Algorithmic Learning Theory*, pp. 130–147. PMLR, 2022.
- Zachary Charles, Zachary Garrett, Zhouyuan Huo, Sergei Shmulyian, and Virginia Smith. On large-cohort training for federated learning. *Advances in Neural Information Processing Systems*, 34, 2021.
- Michael Crawshaw, Yajie Bao, and Mingrui Liu. Federated learning with client subsampling, data heterogeneity, and unbounded smoothness: A new algorithm and lower bounds. *Advances in Neural Information Processing Systems*, 36, 2024.
- Rudrajit Das, Anish Acharya, Abolfazl Hashemi, Sujay Sanghavi, Inderjit S Dhillon, and Ufuk Topcu. Faster non-convex federated learning via global and local momentum. In *Uncertainty in Artificial Intelligence*, pp. 496–506. PMLR, 2022.
- Margalit Glasgow, Honglin Yuan, and Tengyu Ma. Sharp bounds for federated averaging (local sgd) and continuous perspective. *arXiv preprint arXiv:2111.03741*, 2021.
- Xinran Gu, Kaifeng Lyu, Longbo Huang, and Sanjeev Arora. Why (and when) does local SGD generalize better than SGD? In *The Eleventh International Conference on Learning Representations*, 2023. URL <https://openreview.net/forum?id=svCcu6Dr1>.
- Farzin Haddadpour and Mehrdad Mahdavi. On the convergence of local descent methods in federated learning. *arXiv preprint arXiv:1910.14425*, 2019.
- Tzu-Ming Harry Hsu, Hang Qi, and Matthew Brown. Measuring the effects of non-identical data distribution for federated visual classification. *arXiv preprint arXiv:1909.06335*, 2019.
- Divyansh Jhunjhunwala, PRANAY SHARMA, Aushim Nagarkatti, and Gauri Joshi. Fedvarp: Tackling the variance due to partial client participation in federated learning. In *The 38th Conference on Uncertainty in Artificial Intelligence*, 2022.
- Peter Kairouz, H Brendan McMahan, Brendan Avent, Aurélien Bellet, Mehdi Bennis, Arjun Nitin Bhagoji, Keith Bonawitz, Zachary Charles, Graham Cormode, Rachel Cummings, et al. Advances and open problems in federated learning. *arXiv preprint arXiv:1912.04977*, 2019.
- Sai Praneeth Karimireddy, Martin Jaggi, Satyen Kale, Mehryar Mohri, Sashank J Reddi, Sebastian U Stich, and Ananda Theertha Suresh. Mime: Mimicking centralized stochastic algorithms in federated learning. *arXiv preprint arXiv:2008.03606*, 2020a.
- Sai Praneeth Karimireddy, Satyen Kale, Mehryar Mohri, Sashank J Reddi, Sebastian U Stich, and Ananda Theertha Suresh. SCAFFOLD: Stochastic controlled averaging for on-device federated learning. In *Proceedings of the International Conference on Machine Learning*, 2020b.
- A Khaled, K Mishchenko, and P Richtárik. Tighter theory for local SGD on identical and heterogeneous data. In *Proceedings of the Twenty Third International Conference on Artificial Intelligence and Statistics*, 2020.
- Anastasia Koloskova, Nicolas Loizou, Sadra Boreiri, Martin Jaggi, and Sebastian U Stich. A unified theory of decentralized SGD with changing topology and local updates. In *International Conference on Machine Learning*, 2020.
- Qinbin Li, Bingsheng He, and Dawn Song. Model-contrastive federated learning. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 10713–10722, 2021.

- Tian Li, Anit Kumar Sahu, Manzil Zaheer, Maziar Sanjabi, Ameet Talwalkar, and Virginia Smith. Federated optimization in heterogeneous networks. In *Proceedings of the Conference on Machine Learning and Systems*, 2020a.
- 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.
- Xiang Li, Kaixuan Huang, Wenhao Yang, Shusen Wang, and Zhihua Zhang. On the convergence of FedAvg on non-IID data. In *International Conference on Learning Representations*, 2020b. URL <https://openreview.net/forum?id=HJxNAnVtDS>.
- Xiangru Lian, Ce Zhang, Huan Zhang, Cho-Jui Hsieh, Wei Zhang, and Ji Liu. Can decentralized algorithms outperform centralized algorithms? a case study for decentralized parallel stochastic gradient descent. In *Advances in Neural Information Processing Systems*, pp. 5336–5346, 2017.
- Grigory Malinovskiy, Dmitry Kovalev, Elnur Gasanov, Laurent Condat, and Peter Richtarik. From local sgd to local fixed-point methods for federated learning. In *International Conference on Machine Learning*, pp. 6692–6701. PMLR, 2020.
- H Brendan McMahan, Eider Moore, Daniel Ramage, Seth Hampson, et al. Communication-efficient learning of deep networks from decentralized data. In *Proceedings of the 20th International Conference on Artificial Intelligence and Statistics*, 2017.
- Sashank Reddi, Zachary Charles, Manzil Zaheer, Zachary Garrett, Keith Rush, Jakub Konečný, Sanjiv Kumar, and H Brendan McMahan. Adaptive federated optimization. *arXiv preprint arXiv:2003.00295*, 2020.
- Sashank J Reddi, Satyen Kale, and Sanjiv Kumar. On the convergence of adam and beyond. *arXiv preprint arXiv:1904.09237*, 2019.
- Sashank J. Reddi, Zachary Charles, Manzil Zaheer, Zachary Garrett, Keith Rush, Jakub Konečný, Sanjiv Kumar, and Hugh Brendan McMahan. Adaptive federated optimization. In *International Conference on Learning Representations*, 2021. URL <https://openreview.net/forum?id=LkFG31B13U5>.
- Sebastian U Stich. Local SGD converges fast and communicates little. In *Proceedings of the International Conference on Learning Representations (ICLR)*, 2019.
- Jianyu Wang and Gauri Joshi. Cooperative SGD: A unified framework for the design and analysis of communication-efficient SGD algorithms. *arXiv preprint arXiv:1808.07576*, 2018.
- Jianyu Wang, Qinghua Liu, Hao Liang, Gauri Joshi, and H Vincent Poor. Tackling the objective inconsistency problem in heterogeneous federated optimization. *Advances in Neural Information Processing Systems*, 33, 2020a.
- Jianyu Wang, Vinayak Tantia, Nicolas Ballas, and Michael Rabbat. SlowMo: Improving communication-efficient distributed SGD with slow momentum. In *Proceedings of the International Conference on Learning Representations (ICLR)*, 2020b.
- Jianyu Wang, Zachary Charles, Zheng Xu, Gauri Joshi, H Brendan McMahan, Maruan Al-Shedivat, Galen Andrew, Salman Avestimehr, Katharine Daly, Deepesh Data, et al. A field guide to federated optimization. *arXiv preprint arXiv:2107.06917*, 2021.
- Jianyu Wang, Anit Kumar Sahu, Gauri Joshi, and Soumya Kar. Matcha: A matching-based link scheduling strategy to speed up distributed optimization. *IEEE Transactions on Signal Processing*, 70:5208–5221, 2022.
- Jiayi Wang, Shiqiang Wang, Rong-Rong Chen, and Mingyue Ji. A new theoretical perspective on data heterogeneity in federated optimization. In *Federated Learning and Analytics in Practice: Algorithms, Systems, Applications, and Opportunities*, 2023a.

- Jiayi Wang, Shiqiang Wang, Rong-Rong Chen, and Mingyue Ji. Rethinking the data heterogeneity in federated learning. In *2023 57th Asilomar Conference on Signals, Systems, and Computers*, pp. 624–628. IEEE, 2023b.
- Blake Woodworth, Kumar Kshitij Patel, and Nathan Srebro. Minibatch vs local sgd for heterogeneous distributed learning. *arXiv preprint arXiv:2006.04735*, 2020a.
- Blake Woodworth, Kumar Kshitij Patel, Sebastian U Stich, Zhen Dai, Brian Bullins, H Brendan McMahan, Ohad Shamir, and Nathan Srebro. Is local SGD better than minibatch SGD? In *International Conference on Machine Learning*, 2020b.
- Feijie Wu, Song Guo, Zhihao Qu, Shiqi He, Ziming Liu, and Jing Gao. Anchor sampling for federated learning with partial client participation. In *International Conference on Machine Learning*, pp. 37379–37416. PMLR, 2023.
- Haibo Yang, Minghong Fang, and Jia Liu. Achieving linear speedup with partial worker participation in non-iid federated learning. In *International Conference on Learning Representations*, 2020.
- 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*, 2019a.
- Hao Yu, Sen Yang, and Shenghuo Zhu. Parallel restarted SGD with faster convergence and less communication: Demystifying why model averaging works for deep learning. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 33, pp. 5693–5700, 2019b.
- Kun Yuan, Qing Ling, and Wotao Yin. On the convergence of decentralized gradient descent. *SIAM Journal on Optimization*, 26(3):1835–1854, 2016.
- Yue Zhao, Meng Li, Liangzhen Lai, Naveen Suda, Damon Civin, and Vikas Chandra. Federated learning with non-IID data. *arXiv preprint arXiv:1806.00582*, 2018.
- Fan Zhou and Guojing Cong. On the convergence properties of a k-step averaging stochastic gradient descent algorithm for nonconvex optimization. In *Proceedings of the 27th International Joint Conference on Artificial Intelligence (IJCAI)*, pp. 3219–3227, 2018.
- Ali Zindari, Ruichen Luo, and Sebastian U Stich. On the convergence of local sgd under third-order smoothness and hessian similarity. In *OPT 2023: Optimization for Machine Learning*, 2023.

## A Experimental Details

On FEMNIST, StackOverflow, and CIFAR-100 datasets, we strictly follow the training setup given in Reddi et al. (2020). Both models are neural networks and hence the objective functions are non-convex. In Figure 3, we used a trained model (obtained after training with mini-batch SGD) as an approximation of the global optimum. Then, a large set of clients are selected to perform local model training from the optimum to calculate the gradient bias. Details on the local training can be found in the following table.

Table 3: Details on local training in Figure 3.

Dataset	Model	Loss function	# of clients	Local optimizer	Local learning rate
FEMNIST	ConvNet	Cross-Entropy	500	GD	0.1
StackOverflow	LSTM	Cross-Entropy	1000	GD	0.5
CIFAR-100	ConvNet	Cross-Entropy	200	GD	0.5

## B Proof of Theorem 2

### B.1 Preliminaries

In this subsection, we will first introduce several useful lemmas, which relate to the properties of the deterministic pseudo-gradients. Before diving into the proof details, we would like to first introduce a lemma, which will be frequently used in the subsequent sections.

**Lemma 1 (Mean Value Theorem).** *Suppose function  $F$  is twice differentiable, then*

$$\nabla F_c(\mathbf{w}) - \nabla F_c(\mathbf{u}) = \mathbf{A}_c(\mathbf{w}, \mathbf{u}) \cdot (\mathbf{w} - \mathbf{u}) \quad (19)$$

where  $\mathbf{A}_c(\mathbf{w}, \mathbf{u}) = \int_0^1 \nabla^2 F_c(\mathbf{u} + s(\mathbf{w} - \mathbf{u})) ds$ . If  $\mu \preceq \nabla^2 F_c \preceq L$ , then it follows that  $\mu \preceq \mathbf{A}_c(\mathbf{w}, \mathbf{u}) \preceq L$  for any  $\mathbf{w}, \mathbf{u}$ .

**Lemma 2 (Convexity, Smoothness & Co-coercivity).** *When each local objective function  $F_c$  is  $L$ -Lipschitz smooth and  $\mu$ -strongly convex, for any  $\mathbf{w}, \mathbf{u} \in \mathbb{R}^d$ , we have*

$$\tilde{\mu} \|\mathbf{w} - \mathbf{u}\|^2 \leq \langle \mathcal{G}(\mathbf{w}) - \mathcal{G}(\mathbf{u}), \mathbf{w} - \mathbf{u} \rangle \leq \tilde{L} \|\mathbf{w} - \mathbf{u}\|^2, \quad (20)$$

$$\|\mathcal{G}(\mathbf{w}) - \mathcal{G}(\mathbf{u})\|^2 \leq \tilde{L} \langle \mathcal{G}(\mathbf{w}) - \mathcal{G}(\mathbf{u}), \mathbf{w} - \mathbf{u} \rangle, \quad (21)$$

where  $\tilde{\mu} = [1 - (1 - \eta\mu)^H]/(\eta H)$  and  $\tilde{L} = [1 + (1 - \eta\mu)^H]/(\eta H)$ .

*Proof.* Let us first focus on the pseudo-gradient on a specific client  $c$ . According to the definition of pseudo-gradient, we have

$$\eta H [\mathcal{G}_c(\mathbf{w}) - \mathcal{G}_c(\mathbf{u})] = \mathbf{w} - \mathbf{u} - (\mathbf{w}_c^{(H)} - \mathbf{u}_c^{(H)}) \quad (22)$$

$$= \mathbf{w} - \mathbf{u} - [\mathbf{w}_c^{(H-1)} - \mathbf{u}_c^{(H-1)} - \eta(\nabla F_c(\mathbf{w}_c^{(H-1)}) - \nabla F_c(\mathbf{u}_c^{(H-1)}))] \quad (23)$$

$$= \mathbf{w} - \mathbf{u} - (\mathbf{I} - \eta \mathbf{D}_c^{(H-1)})(\mathbf{w}_c^{(H-1)} - \mathbf{u}_c^{(H-1)}) \quad (24)$$

where (24) follows from Lemma 1 and  $\mathbf{D}_c$  is a symmetric matrix satisfying  $\mu \preceq \mathbf{D}_c \preceq L$ . Repeating the above procedure, we can obtain that

$$\eta H [\mathcal{G}_c(\mathbf{w}) - \mathcal{G}_c(\mathbf{u})] = \mathbf{w} - \mathbf{u} - \prod_{k=0}^{H-1} (\mathbf{I} - \eta \mathbf{D}_c^{(k)})(\mathbf{w} - \mathbf{u}) \quad (25)$$

$$= \left[ \mathbf{I} - \prod_{k=0}^{H-1} (\mathbf{I} - \eta \mathbf{D}_c^{(k)}) \right] (\mathbf{w} - \mathbf{u}). \quad (26)$$

As a consequence, we have

$$\eta H \langle \mathcal{G}_c(\mathbf{w}) - \mathcal{G}_c(\mathbf{u}), \mathbf{w} - \mathbf{u} \rangle = \|\mathbf{w} - \mathbf{u}\|^2 - \left\langle \prod_{k=0}^{H-1} (\mathbf{I} - \eta \mathbf{D}_c^{(k)}) (\mathbf{w} - \mathbf{u}), \mathbf{w} - \mathbf{u} \right\rangle. \quad (27)$$

Note that, due to Cauchy–Schwarz inequality,

$$\left| \left\langle \prod_{k=0}^{H-1} (\mathbf{I} - \eta \mathbf{D}_c^{(k)}) (\mathbf{w} - \mathbf{u}), \mathbf{w} - \mathbf{u} \right\rangle \right| \leq \prod_{k=0}^{H-1} \|\mathbf{I} - \eta \mathbf{D}_c^{(k)}\| \|\mathbf{w} - \mathbf{u}\|^2 \quad (28)$$

$$\leq (1 - \eta\mu)^H \|\mathbf{w} - \mathbf{u}\|^2. \quad (29)$$

That is,

$$-(1 - \eta\mu)^H \|\mathbf{w} - \mathbf{u}\|^2 \leq \left\langle \prod_{k=0}^{H-1} (\mathbf{I} - \eta \mathbf{D}_c^{(k)}) (\mathbf{w} - \mathbf{u}), \mathbf{w} - \mathbf{u} \right\rangle \leq (1 - \eta\mu)^H \|\mathbf{w} - \mathbf{u}\|^2. \quad (30)$$

It follows that,

$$\frac{1 - (1 - \eta\mu)^H}{\eta H} \|\mathbf{w} - \mathbf{u}\|^2 \leq \langle \mathcal{G}_c(\mathbf{w}) - \mathcal{G}_c(\mathbf{u}), \mathbf{w} - \mathbf{u} \rangle \leq \frac{1 + (1 - \eta\mu)^H}{\eta H} \|\mathbf{w} - \mathbf{u}\|^2. \quad (31)$$

Taking the average over all clients, we complete the proof of (20).

Next, we are going to prove (21). Note that

$$\|\mathbf{w}_c^{(H)} - \mathbf{u}_c^{(H)}\| = \left\| \prod_{k=0}^{H-1} (\mathbf{I} - \eta \mathbf{D}_c^{(k)}) (\mathbf{w} - \mathbf{u}) \right\| \quad (32)$$

$$\leq \prod_{k=0}^{H-1} \|\mathbf{I} - \eta \mathbf{D}_c^{(k)}\| \|\mathbf{w} - \mathbf{u}\| \quad (33)$$

$$\leq (1 - \eta\mu)^H \|\mathbf{w} - \mathbf{u}\|. \quad (34)$$

As a result, we have

$$\|\mathbf{w}^{(H)} - \mathbf{u}^{(H)}\|^2 = \|\mathbb{E}_c \mathbf{w}_c^{(H)} - \mathbb{E}_c \mathbf{u}_c^{(H)}\|^2 \leq \mathbb{E}_c \|\mathbf{w}_c^{(H)} - \mathbf{u}_c^{(H)}\|^2 \quad (35)$$

$$\leq [(1 - \eta\mu)^H]^2 \|\mathbf{w} - \mathbf{u}\|^2. \quad (36)$$

Then, according to the definition of pseudo-gradients, one can obtain

$$\eta^2 H^2 \|\mathcal{G}(\mathbf{w}) - \mathcal{G}(\mathbf{u})\|^2 = \|\mathbf{w} - \mathbf{u} - \mathbf{w}^{(H)} + \mathbf{u}^{(H)}\|^2 \quad (37)$$

$$= \|\mathbf{w} - \mathbf{u}\|^2 + \|\mathbf{w}^{(H)} - \mathbf{u}^{(H)}\|^2 - 2 \langle \mathbf{w} - \mathbf{u}, \mathbf{w}^{(H)} - \mathbf{u}^{(H)} \rangle \quad (38)$$

$$\leq [1 + (1 - \eta\mu)^H] \|\mathbf{w} - \mathbf{u}\|^2 - 2 \langle \mathbf{w} - \mathbf{u}, \mathbf{w}^{(H)} - \mathbf{u}^{(H)} \rangle - (1 - \eta\mu)^H [1 - (1 - \eta\mu)^H] \|\mathbf{w} - \mathbf{u}\|^2 \quad (39)$$

$$= [1 + (1 - \eta\mu)^H] \left[ \|\mathbf{w} - \mathbf{u}\|^2 - \langle \mathbf{w} - \mathbf{u}, \mathbf{w}^{(H)} - \mathbf{u}^{(H)} \rangle \right] - (1 - (1 - \eta\mu)^H) \langle \mathbf{w} - \mathbf{u}, \mathbf{w}^{(H)} - \mathbf{u}^{(H)} \rangle - (1 - \eta\mu)^H [1 - (1 - \eta\mu)^H] \|\mathbf{w} - \mathbf{u}\|^2 \quad (40)$$

$$= [1 + (1 - \eta\mu)^H] \left[ \|\mathbf{w} - \mathbf{u}\|^2 - \langle \mathbf{w} - \mathbf{u}, \mathbf{w}^{(H)} - \mathbf{u}^{(H)} \rangle \right] + (1 - (1 - \eta\mu)^H) \left[ \|\mathbf{w} - \mathbf{u}\|^2 - \langle \mathbf{w} - \mathbf{u}, \mathbf{w}^{(H)} - \mathbf{u}^{(H)} \rangle \right] - [1 + (1 - \eta\mu)^H] [1 - (1 - \eta\mu)^H] \|\mathbf{w} - \mathbf{u}\|^2 \quad (41)$$

Note that  $\eta H \langle \mathbf{w} - \mathbf{u}, \mathcal{G}(\mathbf{w}) - \mathcal{G}(\mathbf{u}) \rangle = \|\mathbf{w} - \mathbf{u}\|^2 - \langle \mathbf{w} - \mathbf{u}, \mathbf{w}^{(H)} - \mathbf{u}^{(H)} \rangle$  and  $\eta H \tilde{L} = 1 + (1 - \eta\mu)^H$ ,  $\eta H \tilde{\mu} = 1 - (1 - \eta\mu)^H$ , we have

$$\|\mathcal{G}(\mathbf{w}) - \mathcal{G}(\mathbf{u})\|^2 \leq (\tilde{L} + \tilde{\mu}) \langle \mathbf{w} - \mathbf{u}, \mathcal{G}(\mathbf{w}) - \mathcal{G}(\mathbf{u}) \rangle - \tilde{L}\tilde{\mu} \|\mathbf{w} - \mathbf{u}\|^2 \quad (42)$$

$$\begin{aligned} &= \tilde{L} \langle \mathbf{w} - \mathbf{u}, \mathcal{G}(\mathbf{w}) - \mathcal{G}(\mathbf{u}) \rangle \\ &\quad - \tilde{\mu} \left[ \tilde{L} \|\mathbf{w} - \mathbf{u}\|^2 - \langle \mathbf{w} - \mathbf{u}, \mathcal{G}(\mathbf{w}) - \mathcal{G}(\mathbf{u}) \rangle \right] \end{aligned} \quad (43)$$

$$\leq \tilde{L} \langle \mathbf{w} - \mathbf{u}, \mathcal{G}(\mathbf{w}) - \mathcal{G}(\mathbf{u}) \rangle \quad (44)$$

where the last inequality is due to the smoothness of the pseudo-gradient (20).  $\square$

## B.2 Main Proof

In the analysis below, we first analyze the training progress within one round. Suppose the current global model is  $\mathbf{w}$  and the next round's global model is  $\mathbf{w}^+$ . Without otherwise stated, the expectation  $\mathbb{E}$  and variance  $\text{Var}$  are conditioned on the current global model  $\mathbf{w}$ . For the ease of writing, we define effective learning rate  $\tilde{\alpha} = \alpha\eta H$ .

First, according to the update rule of FEDAVG, we have

$$\mathbb{E} \|\mathbf{w}^+ - \mathbf{w}^*\|^2 = \mathbb{E} \left\| \mathbf{w} - \tilde{\alpha} \hat{\mathcal{G}}(\mathbf{w}) - \mathbf{w}^* \right\|^2 \quad (45)$$

$$= \left\| \mathbf{w} - \tilde{\alpha} \mathbb{E}[\hat{\mathcal{G}}(\mathbf{w})] - \mathbf{w}^* \right\|^2 + \tilde{\alpha}^2 \text{Var}[\hat{\mathcal{G}}(\mathbf{w})] \quad (46)$$

$$= \|\mathbf{w} - \mathbf{w}^*\|^2 + \tilde{\alpha}^2 \left\| \mathbb{E}[\hat{\mathcal{G}}(\mathbf{w})] \right\|^2 - 2\tilde{\alpha} \langle \mathbf{w} - \mathbf{w}^*, \mathbb{E}[\hat{\mathcal{G}}(\mathbf{w})] \rangle + \tilde{\alpha}^2 \text{Var}[\hat{\mathcal{G}}(\mathbf{w})] \quad (47)$$

Also, note that  $\mathbb{E}[\hat{\mathcal{G}}(\mathbf{w})] = \mathcal{G}(\mathbf{w}) - \mathcal{G}(\mathbf{w}^*) + \mathcal{G}(\mathbf{w}^*) + \delta(\mathbf{w})$ . So one can obtain

$$\begin{aligned} \mathbb{E} \|\mathbf{w}^+ - \mathbf{w}^*\|^2 &\leq \|\mathbf{w} - \mathbf{w}^*\|^2 + 2\tilde{\alpha}^2 \|\mathcal{G}(\mathbf{w}) - \mathcal{G}(\mathbf{w}^*)\|^2 - 2\tilde{\alpha} \langle \mathbf{w} - \mathbf{w}^*, \mathcal{G}(\mathbf{w}) - \mathcal{G}(\mathbf{w}^*) \rangle \\ &\quad + 2\tilde{\alpha}^2 \|\mathcal{G}(\mathbf{w}^*) + \delta(\mathbf{w})\|^2 - 2\tilde{\alpha} \langle \mathbf{w} - \mathbf{w}^*, \mathcal{G}(\mathbf{w}^*) + \delta(\mathbf{w}) \rangle + \tilde{\alpha}^2 \text{Var}[\hat{\mathcal{G}}(\mathbf{w})] \end{aligned} \quad (48)$$

$$\begin{aligned} &\leq (1 - \tilde{\alpha}\tilde{\mu}) \|\mathbf{w} - \mathbf{w}^*\|^2 + 2\tilde{\alpha}^2 \|\mathcal{G}(\mathbf{w}) - \mathcal{G}(\mathbf{w}^*)\|^2 - \tilde{\alpha} \langle \mathbf{w} - \mathbf{w}^*, \mathcal{G}(\mathbf{w}) - \mathcal{G}(\mathbf{w}^*) \rangle \\ &\quad + 2\tilde{\alpha}^2 \|\mathcal{G}(\mathbf{w}^*) + \delta(\mathbf{w})\|^2 - 2\tilde{\alpha} \langle \mathbf{w} - \mathbf{w}^*, \mathcal{G}(\mathbf{w}^*) + \delta(\mathbf{w}) \rangle + \tilde{\alpha}^2 \text{Var}[\hat{\mathcal{G}}(\mathbf{w})] \end{aligned} \quad (49)$$

where the first inequality uses the fact  $\|a + b\|^2 \leq 2\|a\|^2 + 2\|b\|^2$ , and the second inequality comes from the strongly-convexity of the pseudo-gradient. Now let us check the value of the following terms:

$$T_1 = 2\tilde{\alpha} \|\mathcal{G}(\mathbf{w}) - \mathcal{G}(\mathbf{w}^*)\|^2 - \langle \mathbf{w} - \mathbf{w}^*, \mathcal{G}(\mathbf{w}) - \mathcal{G}(\mathbf{w}^*) \rangle - 2 \langle \mathbf{w} - \mathbf{w}^*, \mathcal{G}(\mathbf{w}^*) + \delta(\mathbf{w}) \rangle. \quad (50)$$

According to the co-coercivity of the pseudo-gradient, we have

$$T_1 \leq \left[ 2\tilde{\alpha}\tilde{L} - 1 \right] \langle \mathbf{w} - \mathbf{w}^*, \mathcal{G}(\mathbf{w}) - \mathcal{G}(\mathbf{w}^*) \rangle - 2 \langle \mathbf{w} - \mathbf{w}^*, \mathcal{G}(\mathbf{w}^*) + \delta(\mathbf{w}) \rangle \quad (51)$$

$$\leq \left[ 2\tilde{\alpha}\tilde{L} - 1 \right] \langle \mathbf{w} - \mathbf{w}^*, \mathcal{G}(\mathbf{w}) - \mathcal{G}(\mathbf{w}^*) \rangle + \epsilon \|\mathbf{w} - \mathbf{w}^*\|^2 + \frac{1}{\epsilon} \|\mathcal{G}(\mathbf{w}^*) + \delta(\mathbf{w})\|^2 \quad (52)$$

where the last inequality comes from Young's inequality. When  $\tilde{\alpha}\tilde{\mu} \leq \tilde{\alpha}\tilde{L} \leq 1/4$ , we have

$$T_1 \leq -\frac{\tilde{\mu}}{2} \|\mathbf{w} - \mathbf{w}^*\|^2 + \epsilon \|\mathbf{w} - \mathbf{w}^*\|^2 + \frac{1}{\epsilon} \|\mathcal{G}(\mathbf{w}^*) + \delta(\mathbf{w})\|^2 \leq \frac{2 \|\mathcal{G}(\mathbf{w}^*) + \delta(\mathbf{w})\|^2}{\tilde{\mu}} \quad (53)$$

where the last inequality is obtained by setting  $\epsilon = \tilde{\mu}/2$ . Then, substituting (53) into (49) and noting that  $\tilde{\alpha}\tilde{\mu} \leq \tilde{\alpha}\tilde{L} \leq 1/4$  (that is,  $\tilde{\alpha} \leq 1/(4\tilde{\mu})$ ),

$$\mathbb{E} \|\mathbf{w}^+ - \mathbf{w}^*\|^2 \leq (1 - \tilde{\alpha}\tilde{\mu}) \|\mathbf{w} - \mathbf{w}^*\|^2 + \left(2\tilde{\alpha}^2 + \frac{2\tilde{\alpha}}{\tilde{\mu}}\right) \|\mathcal{G}(\mathbf{w}^*) + \delta(\mathbf{w})\|^2 + \tilde{\alpha}^2 \text{Var}[\hat{\mathcal{G}}(\mathbf{w})] \quad (54)$$

$$\leq (1 - \tilde{\alpha}\tilde{\mu}) \|\mathbf{w} - \mathbf{w}^*\|^2 + \frac{5\tilde{\alpha}}{2\tilde{\mu}} \|\mathcal{G}(\mathbf{w}^*) + \delta(\mathbf{w})\|^2 + \tilde{\alpha}^2 \text{Var}[\hat{\mathcal{G}}(\mathbf{w})] \quad (55)$$

$$\leq (1 - \tilde{\alpha}\tilde{\mu}) \|\mathbf{w} - \mathbf{w}^*\|^2 + \tilde{\alpha}^2 \text{Var}[\hat{\mathcal{G}}(\mathbf{w})] + \frac{5\tilde{\alpha}}{\tilde{\mu}} \|\delta(\mathbf{w})\|^2 + \frac{5\tilde{\alpha}}{\tilde{\mu}} \|\mathcal{G}(\mathbf{w}^*)\|^2 \quad (56)$$

$$\leq (1 - \tilde{\alpha}\tilde{\mu}) \|\mathbf{w} - \mathbf{w}^*\|^2 + \tilde{\alpha}^2 \max_{\mathbf{w}} \text{Var}[\hat{\mathcal{G}}(\mathbf{w})] + \frac{5\tilde{\alpha}}{\tilde{\mu}} \max_{\mathbf{w}} \|\delta(\mathbf{w})\|^2 + \frac{5\tilde{\alpha}}{\tilde{\mu}} \|\mathcal{G}(\mathbf{w}^*)\|^2 \quad (57)$$

After total  $T$  communication rounds and taking the total expectation, we end up with

$$\begin{aligned} \mathbb{E} \|\mathbf{w}^{(T)} - \mathbf{w}^*\|^2 &\leq (1 - \tilde{\alpha}\tilde{\mu})^T \|\mathbf{w}^{(0)} - \mathbf{w}^*\|^2 + \frac{\tilde{\alpha}}{\tilde{\mu}} \max_{\mathbf{w}} \text{Var}[\hat{\mathcal{G}}(\mathbf{w})] \\ &\quad + \frac{5}{\tilde{\mu}^2} \max_{\mathbf{w}} \mathbb{E}_c \|\delta_c(\mathbf{w})\|^2 + \frac{5\rho^2}{\tilde{\mu}^2}. \end{aligned} \quad (58)$$

When  $\eta H \mu \leq 1$ , one can easily validate that

$$\tilde{\mu} = \frac{1 - (1 - \eta\mu)^H}{\eta H} \geq \frac{\mu}{2}. \quad (59)$$

So it follows that

$$\begin{aligned} \mathbb{E} \|\mathbf{w}^{(T)} - \mathbf{w}^*\|^2 &\leq (1 - \frac{1}{2}\alpha\eta H \mu)^T \|\mathbf{w}^{(0)} - \mathbf{w}^*\|^2 + \frac{2\alpha\eta H}{\mu} \max_{\mathbf{w}} \text{Var}[\hat{\mathcal{G}}(\mathbf{w})] \\ &\quad + \frac{20}{\mu^2} \max_{\mathbf{w}} \mathbb{E}_c \|\delta_c(\mathbf{w})\|^2 + \frac{20\rho^2}{\mu^2}. \end{aligned} \quad (60)$$

At last, in order to satisfy  $\tilde{\alpha}\tilde{L} \leq 1/4$ , one can set  $\alpha \leq 1/8$ , such that

$$\tilde{\alpha}\tilde{L} = \alpha\eta H \cdot \frac{1 + (1 - \eta\mu)^H}{\eta H} = \alpha(1 + (1 - \eta\mu)^H) \leq 2\alpha \leq \frac{1}{4}. \quad (61)$$

## C Bound on Iterate Bias

In this section, we will provide an upper bound for the iterate bias (11). According to the local update rules, we have

$$\left\| \mathbf{w}_{c,\text{GD}}^{(H)} - \mathbb{E}[\mathbf{w}_c^{(H)}] \right\| = \left\| \mathbf{w}_{c,\text{GD}}^{(H-1)} - \mathbb{E}[\mathbf{w}_c^{(H-1)}] - \eta \nabla F_c(\mathbf{w}_{c,\text{GD}}^{(H-1)}) + \eta \mathbb{E}[\nabla F_c(\mathbf{w}_c^{(H-1)})] \right\| \quad (62)$$

$$\begin{aligned} &\leq \left\| \mathbf{w}_{c,\text{GD}}^{(H-1)} - \mathbb{E}[\mathbf{w}_c^{(H-1)}] - \eta \nabla F_c(\mathbf{w}_{c,\text{GD}}^{(H-1)}) + \eta \nabla F_c(\mathbb{E}[\mathbf{w}_c^{(H-1)}]) \right\| \\ &\quad + \eta \left\| \mathbb{E}[\nabla F_c(\mathbf{w}_c^{(H-1)})] - \nabla F_c(\mathbb{E}[\mathbf{w}_c^{(H-1)}]) \right\| \end{aligned} \quad (63)$$

$$\begin{aligned} &\leq (1 - \eta\mu) \left\| \mathbf{w}_{c,\text{GD}}^{(H-1)} - \mathbb{E}[\mathbf{w}_c^{(H-1)}] \right\| \\ &\quad + \eta \left\| \mathbb{E}[\nabla F_c(\mathbf{w}_c^{(H-1)})] - \nabla F_c(\mathbb{E}[\mathbf{w}_c^{(H-1)}]) \right\| \end{aligned} \quad (64)$$

For the second term, we have

$$\left\| \mathbb{E}[\nabla F_c(\mathbf{w}_c^{(H-1)})] - \nabla F_c(\mathbb{E}[\mathbf{w}_c^{(H-1)}]) \right\|^2 \leq \mathbb{E} \left\| \nabla F_c(\mathbf{w}_c^{(H-1)}) - \nabla F_c(\mathbb{E}[\mathbf{w}_c^{(H-1)}]) \right\|^2 \quad (65)$$

$$\leq L^2 \mathbb{E} \left\| \mathbf{w}_c^{(H-1)} - \mathbb{E}[\mathbf{w}_c^{(H-1)}] \right\|^2 \quad (66)$$

$$\leq 2\eta^2 L^2 \sigma^2 (H - 1) \quad (67)$$

where the last inequality comes from previous works Khaled et al. (2020); Glasgow et al. (2021). As a result, one can obtain

$$\left\| \mathbf{w}_{c,\text{GD}}^{(H)} - \mathbb{E}[\mathbf{w}_c^{(H)}] \right\| \leq (1 - \eta\mu) \left\| \mathbf{w}_{c,\text{GD}}^{(H-1)} - \mathbb{E}[\mathbf{w}_c^{(H-1)}] \right\| + \sqrt{2}\eta^2 L\sigma(H-1)^{\frac{1}{2}} \quad (68)$$

$$\leq \sqrt{2}\eta^2 L\sigma \sum_{h=0}^{H-1} (1 - \eta\mu)^{H-1-h} h^{\frac{1}{2}} \quad (69)$$

$$\leq \sqrt{2}\eta^2 L\sigma \left[ \sum_{h=0}^{H-1} (1 - \eta\mu)^{2(H-1-h)} \right]^{\frac{1}{2}} \left[ \sum_{h=0}^{H-1} h \right]^{\frac{1}{2}} \quad (70)$$

$$\leq \sqrt{2}\eta^2 L\sigma \left[ \sum_{h=0}^{H-1} (1 - \eta\mu)^{H-1-h} \right]^{\frac{1}{2}} \left[ \sum_{h=0}^{H-1} h \right]^{\frac{1}{2}} \quad (71)$$

$$= \left[ \frac{1 - (1 - \eta\mu)^H}{\eta\mu H} \right]^{\frac{1}{2}} \eta^2 L\sigma H(H-1)^{\frac{1}{2}}. \quad (72)$$

According to the definition of  $\delta(\mathbf{w})$ , we obtain

$$\mathbb{E}_c \|\delta_c(\mathbf{w})\|^2 \leq \mathbb{E}_c \left\| \mathbf{w}_{c,\text{GD}}^{(H)} - \mathbb{E}[\mathbf{w}_c^{(H)}] \right\|^2 \leq \frac{\tilde{\mu}\eta^2 L^2 \sigma^2 (H-1)}{\mu} \leq \eta^2 L^2 \sigma^2 (H-1) \quad (73)$$

where the last inequality comes from the fact that  $\tilde{\mu} \leq \mu$ .

## D Proof of Corollary 1

### D.1 Deterministic Setting

When clients perform local GD in each round, there is no stochastic noise. So the error upper bound (9) can be simplified as follows

$$\left\| \mathbf{w}^{(T)} - \mathbf{w}^* \right\|^2 \leq \left(1 - \frac{1}{2}\alpha\eta H\mu\right)^T \left\| \mathbf{w}^{(0)} - \mathbf{w}^* \right\|^2 + \frac{20\rho^2}{\mu^2}. \quad (74)$$

If  $H\mu \geq L$ , then the maximal learning rate is  $\eta = 1/H\mu$ . When  $\alpha = 1/8$ , the upper bound becomes

$$\left\| \mathbf{w}^{(T)} - \mathbf{w}^* \right\|^2 \leq \left(1 - \frac{1}{16}\right)^T \left\| \mathbf{w}^{(0)} - \mathbf{w}^* \right\|^2 + \frac{20\rho^2}{\mu^2} \quad (75)$$

$$\leq \exp\left(-\frac{T}{16}\right) \left\| \mathbf{w}^{(0)} - \mathbf{w}^* \right\|^2 + \frac{20\rho^2}{\mu^2}. \quad (76)$$

If  $H\mu \leq L$ , then the maximal learning rate is  $\eta = 1/L$ . When  $\alpha = 1/8$ , the upper bound becomes

$$\left\| \mathbf{w}^{(T)} - \mathbf{w}^* \right\|^2 \leq \left(1 - \frac{H\mu}{16L}\right)^T \left\| \mathbf{w}^{(0)} - \mathbf{w}^* \right\|^2 + \frac{20\rho^2}{\mu^2} \quad (77)$$

$$\leq \exp\left(-\frac{\mu HT}{16L}\right) \left\| \mathbf{w}^{(0)} - \mathbf{w}^* \right\|^2 + \frac{20\rho^2}{\mu^2}. \quad (78)$$

We can summarize the above two bounds as follows:

$$\left\| \mathbf{w}^{(T)} - \mathbf{w}^* \right\|^2 \leq \exp\left(-\frac{T}{16\kappa} \min\{\kappa, H\}\right) \left\| \mathbf{w}^{(0)} - \mathbf{w}^* \right\|^2 + \frac{20\rho^2}{\mu^2} \quad (79)$$

$$= \mathcal{O}\left(\exp\left(-\frac{T}{16\kappa} \min\{\kappa, H\}\right) + \rho^2\right). \quad (80)$$

## D.2 Stochastic Setting

Substituting the upper bounds for  $\text{Var}[\mathcal{G}(\mathbf{w})]$  and  $\delta(\mathbf{w})$  into (57) and setting  $\alpha = 1/8$ ,

$$\mathbb{E} \left\| \mathbf{w}^{(t+1)} - \mathbf{w}^* \right\|^2 \leq (1 - \tilde{\alpha}\tilde{\mu}) \left\| \mathbf{w}^{(t)} - \mathbf{w}^* \right\|^2 + \tilde{\alpha}^2 \frac{2\sigma^2}{MH} + \frac{5\tilde{\alpha}^3 \sigma^2 L^2 (H-1)}{\tilde{\mu} \alpha^2 H^2} + \frac{5\tilde{\alpha}}{\tilde{\mu}} \|\mathcal{G}(\mathbf{w}^*)\|^2 \quad (81)$$

$$\leq (1 - \tilde{\alpha}\tilde{\mu}) \left\| \mathbf{w}^{(t)} - \mathbf{w}^* \right\|^2 + \tilde{\alpha}^2 \frac{2\sigma^2}{MH} + \tilde{\alpha}^3 \frac{320\sigma^2 L^2}{\tilde{\mu}H} + \tilde{\alpha} \frac{5\rho^2}{\tilde{\mu}}. \quad (82)$$

After minor rearrangement, we can get

$$\mathbb{E} \left\| \mathbf{w}^{(t+1)} - \mathbf{w}^* \right\|^2 - \frac{5\rho^2}{\tilde{\mu}^2} \leq (1 - \tilde{\alpha}\tilde{\mu}) \left[ \left\| \mathbf{w}^{(t)} - \mathbf{w}^* \right\|^2 - \frac{5\rho^2}{\tilde{\mu}^2} \right] + \tilde{\alpha}^2 \frac{2\sigma^2}{MH} + \tilde{\alpha}^3 \frac{320\sigma^2 L^2}{\tilde{\mu}H}. \quad (83)$$

After total  $T$  communication rounds,

$$\mathbb{E} \left\| \mathbf{w}^{(t+1)} - \mathbf{w}^* \right\|^2 - \frac{5\rho^2}{\tilde{\mu}^2} \leq (1 - \tilde{\alpha}\tilde{\mu})^T \underbrace{\left[ \left\| \mathbf{w}^{(0)} - \mathbf{w}^* \right\|^2 - \frac{5\rho^2}{\tilde{\mu}^2} \right]}_{r_0} + \frac{2\tilde{\alpha}\sigma^2}{\tilde{\mu}MH} + \frac{320\tilde{\alpha}^2\sigma^2 L^2}{\tilde{\mu}^2 H}. \quad (84)$$

If we set  $\tilde{\alpha} = \frac{\nu}{\mu T}$ , where  $\nu = 2 \ln(\max\{r_0\mu^2 MHT/(8\sigma^2), r_0\mu^4 HT^2/(1280L^2\sigma^2)\})$ , then it follows that

$$\mathbb{E} \left\| \mathbf{w}^{(t+1)} - \mathbf{w}^* \right\|^2 - \frac{20\rho^2}{\mu^2} \leq \frac{8\sigma^2\nu}{\mu^2 MHT} + \frac{1280\sigma^2 L^2 \nu}{\mu^4 HT^2} + \exp\left(-\frac{\nu}{2}\right) r_0 \quad (85)$$

$$\leq \frac{8\sigma^2(\nu+1)}{\mu^2 MHT} + \frac{1280\sigma^2 L^2 (\nu+1)}{\mu^4 HT^2} \quad (86)$$

$$= \tilde{\mathcal{O}} \left( \frac{\sigma^2}{MHT} + \frac{\sigma^2}{HT^2} \right) \quad (87)$$

where  $\tilde{\mathcal{O}}$  hides logarithmic factors.

## E Extensions to General Convex Functions

In this section, we are going to extend Theorem 2 to general convex settings for deterministic FEDAVG. This extension can help to show that our conclusion ‘‘FedAvg can have identical convergence rate in homogeneous and heterogeneous settings’’ is not only constrained to the strongly convex settings.

When the average client drift at optimum is zero and there is no stochastic noise, we have

$$\left\| \mathbf{w}^+ - \mathbf{w}^* \right\|^2 = \left\| \mathbf{w} - \mathbf{w}^* - \alpha \mathcal{G}(\mathbf{w}) \right\|^2 \quad (88)$$

$$= \left\| \mathbf{w} - \mathbf{w}^* \right\|^2 - 2\alpha \langle \mathbf{w} - \mathbf{w}^*, \mathcal{G}(\mathbf{w}) - \mathcal{G}(\mathbf{w}^*) \rangle + \alpha^2 \left\| \mathcal{G}(\mathbf{w}) - \mathcal{G}(\mathbf{w}^*) \right\|^2 \quad (89)$$

$$\leq (1 - \alpha\tilde{\mu}) \left\| \mathbf{w} - \mathbf{w}^* \right\|^2 - \frac{\alpha}{\tilde{L}} (1 - \alpha\tilde{L}) \left\| \mathcal{G}(\mathbf{w}) - \mathcal{G}(\mathbf{w}^*) \right\|^2 \quad (90)$$

$$\leq (1 - \alpha\tilde{\mu}) \left\| \mathbf{w} - \mathbf{w}^* \right\|^2 - \frac{\alpha}{2\tilde{L}} \left\| \mathcal{G}(\mathbf{w}) - \mathcal{G}(\mathbf{w}^*) \right\|^2 \quad (91)$$

where the last inequality is due to  $\alpha\tilde{L} \leq 1/2$ . In the general convex setting, we have  $\mu = 0$ . According to the definitions of  $\tilde{\mu}, \tilde{L}$  in Lemma 2, it follows that  $\tilde{\mu} = 0$  and  $\tilde{L} = 2/\eta H$ . Substituting these values into (91), we obtain

$$\left\| \mathbf{w}^+ - \mathbf{w}^* \right\|^2 \leq \left\| \mathbf{w} - \mathbf{w}^* \right\|^2 - \frac{\alpha\eta H}{4} \left\| \mathcal{G}(\mathbf{w}) - \mathcal{G}(\mathbf{w}^*) \right\|^2. \quad (92)$$

After minor rearrangement, we have

$$\left\| \mathcal{G}(\mathbf{w}) - \mathcal{G}(\mathbf{w}^*) \right\|^2 \leq \frac{4}{\alpha\eta H} \left[ \left\| \mathbf{w} - \mathbf{w}^* \right\|^2 - \left\| \mathbf{w}^+ - \mathbf{w}^* \right\|^2 \right]. \quad (93)$$

Taking the average from  $t = 0$  to  $t = T - 1$ ,

$$\frac{1}{T} \sum_{t=0}^{T-1} \left\| \mathcal{G}(\mathbf{w}^{(t)}) - \mathcal{G}(\mathbf{w}^*) \right\|^2 \leq \frac{4 \|\mathbf{w}^{(0)} - \mathbf{w}^*\|^2}{\alpha \eta H T} \quad (94)$$

If we set  $\alpha = 1/4, \eta = 1/L$ , then

$$\frac{1}{T} \sum_{t=0}^{T-1} \left\| \mathcal{G}(\mathbf{w}^{(t)}) - \mathcal{G}(\mathbf{w}^*) \right\|^2 \leq \frac{16L \|\mathbf{w}^{(0)} - \mathbf{w}^*\|^2}{H T}. \quad (95)$$

The above rate is the same as GD for general convex functions and data heterogeneity does not have negative impacts, as  $\rho = \|\mathcal{G}(\mathbf{w})\| = 0$ .

## F Proof of Theorem 3

*Proof.* According to the local update rule, we have

$$\mathbf{w}_c^{(t,h+1)} = \mathbf{w}_c^{(t,h)} - \eta \nabla F_c(\mathbf{w}_c^{(t,h)}) \quad (96)$$

$$= \mathbf{w}_c^{(t,h)} - \eta \left[ \mathbf{A}_c(\mathbf{w}_c^{(t,h)} - \mathbf{w}^*) - \mathbf{b}_c \right] \quad (97)$$

$$= (\mathbf{I} - \eta \mathbf{A}_c) \mathbf{w}_c^{(t,h)} + \eta (\mathbf{A}_c \mathbf{w}^* + \mathbf{b}_c). \quad (98)$$

Subtracting  $\mathbf{w}_c^* = \mathbf{w}^* + \mathbf{A}_c^{-1} \mathbf{b}_c$  on both sides, it follows that

$$\mathbf{w}_c^{(t,h+1)} - \mathbf{w}_c^* = (\mathbf{I} - \eta \mathbf{A}_c) (\mathbf{w}_c^{(t,h)} - \mathbf{w}_c^*) \quad (99)$$

$$= (\mathbf{I} - \eta \mathbf{A}_c)^{h+1} (\mathbf{w}^{(t)} - \mathbf{w}_c^*). \quad (100)$$

Setting  $h = H$ , we have  $\mathbf{w}_c^{(t,H)} = (\mathbf{I} - \eta \mathbf{A}_c)^H (\mathbf{w}^{(t)} - \mathbf{w}_c^*) + \mathbf{w}_c^*$ . Recall the definition of pseudo-gradient (6), we get

$$\mathcal{G}_c(\mathbf{w}^{(t)}) = \frac{1}{\eta H} (\mathbf{w}^{(t)} - \mathbf{w}^{(t,H)}) \quad (101)$$

$$= \frac{1}{\eta H} [\mathbf{I} - (\mathbf{I} - \eta \mathbf{A}_c)^H] (\mathbf{w}^{(t)} - \mathbf{w}_c^*). \quad (102)$$

According to the global update rule of FEDAVG, one can obtain that

$$\mathbf{w}^{(t+1)} = \mathbf{w}^{(t)} - \alpha \eta H \mathbb{E}_c \mathcal{G}_c(\mathbf{w}^{(t)}) \quad (103)$$

$$= \mathbf{w}^{(t)} - \alpha \mathbb{E}_c \left[ (\mathbf{I} - (\mathbf{I} - \eta \mathbf{A}_c)^H) (\mathbf{w}^{(t)} - \mathbf{w}_c^*) \right] \quad (104)$$

$$= \mathbf{w}^{(t)} - \alpha \mathbb{E}_c \left[ (\mathbf{I} - (\mathbf{I} - \eta \mathbf{A}_c)^H) (\mathbf{w}^{(t)} - \mathbf{w}^*) \right] - \alpha \mathbb{E}_c \left[ (\mathbf{I} - (\mathbf{I} - \eta \mathbf{A}_c)^H) (\mathbf{w}^* - \mathbf{w}_c^*) \right]. \quad (105)$$

Subtracting  $\mathbf{w}^*$  on both sides and setting  $\alpha = 1$ , we have

$$\mathbf{w}^{(t+1)} - \mathbf{w}^* = \mathbb{E}_c \left[ (\mathbf{I} - \eta \mathbf{A}_c)^H \right] (\mathbf{w}^{(t)} - \mathbf{w}^*) - \mathbb{E}_c \left[ \underbrace{(\mathbf{I} - (\mathbf{I} - \eta \mathbf{A}_c)^H) (\mathbf{w}^* - \mathbf{w}_c^*)}_{\eta H \mathcal{G}_c(\mathbf{w}^*)} \right] \quad (106)$$

$$= \mathbb{E}_c \left[ (\mathbf{I} - \eta \mathbf{A}_c)^H \right] (\mathbf{w}^{(t)} - \mathbf{w}^*) - \eta H \mathcal{G}(\mathbf{w}^*) \quad (107)$$

where  $\mathcal{G}(\mathbf{w}^*) = \mathbb{E}_c \mathcal{G}_c(\mathbf{w}^*)$ .

Now, we prove that  $\mathcal{G}(\mathbf{w}^*) = 0$  almost surely as  $M \rightarrow \infty$  on this synthetic problem. According to the definition of  $\mathbf{A}_c, \mathbf{b}_c$ , we obtain that

$$\mathcal{G}(\mathbf{w}^*) = \mathbb{E}_c \left[ (\mathbf{I} - (\mathbf{I} - \eta \mathbf{A}_c)^H) \mathbf{A}_c^{-1} \mathbf{b}_c \right] \quad (108)$$

$$= \mathbb{E}_c \left[ \underbrace{(\mathbf{I} - (\mathbf{I} - \eta \mathbf{A}_c)^H) \mathbf{A}_c^{-1} \frac{1}{n} \sum_{i=1}^n \mathbf{x}_{c,i} \epsilon_{c,i}}_{\xi_c} \right]. \quad (109)$$

Since the noise  $\epsilon_{c,i}$  are independent of  $\mathbf{x}_{c,i}$ ,  $\xi_c$  is a zero-mean random variable that depends on client  $c$ . Since we have assumed that all  $\|\mathbf{x}_{c,i}\|$  and  $\epsilon_{c,i}$  have bounded variance, we know that  $\mathbb{E}[\xi_c^2] < \infty$ . Since we have a uniform weighting on the  $M$  clients, it follows  $\mathbb{E}_c[\xi_c] = O(1/\sqrt{M})$  with probability  $1 - o_M(1)$ , and as  $M \rightarrow \infty$ , we have  $\mathbb{E}_c[\xi_c] \rightarrow 0$  almost surely.

Hence, from (107), we conclude that

$$\mathbf{w}^{(t+1)} - \mathbf{w}^* = \left[ \mathbb{E}_c \left[ (\mathbf{I} - \eta \mathbf{A}_c)^H \right] \right]^{t+1} (\mathbf{w}^{(0)} - \mathbf{w}^*), \quad (110)$$

almost surely as  $M \rightarrow \infty$ , which proves the desired result.

□