
DISTRIBUTED GRADIENT DESCENT WITH MANY LOCAL STEPS IN OVERPARAMETERIZED MODELS

Anonymous authors

Paper under double-blind review

ABSTRACT

In distributed training of machine learning models, gradient descent with *local iterative steps* is a very popular method, variants of which are commonly known as Local-SGD or the Federated Averaging (FedAvg). In this method, gradient steps based on local datasets are taken independently in distributed compute nodes to update the local models, which are then aggregated intermittently. Although the existing convergence analysis suggests that with heterogeneous data, FedAvg encounters quick performance degradation as the number of local steps increases, [it is shown to work quite well in practice, especially in the distributed training of large language models](#). In this work we try to explain this good performance from a viewpoint of implicit bias in Local Gradient Descent (Local-GD) with a large number of local steps. In overparameterized regime, the gradient descent at each compute node would lead the model to a specific direction locally. We characterize the dynamics of the aggregated global model and compare it to the centralized model trained with all of the data in one place. In particular, we analyze the implicit bias of gradient descent on linear models, for both regression and classification tasks. Our analysis shows that the aggregated global model converges exactly to the centralized model for regression tasks, and converges (in direction) to the same feasible set as centralized model for classification tasks. We further propose a Modified Local-GD with a refined aggregation and theoretically show it converges to the centralized model in direction for linear classification. We empirically verified our theoretical findings in linear models and also conducted experiments on distributed fine-tuning of pretrained neural networks to further apply our theory.

1 INTRODUCTION

In this era of large machine learning models, distributed training is an essential part of machine learning pipelines. It can happen in a data center with thousands connected compute nodes [Sergeev & Del Balso \(2018\)](#); [Huang et al. \(2019\)](#), or across several data centers and millions of mobile devices in federated learning [Konečný et al. \(2016\)](#); [Kairouz et al. \(2019\)](#). In such a network, the communication cost is usually the bottleneck in the whole system. To alleviate communication burden, and also to preserve privacy to some extent, one common strategy is to perform multiple local updates before sending the information to other nodes, which is called Local Gradient Descent (Local-GD) [Stich \(2019\)](#); [Lin et al. \(2019\)](#). It is also a standard algorithm in federated learning, varied by partial device participation and privacy constraints, and known as FedAvg [McMahan et al. \(2017\)](#). While local updates can reduce communication cost, the number of local steps is usually considered to be small [Stich \(2019\)](#); [Li et al. \(2020b\)](#). When data distributions across machines are heterogeneous, a large number of local steps would result in local iterates to diverge significantly (called client-drift), and the aggregated values to oscillate and be far away from the optimum global model.

However, in practical implementation of distributed training on large models, the performance of vanilla FedAvg is surprisingly good even with heterogeneous data distribution [McMahan et al. \(2017\)](#); [Charles et al. \(2021\)](#). In fact SCAFFOLD [Karimireddy et al. \(2020\)](#), an algorithm designed to mitigate the effect of heterogeneity theoretically, is shown to have similar empirical performance as FedAvg [Reddi et al. \(2021\)](#); [Wu et al. \(2023\)](#). There are some works trying to explain the effectiveness of FedAvg from different theoretical aspects, such as representation learning [Collins et al. \(2022\)](#), refined theoretical assumption [Wang et al. \(2024\)](#) etc. Also, the number of local steps can be very large in real-world systems, for example, performing 500 local steps in distributed training of large language

models (LLM) Douillard et al. (2023); Jaghouar et al. (2024). These practical experiences motivates us to consider the following question:

Q: Can we establish rigorous conditions, independent of data distribution, under which Local-GD performs well with a very large number of local steps?

In this work we answer this question in affirmative by considering overparameterized models on regression and classification tasks. Our main tool is to analyze the *implicit bias* of gradient descent to characterize the dynamics of aggregated models with many local steps. In a network with M compute nodes, the goal is to train a global model to fit in the distributed datasets:

$$\min_{w \in \mathbb{R}^d} f(w) \quad \text{with } f(w) \equiv \frac{1}{M} \sum_{i=1}^M f_i(w|D_i), \quad (1)$$

where $w \in \mathbb{R}^d$ is the single model to be trained and $f_i(w|D_i)$ is the local objective function, and D_i is the local distribution of d -dimensional samples and corresponding labels $\{x_{ij}, y_{ij}\}_{j=1}^N$.

To reduce the communication frequency, Local-GD chooses to do L local gradient descent steps before sending the local model to a central node. The detailed algorithm of Local-GD is described in Algorithm 1 and 2. In the existing convergence analysis of Local-GD, the number of local steps L should not be very large. **For example, with strongly convex and smooth loss functions, the number of local steps should not be larger than $O(\sqrt{T})$ for i.i.d data Stich (2019) and non-i.i.d. data Li et al. (2020b).** However, such analysis is developed for general/classical models and does not consider the special properties of overparameterized models. In this work we specifically focus on linear models for both regression and classification tasks and take the overparameterized regime into account. That is, the dimension d is larger than the total number of samples, i.e. $d > MN$. While modern machine learning concerns primarily large nonlinear models, it is instructive to explore the intrinsic property of Local-GD in simpler linear setting and establish the connection to other areas. For example, the leading theories of deep learning, such as implicit bias of optimization algorithms, or double descent Belkin et al. (2018; 2019), were built for linear models first. Moreover, fine-tuning on pretrained large models has gradually become the popular paradigm in practical machine learning pipeline. It is widely used to fine-tune the final linear layer or add a few linear layers to pretrained models in transfer learning Donahue et al. (2014); Kornblith et al. (2019) and deployment of LLM Devlin (2018); Jiang et al. (2020).

As stated, to characterize the behavior of Local-GD with large number of local steps in overparameterized models, we leverage the implicit bias of gradient descent, which is an active area in theoretical explanation of modern large models Soudry et al. (2018); Gunasekar et al. (2018a); Ji & Telgarsky (2019a); Chizat & Bach (2020); Frei et al. (2024). With a very large number of local steps, the local optimization problem can be exactly solved for linear regression and classification models. In overparameterized regime, gradient descent would converge to a specific solution. After aggregation of these specific local solutions, we can characterize the dynamic of the global model and finally compare it to the centralized model trained on a collection of distributed datasets at one place.

Specifically, in linear regression minimizing a squared loss, the local models would fit to the corresponding local datasets, and converge to the solution with minimum distance to initial aggregated global model at each communication round. We can obtain the closed form of this solution and calculate the global model after aggregation. We prove that it exactly converges to the centralized model (the model trained by gradient descent if all data were in one place) as the number of rounds of communication increases.

The analysis of linear classification (halfspace learning) is more involved and proceeds according to the following steps. First, it turns out that when minimizing an exponential loss with a weakly regularized term, the aggregated global model is equivalent to a model aggregated from local models obtained by solving *local max-margin* problems. Subsequently we relate the update of global model aggregated from solutions of local max-margin problems to *Parallel Projection Method (PPM)*, an iterative algorithm used for finding a point in the intersection of multiple constraint sets by projecting onto each constraint set in parallel Gilbert (1972); de Pierro & Iusem (1984); Combettes (1994; 1996). Using properties of PPM, we can characterize the dynamics of the aggregated global model. We prove that it converges to a global feasible set, which is the intersection of constraint sets in local max-margin problems. The centralized model trained with all of the data also converges to the global feasible set. To further explain the similar performance obtained by global model and centralized model, we propose a *modified* Local-GD with a different aggregation method from vanilla Local-GD (Algorithm 3). We theoretically prove that the aggregated global model obtained from Modified Local-GD exactly

108 converges to the centralized model in direction. We show the vanilla Local-GD actually converges to
 109 the same point as the modified Local-GD experimentally. For both linear regression and classification,
 110 our results show that the aggregated global model would converge to the centralized model even with
 111 a very large number of local steps on heterogeneous data.

112 In summary, the contribution of this work is as follows:

- 114 • We established the theoretical performance of Local-GD with a large number of local steps in
 115 overparameterized models. We analyzed the implicit bias of Local-GD, for single communication
 116 round of linear regression, and for whole algorithmic process of classification, respectively. As far
 117 as we know, this is the first attempt to analyze implicit bias of gradient descent in distributed setting.
 118
- 119 • We obtained closed form of the aggregated global model in linear regression and analyzed its dynam-
 120 ics. We proved that it exactly converges to the centralized model as communication rounds increase.
 121
- 122 • We related the Local-GD for linear classification to Parallel Projection Method and characterized the
 123 dynamics based on the properties of projections. We proved the aggregated global model converges
 124 to a global feasible set same as the centralized model.
- 125 • We further proposed a Modified Local-GD with a different aggregation method and proved it
 126 converges exactly to the centralized model in direction.
- 127 • We experimentally verify our theoretical findings on synthetic datasets and real datasets with linear
 128 models. We further conducted experiments on fine-tuning the final linear layer of neural networks
 129 to show the broader impact of our work.
 130

131 Our main technical challenge comes while analyzing classification. In linear regression, the implicit
 132 bias for a single round of communication is directly derived from the gradient on squared loss (each
 133 gradient step is on the row space of local data). In contrast, for classification we have to consider
 134 the whole algorithmic process of both Local-GD and Parallel Projection Method and then derive
 135 the equivalence between them. Compared to the continual learning work Evron et al. (2023) where
 136 overparameterized models are handled sequentially, the challenge is that we need to handle the parallel
 137 projections happening *simultaneously* from the same initial point. Due to space limit, we give more
 138 additional references and discussion on Related Works in Appendix A.

140 **Algorithm 1** LOCAL-GD.

141 1: **Input:** learning rate η .
 142 2: Initialize w_0^0
 143 3: **for** $k=0$ to $K-1$ **do**
 144 4: The aggregator sends global model w_0^k to all compute nodes.
 145 5: **for** $i=1$ to $i=M$ **do**
 146 6: compute node i updates local model starting from w_0^k : $w_i^{k+1} = \text{LocalUpdate}(w_0^k)$.
 147 7: compute node i sends back the updated local model w_i^{k+1} .
 148 8: **end for**
 149 9: The aggregator aggregates all the local models: $w_0^{k+1} = \frac{1}{M} \sum_{i=1}^M w_i^{k+1}$.
 150 10: **end for**
 151 11: **Output:** w_0^K .

154 **Algorithm 2** LocalUpdate(w_0^k) in general Local-GD.

155 1: **Input:** an initial point w_0^k , the number of local steps L , and the learning rate η .
 156 2: Initialize $w_i^{k,0} = w_0^k$.
 157 3: **for** $l=0$ to $L-1$ **do**
 158 4: $w_i^{k,l+1} = w_i^{k,l} - \eta \nabla f_i(w_i^{k,l})$.
 159 5: **end for**
 160 6: **Output:** LocalUpdate(w_0^k) := $w_i^{k,L}$.

2 LOCAL-GD IN LINEAR REGRESSION: A WARM-UP

2.1 SETTING

In this section we first consider linear regression in overparameterized regime. The behavior of linear regression is very well-understood in high-dimensional statistics; and we can clearly convey our key message based on this fundamental setting.

At each compute node i , the dataset S_i consists of N tuples of samples and their corresponding labels, $(x, y) \in \mathbb{R}^d \times \mathbb{R}$. We assume the label y_{ij} is generated by

$$y_{ij} = x_{ij}^T w_i^* + z_{ij} \quad (2)$$

where $w_i^* \in \mathbb{R}^d$ is the ground truth model at i -th compute node, and z_{ij} is the added noise. Denote $X_i = [x_{i1}, x_{i2}, \dots, x_{iN}]^T \in \mathbb{R}^{N \times d}$ as the data matrix at i -th compute node, and $y_i = [y_{i1}, y_{i2}, \dots, y_{iN}] \in \mathbb{R}^N$ as the label vector, $z_i \in \mathbb{R}^N$ as the noise vector. In heterogeneous setting, the w_i^* can be very different to each other. Note that the convergence to centralized model does not rely on the generative model. We just make this assumption on generative model for deriving a more clear form of the aggregated global model.

Algorithm. At each round, the aggregator sends the global model w_0 to all the compute nodes. Each compute node minimizes the squared loss $f_i(w_i) = \frac{1}{2N} \|y_i - X_i w_i\|^2$ by a large number of gradient descent steps *until convergence*. Then each compute node sends back the local model and the aggregator aggregates all the local models to get the updated global model. [The detailed algorithm is Local-GD in Algorithm 1 with \$f_i\(w_i\)\$ replaced in LocalUpdate \(Algorithm 2\)](#). Since minimizing squared loss is a quadratic problem, it is expected to reach convergence locally with a small number of gradient descent steps.

2.2 IMPLICIT BIAS OF LOCAL GD IN LINEAR REGRESSION

For each local problem, when the dimension of the model is larger than the number of samples at each compute node ($d > N$), i.e., locally overparameterized, there are multiple solutions corresponding to zero squared loss. However, gradient descent will lead the model converge to a specific solution, which corresponds to a minimum Euclidean distance to the initial point Gunasekar et al. (2018a); Evron et al. (2022). Formally, the solution w_i^{k+1} obtained at k -th round and i -th node will converge to the solution of the optimization problem

$$\min_{w_i} \|w_i - w_0^k\|^2 \quad \text{s.t.} \quad X_i w_i = y_i. \quad (3)$$

We can obtain the closed form solution of this optimization problem as (see Proof of Lemma 1 in Appendix C.1)

$$\begin{aligned} w_i^{k+1} &= (I - X_i^T (X_i X_i^T)^{-1} X_i) w_0^k + X_i^T (X_i X_i^T)^{-1} y_i \\ &= (I - X_i^T (X_i X_i^T)^{-1} X_i) w_0^k + X_i^T (X_i X_i^T)^{-1} X_i w_i^* + X_i^T (X_i X_i^T)^{-1} z_i. \end{aligned} \quad (4)$$

Denote $P_i \triangleq X_i^T (X_i X_i^T)^{-1} X_i$ and $X_i^\dagger \triangleq X_i^T (X_i X_i^T)^{-1}$. The local model can be rewritten as $w_i^{k+1} = (I - P_i) w_0^k + P_i w_i^* + X_i^\dagger z_i$. We observe that P_i is the projection operator to the row space of X_i , and X_i^\dagger is the pseudo inverse of X_i . After one round of iterations, the local model is actually an interpolation between the initial global model w_0^k at this round and the ground-truth model w_i^* , plus a noise term. We then obtain the closed form of global model by aggregation. After many rounds of communication, we can obtain the final trained global model from Local-GD.

Lemma 1. *When the local overparameterized linear regression problems are exactly solved by gradient descent, then after K rounds of communication, the global model w_0^K obtained from Local-GD is*

$$w_0^K = (I - \bar{P})^K w_0^0 + \sum_{k=0}^{K-1} (I - \bar{P})^k (\bar{Q} + \bar{Z}), \quad (5)$$

where $\bar{P} = \frac{1}{M} \sum_{i=1}^M P_i$, $\bar{Q} = \frac{1}{M} \sum_{i=1}^M P_i w_i^*$, $\bar{Z} = \frac{1}{M} \sum_{i=1}^M X_i^\dagger z_i$.

Note that $\bar{P}, \bar{Q}, \bar{Z}$ are constant after the data is generated. Since we only know the $\{X_i, y_i\}_{i=1}^M$ in the training process, we can also write it as

$$w_0^K = (I - \bar{P})^K w_0^0 + \sum_{k=0}^{K-1} (I - \bar{P})^k \bar{Y}, \quad (6)$$

where $\bar{Y} = \frac{1}{M} \sum_{i=1}^M X_i^\dagger y_i$. Then we can directly get the final model from the training set.

Singularity of \bar{P} . If \bar{P} is invertible, we can further simplify the form of global model. However, since $P_i \in \mathbb{R}^{d \times d}$ is the projection operator onto row space of X_i , its rank is at most N . The \bar{P} is the average of P_i s, thus its rank is at most MN . Note that we consider the overparameterized regime both locally and globally, i.e., $d \gg MN$. Then \bar{P} is singular, and the sum $\sum_{k=0}^{K-1} (I - \bar{P})^k$ approaches KI when d becomes very large. We cannot get more properties of the final global model from (6), but we can compare it to the centralized model trained with all of the data.

2.3 CONVERGENCE TO CENTRALIZED MODEL

Let $X_c = [X_1^T, \dots, X_M^T]^T \in \mathbb{R}^{MN \times d}$ be the data matrix consisting of all the local data, and $y_c = [y_1^T, \dots, y_M^T]^T \in \mathbb{R}^{MN \times 1}$ be the label vector consisting of the local labels. If we train the centralized model from initial point 0 with squared loss, then the gradient descent will lead the model to the solution of the optimization problem

$$\min_w \|w\|^2 \quad \text{s.t.} \quad X_c w = y_c \quad (7)$$

We can write the closed form of centralized model as $w_c = X_c^T (X_c X_c^T)^{-1} y_c$.

Due to the constraint in problem (7), for each compute node i , we have $X_i w_c = y_i$. We replace y_i in the local model (4), then we have

$$w_i^{k+1} - w_c = (I - P_i)(w_0^k - w_c). \quad (8)$$

The right-hand side is projecting the difference between global model and centralized model onto null space of X_i . After averaging all the local models at the aggregator, we have

$$w_0^{k+1} - w_c = (I - \bar{P})(w_0^k - w_c). \quad (9)$$

In the training process the difference between global model and centralized model is iteratively projected onto the null space of span of row spaces of X_i s. It implies that the difference on the span of data matrix gradually decreases until zero. Based on the evolution of the difference, we can prove the following theorem:

Theorem 1. *For the linear regression problem, suppose the initial point w_0^0 is 0 and $d \gg MN$, then the global model obtained by Local-GD, w_0^K , converges to the centralized solution w_c as the number of communication rounds $K \rightarrow \infty$.*

The proof is deferred in Appendix C.2. The key step is to show the initial difference is already in the data space, and no residual in the null space of row spaces of X_i s.

Due to the linearity of the regression problem, we can theoretically show the global model can exactly converge to the centralized model with implicit bias on overparameterized regime. Note that the proof does not rely on the generative model and assumption on data heterogeneity. It implies that, even if we use a large number of local steps to exactly solve the local problems on very heterogeneous data, the performance of Local-GD is equivalent to train a model with all the data in one place.

3 LOCAL-GD IN LINEAR CLASSIFICATION: RELATION TO PPM

3.1 SETTING

In this section we investigate a binary classification task with linear models. Different from the linear regression problem, it is hard to obtain closed form solution on classification tasks. Thus we need to develop new techniques to handle this case.

Suppose, for each compute node i , the dataset S_i consists of N tuples of samples and their corresponding labels, $(x, y) \in \mathbb{R}^d \times \{+1, -1\}$. Similarly, we denote $X_i \in \mathbb{R}^{N \times d}$ as the data matrix at i -th compute node, and $y_i \in \{+1, -1\}^N$ as the label vector. We do not assume the generative model in classification task, but we need an assumption of separable datasets.

Assumption 1. Each local dataset S_i is separable, i.e., there are non-empty local feasible sets,

$$C_i \triangleq \{w \in \mathbb{R}^d \mid y_{ij} x_{ij}^T w \geq 1, \text{ for } j = 1, \dots, N\}, \quad (10)$$

and there is a non-empty global feasible set,

$$\bar{C} \triangleq \bigcap_{i=1}^m C_i \neq \emptyset. \quad (11)$$

This assumption makes sure that the datasets are locally and globally separable.

Algorithm. At each round, the aggregator sends the global model w_0 to all the compute nodes. Each compute node minimizes an exponential loss with a weakly regularized term by many gradient descent steps *until convergence*. That is, each compute node solves the following problem:

$$\min_{w \in \mathbb{R}^d} f_i(w) = \sum_{j=1}^N \exp(-y_{ij} x_{ij}^T w) + \frac{\lambda}{2} \|w - w_0^k\|^2 \quad (12)$$

where λ is a regularization parameter close to 0.

Then each compute node sends back the local model and the aggregator aggregates all the local models to get the updated global model. [The detailed algorithm for linear classification is Local-GD in Algorithm 1 with \$f_i\(w_i\)\$ replaced in LocalUpdate \(Algorithm 2\).](#)

Regularization methods are very common in distributed learning to force the local models move not too far from global model Li et al. (2020a; 2021); T Dinh et al. (2020). Here we consider the weakly regularized term, $\lambda \rightarrow 0$, to give theoretical insights of Local-GD on classification tasks. Experimentally the λ is set to be extremely small that does not affect the minimization of exponential loss. Since the local problem is a strongly convex problem, with many local gradient descent steps it will be exactly solved.

3.2 IMPLICIT BIAS OF GRADIENT DESCENT IN LINEAR CLASSIFICATION

One can derive the implicit bias of classification at a single local node after a large number of local steps. However, in contrast to linear regression, we cannot easily aggregate the local solutions after a round of communication to a closed form. At each round, the local model is updated from the previously aggregated global model, which is related to previous local updates. To mitigate this, we consider the whole algorithmic process of Local-GD on classification and use another auxiliary sequence of global models, denoted as $\bar{w}_0^k, k = 0, 1, 2, \dots$. Starting from an initial point \bar{w}_0^0 , the central node sends global model \bar{w}_0^k to all the compute nodes at k -th iteration round. Each compute node solves the following Local Max-Margin problem to obtain \bar{w}_i^{k+1} :

$$\bar{w}_i^{k+1} = \arg \min_{w \in \mathbb{R}^d} \|w - \bar{w}_0^k\| \quad \text{s.t.} \quad y_{ij} x_{ij}^T w \geq 1 \quad j = 1, 2, \dots, N. \quad (13)$$

Then the compute node sends the local model back. The central node averages the local models to get $\bar{w}_0^{k+1} = \frac{1}{M} \sum_{i=1}^M \bar{w}_i^{k+1}$.

We can show the solution w_0^K obtained in Local-GD converges in direction to the global model from Local Max-Margin problems \bar{w}_0^K .

Lemma 2. *For almost all datasets sampled from a continuous distribution satisfying Assumption 1, with initialization $w_0^0 = \bar{w}_0^0 = 0$, we have $w_0^k \rightarrow \ln(\frac{1}{\lambda}) \bar{w}_0^k$, and the residual $\|w_0^k - \ln(\frac{1}{\lambda}) \bar{w}_0^k\| = O(k \ln \ln \frac{1}{\lambda})$, as $\lambda \rightarrow 0$. It implies that at any round $k = o\left(\frac{\ln(1/\lambda)}{\ln \ln(1/\lambda)}\right)$, w_0^k converges in direction to \bar{w}_0^k :*

$$\lim_{\lambda \rightarrow 0} \frac{w_0^k}{\|w_0^k\|} = \frac{\bar{w}_0^k}{\|\bar{w}_0^k\|}. \quad (14)$$

The proof is deferred in Appendix D. The framework is similar to the continual learning work Evron et al. (2023), but we need to handle the parallel local updates for each dataset from the same initial

model and the aggregation, which is different from the sequential updates where for each dataset the model is trained from the previous model and there is no need to do aggregation.

Based on this equivalence between Local-GD for linear classification and Local Max-Margin scheme, we can further analyze the performance of Local-GD with a large number of local steps. Instead of a closed-form solution for the Local Max-Margin problem (13), we treat it as a projection of the aggregated global model onto a convex set $C_i: \bar{w}_i^{k+1} = P_i(\bar{w}_0^k)$, which is formed by the constraints in (13) and exactly the local feasible set defined in Assumption 1. Here we slightly overload the notation P_i , which was used as the projection matrix in linear regression since the readers can get a sense of the same effect of them in Local-GD. The aggregation is actually to average the local projected points: $\bar{w}_0^{k+1} = \frac{1}{M} \sum_{i=1}^M P_i(\bar{w}_0^k)$.

The sequence of Local Max-Margin schemes is therefore projections to local (convex) feasible sets followed by aggregation, which is the Parallel Projection Method (PPM) in literature Gilbert (1972); Combettes (1994). Using Lemma 2, we establish the relation between Local-GD and PPM: the model from Local-GD converges to the model from PPM in direction.

3.3 CONVERGENCE TO GLOBAL FEASIBLE SET

Now we use the properties of PPM to characterize the performance of Local-GD in classification. In Combettes (1994), the convergence of PPM has been provided for a relaxed version. The direct average considered in this work can be seen as a special case of the relaxed version, and the following lemma holds.

Lemma 3 (Theorem 1 and Proposition 8, Combettes (1994)). *Suppose all the local feasible sets $C_i, i = 1, 2, \dots$ are closed and convex, and the intersection \bar{C} is not empty. Then for any initial point \bar{w}_0^0 , the global model \bar{w}_0 generated by PPM converges to a point in the global feasible set \bar{C} .*

This lemma guarantees that \bar{w}_0^K will converge to the intersection of the convex sets after many rounds of iteration, however we are not sure which exact point it would converge to.

Similar to linear regression case, we also compare the global model obtained from Local-GD to the centralized model trained with all of data in one place. From the implicit bias of gradient descent on exponential-tailed loss Soudry et al. (2018), the centralized model trained with exponential loss will converge in direction to the solution of a Max-Margin problem:

$$\min_{w \in \mathbb{R}^d} \|w\| \quad \text{s.t.} \quad y_{ij} x_{ij}^T w \geq 1, \quad i = 1, 2, \dots, M, \quad j = 1, 2, \dots, N. \quad (15)$$

This problem is actually the problem of hard margin support vector machine (SVM). The constraints in equation 15 include all the local datasets, and form the global feasible set \bar{C} . That is, the centralized model would converge to the minimum norm solution in global feasible set in direction.

Combining Lemma 2, Lemma 3 and result of centralized model, we immediately have:

Theorem 2. *For linear classification problem with exponential loss, suppose initial point is $w_0^0 = 0$. The aggregated global model w_0^K obtained by Local-GD with a large number of local steps converges in direction to one point in the global feasible set \bar{C} , while the centralized model converges in direction to the minimum norm point in the same set.*

The main difference from linear regression is that we cannot guarantee the global model obtained by Local-GD to converge exactly to the centralized model in classification, but show that it converges to the same global feasible set as the centralized solution. Nevertheless, in experiments the test accuracy of the Local-GD model is very similar to that of centralized model. To theoretically support that the Local-GD model converges to the centralized model, we propose a slightly Modified Local-GD by just changing the aggregation method, and showing that it converges to the centralized model exactly.

4 MODIFIED LOCAL-GD: CONVERGENCE TO CENTRALIZED MODEL

Previously, we established the connection between Local-GD and PPM in linear classification. In Combettes (1996) it was shown that if the aggregation method is modified to incorporate the influence of the initial point \bar{w}_0^0 in PPM, then the sequence generated by PPM will converge to a specific point

in global feasible set \bar{C} with minimum distance to this initial point. Denote $P_c(\cdot)$ as the projection operator onto the global feasible set \bar{C} . Formally we have the following lemma.

Lemma 4 (Theorem 5.3, Combettes (1996)). *Suppose \bar{C} is not empty. For any initial point \bar{w}_0^0 , when the local models are aggregated as*

$$\bar{w}_0^{k+1} = (1 - \alpha^{k+1})\bar{w}_0^0 + \alpha^{k+1} \left(\frac{1}{M} \sum_{i=1}^M P_i(\bar{w}_0^k) \right), \quad (16)$$

where $\{\alpha^k\}$ satisfy (i) $\lim_{k \rightarrow \infty} \alpha^k = 1$, (ii) $\sum_{k \geq 0} (1 - \alpha^k) = \infty$, (iii) $\sum_{k \geq 0} |\alpha^{k+1} - \alpha^k| < \infty$, then the global model generated by PPM will converge to the point $P_c(\bar{w}_0^0)$.

That is the sequence generated by PPM would converge to the point in global feasible set, \bar{C} , with minimum distance to \bar{w}_0^0 . The modified aggregation method is a linear combination of initial point and current average of local projected points. One example of the sequence $\{\alpha^k\}$ satisfying the conditions is $\alpha^k = 1 - \frac{1}{k+1}$.

If we start from $\bar{w}_0^0 = 0$, then the point $P_c(\bar{w}_0^0)$ is exactly the minimum norm point in the global feasible set. It shows the PPM can exactly converge to the minimum norm point as the centralized model. Based on this result, we propose a Modified Local-GD algorithm shown in Algorithm 3, which only differs from Local-GD in the aggregation method.

Algorithm 3 MODIFIED LOCAL-GD.

```

1: Input: learning rate  $\eta$ .
2: Initialize  $w_0^0$ 
3: for  $k=0$  to  $K-1$  do
4:   The central node sends global model  $w_0^k$  to all compute nodes.
5:   for  $i=1$  to  $i=M$  do
6:     compute node  $i$  updates local model starting from  $w_0^k$ :  $w_i^{k+1} = \text{LocalUpdate}(w_0^k)$ .
7:     compute node  $i$  sends back the updated local model  $w_0^{k+1}$ .
8:   end for
9:   The central node aggregates all the local models:  $w_0^{k+1} = (1 - \alpha^k)w_0^0 + \alpha^k \left( \frac{1}{M} \sum_{i=1}^M w_i^k \right)$ .
10: end for
11: Output:  $w_0^K$ .

```

We still need to prove a lemma analogous to Lemma 2 to establish the equivalence between Modified Local-GD and Modified PPM, which is omitted here due to space limit (Please refer to Appendix E and the proof is very similar to proof in Lemma 2). From the equivalence, Lemma 4, and result of the centralized model, we can have the following theorem:

Theorem 3. *For linear classification problem, suppose the initial point is $w_0^0 = 0$. Then the global model w_0^K obtained by Modified Local-GD (Algorithm 3) converges in direction to the centralized model obtained from (15).*

Unlike the vanilla Local-GD, which is only guaranteed to converge to the global feasible set, the Modified Local-GD is guaranteed to converge to the centralized model in direction. Unlike linear regression, the convergence is established *in direction* since the solution on exponential loss could go to infinity.

Note that if we start from $\bar{w}_0^0 = 0$, the aggregation in Modified Local-GD becomes $w_0^{k+1} = \frac{k}{k+1} \left(\frac{1}{M} \sum_{i=1}^M w_i^k \right)$, which is just a *scaling* of vanilla aggregation with a parameter less than 1. Thus we can see experimentally they usually converge to the same point and Modified Local-GD converges slightly slower. In summary, Modified Local-GD theoretically illustrates that the global model trained from Local-GD could obtain similar performance as the centralized model.

5 EXPERIMENTS

Linear Regression. We simulated 10 compute nodes, each with 50 training samples. The label vector y_i at i -th compute node is exactly generated as (2), where ground truth model w_i^* is Gaussian

vector with each element following $\mathcal{N}(0,4)$. Each ground truth model at different compute nodes is independently generated, thus the datasets can be very different from each other. The data matrix X_i also follows Gaussian distribution, with each element being $\mathcal{N}(0,1)$, and z_i is a Gaussian vector with $\mathcal{N}(0,0.04)$. In Local-GD, the number of local steps is $L=200$, number of rounds is also $R=200$, and the learning rate $\eta=0.0001$. Actually it just take a few local steps to converge locally at each round, but we set a large number of local steps to show it can be large at $O(\sqrt{T})$, where $T=L*R$ is the number of total iterations. We tested the global model (G) from Local-GD on squared loss, centralized model (C) trained from global dataset on squared loss, closed form of global model (G-Closed) in (6), closed form of centralized model (C-Closed) as solution of problem (7). The centralized model is trained 10000 steps with learning rate 0.0001.

Fig. 1(a) displays the difference between global model and trained centralized model, and difference between global model and closed form of global model at each round when dimension is $d=1500$, which is locally and globally overparameterized. The difference between two models is $\|w_1 - w_2\|/d$. We can see the difference between global model and its closed form is always 0 during the training process, verifying the correctness of the derived closed form (6). The global model can gradually converge to the centralized model with more communication rounds.

Fig. 1(b) displays the difference between global model and centralized model, global model and its closed form, and centralized model and its closed form, with respect to model dimension. Since it is always locally overparameterized, the difference between global model and the closed form is always zero. The difference between global model and centralized model has an obvious peak around 500, which is the number of total samples. The phenomenon that global model converges exactly to centralized model only happens when the model is sufficiently overparameterized. Fig. 1(c) shows the generalization error of global model and centralized model in linear regression. Since the data matrix is Gaussian, the generalization error of model w can be computed as $\frac{1}{M} \sum_{i=1}^M \|w - w_i^*\|^2$. We plot the generalization error divided by d . It is shown the global model and centralized model can get the same performance when model is sufficiently overparameterized.

Classification. For linear classification, we also have 10 compute nodes, with 50 samples at each. The dataset is generated as $y_{ij} = \text{sign}(x_{ij}^T w_i^*)$, where ground truth model is $w_i^* = w^* + z_i$, and w^* is a Gaussian vector randomly chosen, z_i is a Gaussian noise. The data matrix X_i is still a Gaussian matrix. This setting makes sure the datasets across compute nodes are different from each other, meanwhile they are not totally different such that there may be a non-empty global feasible set. The global model is trained exactly as Local-GD for linear classification, where the λ is 0.0001. Actually we can use the standard logistic regression without regularization to obtain the same performance. But aligning with theoretical proof, we still use exponential loss with a very weak regularization. We tested global model (G), global model from Modified Local-GD (G-Mod), centralized model (C) from minimizing exponential loss on all the data, centralized SVM model (S) solved from problem (15) via standard scikit-learn package. Note that centralized model and SVM model are the final trained model in the plots. In Local-GD, the number of local steps is $L=150$, the number of communication rounds is $R=120$, and the learning rate is $\eta=0.01$. The centralized model is trained with same learning rate for 20000 steps. Since our theory claimed the convergence is established in direction, the difference computed here for two models is defined after normalization $\|w_1/\|w_1\| - w_2/\|w_2\|\|$.

Fig. 1(d) shows the difference between these models with respect to the number of rounds R when dimension is $d=1500$. We can see both global model and modified global model converges to the centralized model in direction, and the centralized model is close to the SVM model but there is small gap. Fig. 1(e) displays the difference with respect to dimension d . It is seen the difference between global model and centralized model gradually decreases with larger dimensions. The modified global model is almost the same as the centralized model but the gap is slightly larger since it converges slower than vanilla global model with same number of rounds. Fig. 1(f) shows the difference from SVM model with dimension. The gap between the models to SVM model also decreases with larger d . Finally Fig. 1(g) plots the test accuracy of these models. The test datasets are also constructed by the same generation of training set with different data matrix. Although the accuracy decreases with larger dimension (relatively fewer samples), the performance of global models and centralized models are always similar.

Fine-Tuning of Pretrained Neural Network. We further fine-tuned the ResNet50 model pretrained with ImageNet dataset on CIFAR10 dataset. Only the final linear layer is trained during the process, while the rest of model is fixed. The 50000 samples are distributed on 10 compute nodes. For i -th

486
487
488
489
490
491
492
493
494
495
496
497
498
499
500
501
502
503
504
505
506
507
508
509
510
511
512
513
514
515
516
517
518
519
520
521
522
523
524
525
526
527
528
529
530
531
532
533
534
535
536
537
538
539

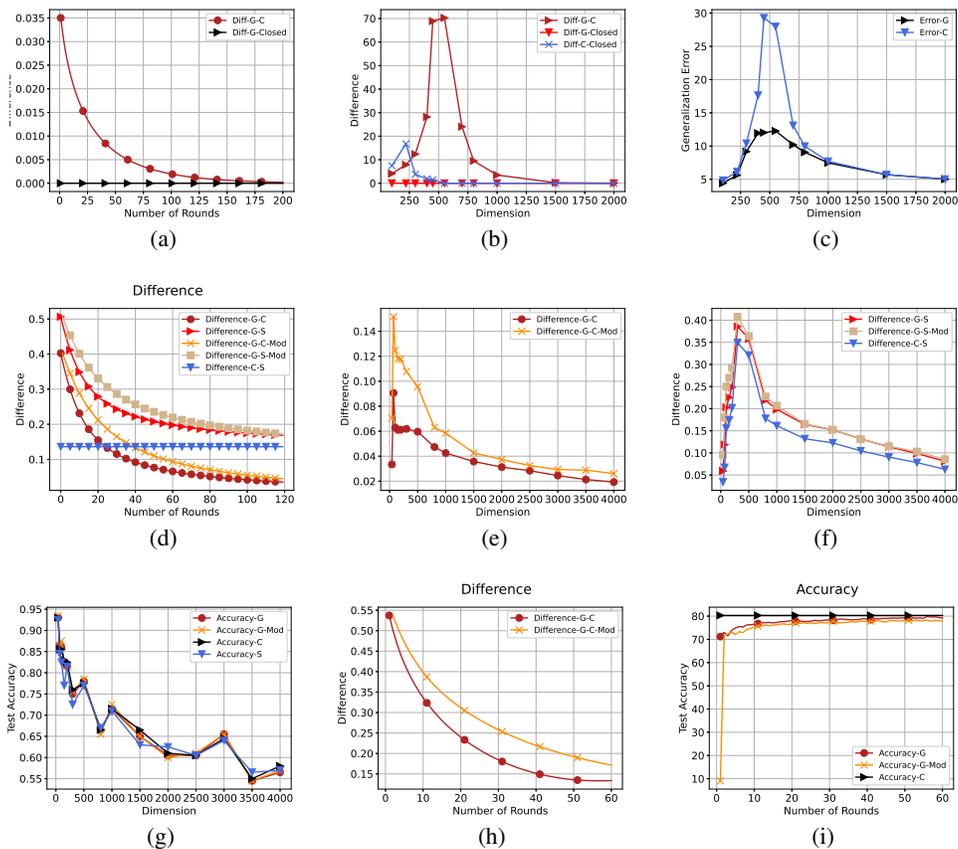


Figure 1: From left to right, from up to bottom (LR: Linear Regression, LC: Linear Classification, NN: Neural Network): (a) Difference between models with communication rounds in LR. (b) Difference between models with dimension in LR. (c) Generalization error with dimension in LR. (d) Difference between global model and centralized model with R in LC. (e) Difference between global model and centralized model with d in LC. (f) Difference from SVM model with d in LC. (g) Test Accuracy in LC. (h) Difference between global model and centralized model with communication rounds in NN. (i) Test accuracy with communication rounds in NN.

compute node, the half of local dataset belongs to the same class, and the other half consists of rest of 9 classes evenly, which forms a heterogeneous data distribution. The centralized model is trained with the whole CIFAR10 dataset. The models are trained with cross entropy loss and SGD. The learning rate is 0.01 and the batch size is 128. The number of local steps is $L = 60$ and number of communication rounds is $R = 60$. The centralized model is trained with the same learning rate for 3600 steps. We plot the difference between the linear layer and test accuracy with number of rounds in Fig. 1 (h) and (i). Again the difference is defined in direction. We can see the difference gradually decreases to a small error floor and the accuracy of global models and centralized model is very similar at last.

6 CONCLUSIONS

In this work we analyzed the implicit bias in distributed setting, and characterized the dynamics of global model trained from Local-GD with many local steps based on the implicit bias. We showed that the global model can converge to centralized model for both linear regression and classification tasks, providing a new perspective why Local-GD (FedAvg) works well in practice even with a large number of local steps on heterogeneous data. One potential future work is to extend the analysis of Local-GD to neural network using the developed implicit bias of deeper models Chizat & Bach (2020); Gunasekar et al. (2018b); Ji & Telgarsky (2019b); Kou et al. (2024).

-
- 540 REFERENCES
541
542 Heinz H Bauschke and Patrick L Combettes. *Convex Analysis and Monotone Operator Theory in*
543 *Hilbert Spaces*. Springer, 2011.
- 544 Mikhail Belkin, Siyuan Ma, and Soumik Mandal. To understand deep learning we need to understand
545 kernel learning. In *International Conference on Machine Learning*, pp. 541–549. PMLR, 2018.
- 546
547 Mikhail Belkin, Daniel Hsu, Siyuan Ma, and Soumik Mandal. Reconciling modern machine-learning
548 practice and the classical bias–variance trade-off. *Proceedings of the National Academy of Sciences*,
549 116(32):15849–15854, 2019.
- 550 Matias D Cattaneo, Jason M Klusowski, and Boris Shigida. On the implicit bias of adam. *arXiv*
551 *preprint arXiv:2309.00079*, 2023.
- 552
553 Zachary Charles, Zachary Garrett, Zhouyuan Huo, Sergei Shmulyian, and Virginia Smith. On
554 large-cohort training for federated learning. *Advances in neural information processing systems*,
555 34:20461–20475, 2021.
- 556 Hong-You Chen and Wei-Lun Chao. Fedbe: Making bayesian model ensemble applicable to federated
557 learning. In *International Conference on Learning Representations*, 2021.
- 558
559 Lenaic Chizat and Francis Bach. Implicit bias of gradient descent for wide two-layer neural networks
560 trained with the logistic loss. In *Conference on learning theory*, pp. 1305–1338. PMLR, 2020.
- 561 Liam Collins, Hamed Hassani, Aryan Mokhtari, and Sanjay Shakkottai. Fedavg with fine tuning:
562 Local updates lead to representation learning. *Advances in Neural Information Processing Systems*,
563 35:10572–10586, 2022.
- 564
565 Patrick L Combettes. Inconsistent signal feasibility problems: Least-squares solutions in a product
566 space. *IEEE Transactions on Signal Processing*, 42(11):2955–2966, 1994.
- 567
568 Patrick L Combettes. The convex feasibility problem in image recovery. In *Advances in imaging and*
569 *electron physics*, volume 95, pp. 155–270. Elsevier, 1996.
- 570
571 Patrick L Combettes. Convex set theoretic image recovery by extrapolated iterations of parallel
572 subgradient projections. *IEEE Transactions on Image Processing*, 6(4):493–506, 1997.
- 573
574 Michael Crawshaw, Yajie Bao, and Mingrui Liu. Federated learning with client subsampling, data
575 heterogeneity, and unbounded smoothness: A new algorithm and lower bounds. *Advances in Neural*
576 *Information Processing Systems*, 36, 2023.
- 577
578 Alvaro Rodolfo de Pierro and Alfredo Noel Iusem. *A parallel projection method of finding a common*
579 *point of a family of convex sets*. Inst. de matemática pura e aplicada, Conselho nacional de
580 desenvolvimento . . . , 1984.
- 581
582 Yuyang Deng, Mohammad Mahdi Kamani, and Mehrdad Mahdavi. Local sgd optimizes overparam-
583 eterized neural networks in polynomial time. In *International Conference on Artificial Intelligence*
584 *and Statistics*, pp. 6840–6861. PMLR, 2022.
- 585
586 Jacob Devlin. Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv*
587 *preprint arXiv:1810.04805*, 2018.
- 588
589 Jeff Donahue, Yangqing Jia, Oriol Vinyals, Judy Hoffman, Ning Zhang, Eric Tzeng, and Trevor Darrell.
590 Decaf: A deep convolutional activation feature for generic visual recognition. In *International*
591 *conference on machine learning*, pp. 647–655. PMLR, 2014.
- 592
593 Arthur Douillard, Qixuan Feng, Andrei A Rusu, Rachita Chhaparia, Yani Donchev, Adhiguna Kuncoro,
594 Marc’ Aurelio Ranzato, Arthur Szlam, and Jiajun Shen. Diloco: Distributed low-communication
595 training of language models. *arXiv preprint arXiv:2311.08105*, 2023.
- 596
597 Itay Evron, Edward Moroshko, Rachel Ward, Nathan Srebro, and Daniel Soudry. How catastrophic can
598 catastrophic forgetting be in linear regression? In *Conference on Learning Theory*, pp. 4028–4079.
599 PMLR, 2022.

594 Itay Evron, Edward Moroshko, Gon Buzaglo, Maroun Khriesh, Badea Marjeh, Nathan Srebro,
595 and Daniel Soudry. Continual learning in linear classification on separable data. *arXiv preprint*
596 *arXiv:2306.03534*, 2023.

597 Spencer Frei, Gal Vardi, Peter Bartlett, and Nati Srebro. The double-edged sword of implicit bias:
598 Generalization vs. robustness in relu networks. *Advances in Neural Information Processing Systems*,
599 36, 2024.

600 Peter Gilbert. Iterative methods for the three-dimensional reconstruction of an object from projections.
601 *Journal of theoretical biology*, 36(1):105–117, 1972.

602 Daniel Goldfarb and Paul Hand. Analysis of catastrophic forgetting for random orthogonal
603 transformation tasks in the overparameterized regime. In *International Conference on Artificial*
604 *Intelligence and Statistics*, pp. 2975–2993. PMLR, 2023.

605 Suriya Gunasekar, Jason Lee, Daniel Soudry, and Nathan Srebro. Characterizing implicit bias in
606 terms of optimization geometry. In *International Conference on Machine Learning*, pp. 1832–1841.
607 PMLR, 2018a.

608 Suriya Gunasekar, Jason D Lee, Daniel Soudry, and Nati Srebro. Implicit bias of gradient descent
609 on linear convolutional networks. *Advances in neural information processing systems*, 31, 2018b.

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

612 Baihe Huang, Xiaoxiao Li, Zhao Song, and Xin Yang. Fl-ntk: A neural tangent kernel-based
613 framework for federated learning analysis. In *International Conference on Machine Learning*, pp.
614 4423–4434. PMLR, 2021.

615 Yanping Huang, Youlong Cheng, Ankur Bapna, Orhan Firat, Dehao Chen, Mia Chen, HyoukJoong
616 Lee, Jiquan Ngiam, Quoc V Le, Yonghui Wu, et al. Gpipe: Efficient training of giant neural
617 networks using pipeline parallelism. *Advances in neural information processing systems*, 32, 2019.

618 Sami Jaghouar, Jack Min Ong, and Johannes Hagemann. Opendiloco: An open-source framework
619 for globally distributed low-communication training. *arXiv preprint arXiv:2407.07852*, 2024.

620 Divyansh Jhunjhunwala, Shiqiang Wang, and Gauri Joshi. Fedexp: Speeding up federated averaging
621 via extrapolation. In *The Eleventh International Conference on Learning Representations*, 2023.

622 Ziwei Ji and Matus Telgarsky. The implicit bias of gradient descent on nonseparable data. In
623 *Conference on Learning Theory*, pp. 1772–1798. PMLR, 2019a.

624 Ziwei Ji and Matus Telgarsky. Gradient descent aligns the layers of deep linear networks. In
625 *International Conference on Learning Representations*, 2019b.

626 Haoming Jiang, Pengcheng He, Weizhu Chen, Xiaodong Liu, Jianfeng Gao, and Tuo Zhao. Smart:
627 Robust and efficient fine-tuning for pre-trained natural language models through principled
628 regularized optimization. In *Proceedings of the 58th Annual Meeting of the Association for*
629 *Computational Linguistics*, pp. 2177–2190, 2020.

630 Peter Kairouz, H Brendan McMahan, Brendan Avent, Aurélien Bellet, Mehdi Bennis, Arjun Nitin
631 Bhagoji, Kallista Bonawitz, Zachary Charles, Graham Cormode, Rachel Cummings, et al. Advances
632 and open problems in federated learning. *arXiv preprint arXiv:1912.04977*, 2019.

633 Sai Praneeth Karimireddy, Satyen Kale, Mehryar Mohri, Sashank Reddi, Sebastian Stich, and
634 Ananda Theertha Suresh. Scaffold: Stochastic controlled averaging for federated learning. In
635 *International Conference on Machine Learning*, pp. 5132–5143. PMLR, 2020.

636 Ahmed Khaled, Konstantin Mishchenko, and Peter Richtárik. Tighter theory for local sgd on identical
637 and heterogeneous data. In *International Conference on Artificial Intelligence and Statistics*, pp.
638 4519–4529. PMLR, 2020.

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

-
- 648 Simon Kornblith, Jonathon Shlens, and Quoc V Le. Do better imagenet models transfer better?
649 In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp.
650 2661–2671, 2019.
- 651 Yiwen Kou, Zixiang Chen, and Quanquan Gu. Implicit bias of gradient descent for two-layer relu
652 and leaky relu networks on nearly-orthogonal data. *Advances in Neural Information Processing*
653 *Systems*, 36, 2024.
- 654 Tian Li, Anit Kumar Sahu, Manzil Zaheer, Maziar Sanjabi, Ameet Talwalkar, and Virginia Smith.
655 Federated optimization in heterogeneous networks. *Proceedings of Machine learning and systems*,
656 2:429–450, 2020a.
- 657 Tian Li, Shengyuan Hu, Ahmad Beirami, and Virginia Smith. Ditto: Fair and robust federated learning
658 through personalization. In *International Conference on Machine Learning*, pp. 6357–6368. PMLR,
659 2021.
- 660 Xiang Li, Kaixuan Huang, Wenhao Yang, Shusen Wang, and Zhihua Zhang. On the convergence
661 of FedAvg on non-IID data. In *International Conference on Learning Representations*, 2020b.
- 662 Sen Lin, Peizhong Ju, Yingbin Liang, and Ness Shroff. Theory on forgetting and generalization of con-
663 tinual learning. In *International Conference on Machine Learning*, pp. 21078–21100. PMLR, 2023.
- 664 Tao Lin, Sebastian U Stich, Kumar Kshitij Patel, and Martin Jaggi. Don’t use large mini-batches, use
665 local SGD. In *International Conference on Learning Representations*, 2019.
- 666 Brendan McMahan, Eider Moore, Daniel Ramage, Seth Hampson, and Blaise Aguera y Arcas.
667 Communication-efficient learning of deep networks from decentralized data. In *Artificial*
668 *intelligence and statistics*, pp. 1273–1282. PMLR, 2017.
- 669 Mor Shpigel Nacson, Nathan Srebro, and Daniel Soudry. Stochastic gradient descent on separable data:
670 Exact convergence with a fixed learning rate. In *The 22nd International Conference on Artificial*
671 *Intelligence and Statistics*, pp. 3051–3059. PMLR, 2019.
- 672 Tiancheng Qin, S Rasoul Etesami, and César A Uribe. Faster convergence of local sgd for
673 over-parameterized models. *arXiv preprint arXiv:2201.12719*, 2022.
- 674 Sashank J Reddi, Zachary Charles, Manzil Zaheer, Zachary Garrett, Keith Rush, Jakub Konečný,
675 Sanjiv Kumar, and Hugh Brendan McMahan. Adaptive federated optimization. In *International*
676 *Conference on Learning Representations*, 2021.
- 677 Hamza Reguieg, Mohammed El Hanjri, Mohamed El Kamili, and Abdellatif Kobbane. A comparative
678 evaluation of fedavg and per-fedavg algorithms for dirichlet distributed heterogeneous data. In *2023*
679 *10th International Conference on Wireless Networks and Mobile Communications (WINCOM)*,
680 pp. 1–6. IEEE, 2023.
- 681 Alexander Sergeev and Mike Del Balso. Horovod: fast and easy distributed deep learning in tensorflow.
682 *arXiv preprint arXiv:1802.05799*, 2018.
- 683 Bingqing Song, Prashant Khanduri, Xinwei Zhang, Jinfeng Yi, and Mingyi Hong. Fedavg converges
684 to zero training loss linearly for overparameterized multi-layer neural networks. In *International*
685 *Conference on Machine Learning*, pp. 32304–32330. PMLR, 2023.
- 686 Daniel Soudry, Elad Hoffer, Mor Shpigel Nacson, Suriya Gunasekar, and Nathan Srebro. The implicit
687 bias of gradient descent on separable data. *The Journal of Machine Learning Research*, 19(1):
688 2822–2878, 2018.
- 689 Sebastian U Stich. Local SGD converges fast and communicates little. In *International Conference*
690 *on Learning Representations*, 2019.
- 691 Canh T Dinh, Nguyen Tran, and Josh Nguyen. Personalized federated learning with moreau envelopes.
692 *Advances in Neural Information Processing Systems*, 33:21394–21405, 2020.
- 693 Jianyu Wang, Qinghua Liu, Hao Liang, Gauri Joshi, and H Vincent Poor. Tackling the objective
694 inconsistency problem in heterogeneous federated optimization. *Advances in neural information*
695 *processing systems*, 33:7611–7623, 2020.

702 Jianyu Wang, Rudrajit Das, Gauri Joshi, Satyen Kale, Zheng Xu, and Tong Zhang. On the unreasonable
703 effectiveness of federated averaging with heterogeneous data. *Transactions on Machine Learning*
704 *Research*, 2024.

705 Feijie Wu, Song Guo, Zhihao Qu, Shiqi He, Ziming Liu, and Jing Gao. Anchor sampling for federated
706 learning with partial client participation. In *International Conference on Machine Learning*, pp.
707 37379–37416. PMLR, 2023.

708 Shuo Xie and Zhiyuan Li. Implicit bias of adamw: ℓ_∞ -norm constrained optimization. In *International*
709 *Conference on Machine Learning*, pp. 54488–54510. PMLR, 2024.

710 Hao Yu, Sen Yang, and Shenghuo Zhu. Parallel restarted sgd with faster convergence and less
711 communication: Demystifying why model averaging works for deep learning. In *Proceedings of*
712 *the AAAI conference on artificial intelligence*, pp. 5693–5700, 2019.

713
714
715
716
717
718
719
720
721
722
723
724
725
726
727
728
729
730
731
732
733
734
735
736
737
738
739
740
741
742
743
744
745
746
747
748
749
750
751
752
753
754
755

756 A RELATED WORK

757
758 **Convergence of Local-GD.** When data distribution is homogeneous, many works have been done
759 to establish convergence analysis for Local (Stochastic) GD Stich (2019); Yu et al. (2019); Khaled et al.
760 (2020). With a “properly” small number of local steps, the dominating convergence rate is not affected.
761 Further various assumptions have been made to handle data heterogeneity and develop convergence
762 analysis Li et al. (2020b); Karimireddy et al. (2020); Khaled et al. (2020); Reddi et al. (2021); Wang
763 et al. (2020); Crawshaw et al. (2023). For strongly convex and smooth loss functions, the number
764 of local steps should not be larger than $O(\sqrt{T})$ for i.i.d data Stich (2019) and non-i.i.d. data Li et al.
765 (2020b). However, in practice Local-GD (FedAvg) works well in many applications McMahan et al.
766 (2017); Charles et al. (2021), even in training large language models Douillard et al. (2023); Jaghour
767 et al. (2024). In Wang et al. (2024), the authors argue that the previous theoretical assumption does not
768 align with practice and proposed a client consensus hypothesis to explain the effectiveness of FedAvg
769 in heterogeneous data. But they do not consider the impact of overparameterization on distributed
770 training. There are some works incorporating the property of zero training loss of overparameterized
771 neural networks into the conventional convergence analysis of FedAvg Huang et al. (2021); Deng
772 et al. (2022); Song et al. (2023); Qin et al. (2022). However, they do not guarantee which point FedAvg
773 can converge to. Our work is different from these works as: 1. We analyze which point the Local-GD
774 can converge to, which is a more elementary problem before obtaining the convergence rate; 2. We
775 use implicit bias as a technical tool to analyze the overparameterized FL.

776 **Implicit Bias.** Soudry et al. (2018) is the first work to show the gradient descent converges to
777 a max-margin direction on linearly separable data with a linear model and exponentially-tailed
778 loss function. Ji & Telgarsky (2019a) has provided an alternative analysis and extended this to
779 non-separable data. The theory of implicit bias has been further developed, for example, for wide
780 two-layer neural networks Chizat & Bach (2020), deep linear models Ji & Telgarsky (2019b), linear
781 convolutional networks Gunasekar et al. (2018b), two-layer ReLU networks Kou et al. (2024) etc.
782 Beyond gradient descent, more algorithms have been considered, including gradient descent with
783 momentum Gunasekar et al. (2018a), SGD Nacson et al. (2019), Adam Cattaneo et al. (2023), AdamW
784 Xie & Li (2024). Recently, implicit bias has also been used to characterize the dynamics of continual
785 learning, on linear regression Evron et al. (2022); Goldfarb & Hand (2023); Lin et al. (2023), and linear
786 classification Evron et al. (2023). In Evron et al. (2023), gradient descent on continually learned tasks
787 is related to Projections onto Convex Sets (POCS) and shown to converge to a *sequential* max-margin
788 scheme. In our work we consider the implicit bias of gradient descent in distributed setting, which
789 is related to a different parallel projection scheme by projecting onto constraint sets *simultaneously*.

790 **Parallel Projection.** Parallel projection methods are a family of algorithms to find a common point
791 across multiple constraint sets by projecting onto these sets in parallel. These methods are widely
792 used in feasibility problems in signal processing and image reconstruction Bauschke & Combettes
793 (2011). The straightforward average of multiple projections is known as the simultaneous iterative
794 reconstruction technique (SIRT) in Gilbert (1972). Then de Pierro & Iusem (1984) studied the
795 convergence of PPM for a relaxed version, and Combettes (1994) further generalized the result to
796 inconsistent feasibility problems. In Combettes (1997), an extrapolated parallel projection method was
797 proposed to accelerate the convergence. We note that Jhunjhunwala et al. (2023) used this extrapolation
798 to accelerate FedAvg. However, it was just inspired by the similarity between parallel projection
799 method and FedAvg, while in this work we rigorously prove the relation between PPM and FedAvg
800 using implicit bias of gradient descent.

800 B ADDITIONAL EXPERIMENTS

801 B.1 LINEAR CLASSIFICATION WITH DIRICHLET DISTRIBUTION

802
803
804 In federated learning, the Dirichlet distribution is usually used to generate heterogeneous datasets
805 across the compute nodes Hsu et al. (2019); Chen & Chao (2021); Reguieg et al. (2023). For binary
806 classification problem, the Dirichlet distribution $\text{Dir}(\alpha)$ is used to unbalance the positive and negative
807 samples. In the experiments we have 10 compute nodes. We generate 500 samples as $y_i = \text{sign}(x_i^T w^*)$
808 for $i \in [500]$ and use $\text{Dir}(\alpha)$ to distribute the 500 samples across 10 compute nodes. Note that the
809 number of samples at each compute node is not necessarily identical. Fig. 2 shows performance of
Local-GD for linear classification with different parameter α in Dirichlet distribution. The λ is set to

810
811
812
813
814
815
816
817
818
819
820
821
822
823
824
825
826
827
828
829
830
831
832
833
834
835
836
837
838
839
840
841
842
843
844
845
846
847
848
849
850
851
852
853
854
855
856
857
858
859
860
861
862
863

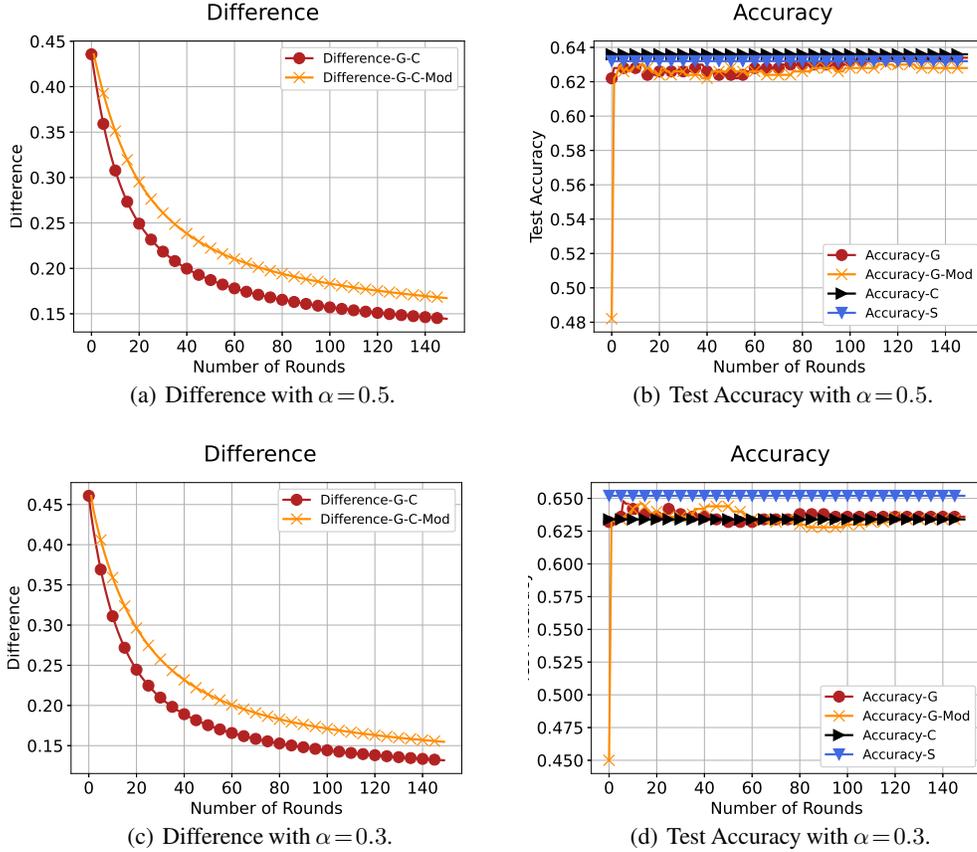


Figure 2: Local-GD on linear classification with Dirichlet distribution.

be 0.0001 and model dimension is fixed as $d = 1500$. The number of local steps L is 150 and number of communication rounds R is 150. The learning rate is 0.01. The centralized model is trained with the same learning rate for 22500 steps. We can see the global model and modified global model still converge to the centralized model in direction and get similar test accuracy.

C PROOFS IN SECTION 2

C.1 PROOF OF LEMMA 1

At each compute node, the local model converges to the solution of problem

$$\min_{w_i} \|w_i - w_0^k\|^2 \quad \text{s.t.} \quad X_i w_i = y_i. \quad (17)$$

Using Lagrange multipliers, we can write the Lagrangian as

$$\frac{1}{2} \|w_i - w_0^k\|^2 + \beta^T (X_i w_i - y_i) \quad (18)$$

Setting the derivative to 0, we know the optimal \tilde{w}_i satisfies

$$\tilde{w}_i - w_0^k + X_i^T \beta = 0, \quad (19)$$

and then

$$\tilde{w}_i = w_0^k - X_i^T \beta. \quad (20)$$

Also by the constraint $y_i = X_i \tilde{w}_i$, we can get

$$y_i = X_i w_0^k - (X_i X_i^T) \beta. \quad (21)$$

Since the model is overparameterized ($d > N$), $X_i X_i^T \in \mathbb{R}^{d \times d}$ is invertible. Then we have

$$\beta = -(X_i X_i^T)^{-1} (y_i - X_i w_0^k). \quad (22)$$

Plugging the β back, we can get the closed form solution as

$$\tilde{w}_i = w_0^k + X_i^T (X_i X_i^T)^{-1} (y_i - X_i w_0^k). \quad (23)$$

We update the local model $w_i^{k+1} = \tilde{w}_i$.

We can also write the closed form solution as

$$\begin{aligned} w_i^{k+1} &= w_0^k + X_i^T (X_i X_i^T)^{-1} (y_i - X_i w_0^k) \\ &= (I - X_i^T (X_i X_i^T)^{-1} X_i) w_0^k + X_i^T (X_i X_i^T)^{-1} y_i \end{aligned} \quad (24)$$

If we plug in the generative model $y_i = X_i w_i^* + z_i$, then the solution is

$$\begin{aligned} w_i^{k+1} &= (I - X_i^T (X_i X_i^T)^{-1} X_i) w_0^k + X_i^T (X_i X_i^T)^{-1} X_i w_i^* + X_i^T (X_i X_i^T)^{-1} z_i \\ &= (I - P_i) w_0^k + P_i w_i^* + X_i^\dagger z_i. \end{aligned} \quad (25)$$

where $P_i = X_i^T (X_i X_i^T)^{-1} X_i$ is the projection operator to the row space of X_i , and $X_i^\dagger = X_i^T (X_i X_i^T)^{-1}$ is the pseudo inverse of X_i . It is an interpolation between the initial global model w_0^k and the local true model w_i^* , plus a noise term.

After aggregating all the local models, the global model is

$$\begin{aligned} w_0^{k+1} &= \frac{1}{m} \sum_{i=1}^m (I - P_i) w_0^k + \frac{1}{m} \sum_{i=1}^m P_i w_i^* + \frac{1}{m} \sum_{i=1}^m X_i^\dagger z_i \\ &= (I - \bar{P}) w_0^k + \bar{Q} + \bar{Z}, \end{aligned} \quad (26)$$

where $\bar{P} = \frac{1}{m} \sum_{i=1}^m P_i$, $\bar{Q} = \sum_{i=1}^m P_i w_i^*$, $\bar{Z} = \frac{1}{m} \sum_{i=1}^m X_i^\dagger z_i$.

After K rounds of communication, the global model is

$$w_0^K = (I - \bar{P})^K w_0^0 + \sum_{k=0}^{K-1} (I - \bar{P}) (\bar{Q} + \bar{Z}). \quad (27)$$

If we start from $w_0^0 = 0$, then the solution will converge to $\sum_{k=0}^{K-1} (I - \bar{P}) (\bar{Q} + \bar{Z})$.

C.2 PROOF OF THEOREM 1

We know the difference between global model and centralized model is iteratively projected onto the null space of span of row spaces of X_i s:

$$w_0^{k+1} - w_c = (I - \bar{P}) (w_0^k - w_c). \quad (28)$$

We can formally describe it as follows. Since the problem is overparameterized globally, we can assume each X_i has full rank N . We apply singular value decomposition (SVD) to X_i as $X_i = U_i \Sigma_i V_i^T$, where $U_i \in \mathbb{R}^{N \times N}$, $V_i \in \mathbb{R}^{d \times N}$. Then $P_i = X_i^T (X_i X_i^T)^{-1} X_i = V_i V_i^T$, which is the projection matrix to the row space of X_i .

We apply eigenvalue decomposition on \bar{P} to get $\bar{P} = Q \Sigma Q^T$, where $Q \in \mathbb{R}^{d \times n'}$ and n' is the rank of \bar{P} . It satisfies $N \leq n' \leq MN$. Since \bar{P} is a linear combination of P_i s, the space of column space of Q is the space spanned by all the vectors $v_{ij}, i = 1, \dots, M, j = 1, \dots, N$.

We also construct a matrix $Q' \in \mathbb{R}^{d \times (d - n')}$, which consists of orthonormal vectors perpendicular to Q . We can project the difference onto column space of Q and Q' respectively.

$$\begin{aligned} Q^T (w_0^{k+1} - w_c) &= Q^T (I - Q \Sigma Q^T) (w_0^k - w_c) = (I - \Sigma) Q^T (w_0^k - w_c) \\ Q'^T (w_0^{k+1} - w_c) &= Q'^T (I - Q \Sigma Q^T) (w_0^k - w_c) = Q'^T (w_0^k - w_c) \end{aligned} \quad (29)$$

After K rounds of communication, we can decompose $w_0^K - w_c$ into two parts:

$$w_0^K - w_c = QQ^T(w_0^K - w_c) + Q'Q'^T(w_0^K - w_c). \quad (30)$$

Then we can obtain

$$\begin{aligned} w_0^K - w_c &= QQ^T(w_0^K - w_c) + Q'Q'^T(w_0^K - w_c) \\ &= Q(I - \Sigma)^K Q^T(w_0^0 - w_c) + Q'Q'^T(w_0^0 - w_c). \end{aligned}$$

It shows the initial difference on the column space of Q continues to decrease until zero if K is sufficiently large. And the initial difference on the null space of Q remains constant.

To show the difference $w_0^K - w_c$ goes to zero entirely, we just need to choose an initial point such that initial difference is on the column space of Q . When we choose $w_0^0 = 0$, the initial difference is w_c itself. Moreover, the centralized solution $w_c = X_c^T(X_c X_c^T)^{-1}y_c$ exactly lies in the data space spanned by vectors $\{v_{ij}\}_{i=1, j=1}^{M, N}$ since it is a linear combination of columns of X_c^T . So if we start from $w_0^0 = 0$, then $w_0^K - w_c$ will go to zero when K is sufficiently large.

D PROOFS IN SECTION 3

In the proofs of linear classification, for ease of notation, we redefine the samples $y_{ij}x_{ij}$ to x_{ij} to subsume the labels.

D.1 PROOFS OF LEMMA 2

We assume $\|w_0^k - \ln(\frac{1}{\lambda})\bar{w}_0^k\| = O(k \ln \ln \frac{1}{\lambda})$. In this case, since $\ln \frac{1}{\lambda}$ grows faster, when $\lambda \rightarrow 0$, we can have $\lim_{\lambda \rightarrow 0} \frac{w_0^k}{\|w_0^k\|} = \frac{\bar{w}_0^k}{\|\bar{w}_0^k\|}$ for any k at order $o\left(\frac{\ln(1/\lambda)}{\ln \ln(1/\lambda)}\right)$. We will prove it by induction. We define global and local residuals as $r^k = w_0^k - \ln(\frac{1}{\lambda})\bar{w}_0^k$ and $r_i^k = w_i^k - \ln(\frac{1}{\lambda})\bar{w}_i^k$.

When $k=0$, since $w_0^0 = \bar{w}_0^0 = 0$, $r_i^0 = 0$ and the assumption trivially holds.

When $k \geq 1$, we have

$$\begin{aligned} \|r^k\| &= \left\| w_0^k - \ln\left(\frac{1}{\lambda}\right)\bar{w}_0^k \right\| = \frac{1}{M} \left\| \sum_{i=1}^M w_i^k - \ln\left(\frac{1}{\lambda}\right)\bar{w}_i^k \right\| \\ &\leq \frac{1}{M} \sum_{i=1}^M \left\| w_i^k - \ln\left(\frac{1}{\lambda}\right)\bar{w}_i^k \right\| = \frac{1}{M} \sum_{i=1}^M \|r_i^k\|. \end{aligned} \quad (31)$$

where the inequality is triangle inequality. We then focus on the local residual r_i^k . We choose an $O(1)$ vector \tilde{w}_i^k and a sign $s_i^k \in \{-1, +1\}$ to show

$$\begin{aligned} \|r_i^k\| &= \left\| w_i^k - \left[\left(\ln\left(\frac{1}{\lambda}\right) + s_i^k \ln \ln\left(\frac{1}{\lambda}\right) \right) \bar{w}_i^k + \tilde{w}_i^k \right] + s_i^k \ln \ln\left(\frac{1}{\lambda}\right) \bar{w}_i^k + \tilde{w}_i^k \right\| \\ &\leq \left\| w_i^k - \left[\left(\ln\left(\frac{1}{\lambda}\right) + s_i^k \ln \ln\left(\frac{1}{\lambda}\right) \right) \bar{w}_i^k + \tilde{w}_i^k \right] \right\| + \ln \ln\left(\frac{1}{\lambda}\right) \|\bar{w}_i^k\| + \|\tilde{w}_i^k\| \end{aligned} \quad (32)$$

Recall the w_i^k is the solution of optimization problem

$$\operatorname{argmin}_{w_i} f_i(w_i) = \sum_{j=1}^N \exp(-x_{ij}^T w_i) + \frac{\lambda}{2} \|w_i - w_0^{k-1}\|^2, \quad (33)$$

and the loss function $f_i(w_i)$ is a λ -strongly convex function. Thus we have

$$\|w_i^k - w\| \leq \frac{1}{\lambda} \|\nabla f_i(w)\|, \quad \text{for any } w. \quad (34)$$

Then back to 32, we have

$$\|r_i^k\| \leq \frac{1}{\lambda} \left\| \underbrace{\nabla f_i \left[\left(\ln\left(\frac{1}{\lambda}\right) + s_i^k \ln \ln\left(\frac{1}{\lambda}\right) \right) \bar{w}_i^k + \tilde{w}_i^k \right]}_{\|A_i\|} \right\| + \ln \ln\left(\frac{1}{\lambda}\right) \|\bar{w}_i^k\| + \|\tilde{w}_i^k\|. \quad (35)$$

Next we need to show the first term A_i is at $O((k-1)\ln \ln(\frac{1}{\lambda}))$, and also since $\|\bar{w}_i^k\|$ and $\|\tilde{w}_i^k\|$ are $O(1)$ vectors, then $\|r_i^k\|$ is at order $O(k \ln \ln(\frac{1}{\lambda}))$. After averaging, $\|r^k\|$ is also at order $O(k \ln \ln(\frac{1}{\lambda}))$. This confirms the assumption made for induction.

Now we focus on the term A_i . The gradient of function $f_i(w)$ is

$$\nabla f_i(w_i) = \sum_j -x_{ij} \exp(-x_{ij}^T w_i) + \lambda(w_i - w_0^{k-1}). \quad (36)$$

The term A_i is

$$\begin{aligned} A_i &= \frac{1}{\lambda} \nabla f_i \left[\left(\ln\left(\frac{1}{\lambda}\right) + s_i^k \ln \ln\left(\frac{1}{\lambda}\right) \right) \bar{w}_i^k + \tilde{w}_i^k \right] \\ &= -\frac{1}{\lambda} \sum_j x_{ij} \exp\left(x_{ij}^T \ln\left(\lambda \ln^{-s_i^k}\left(\frac{1}{\lambda}\right)\right) \bar{w}_i^k\right) \exp(-x_{ij}^T \tilde{w}_i^k) + \left(\ln\left(\frac{1}{\lambda}\right) + s_i^k \ln \ln\left(\frac{1}{\lambda}\right)\right) \bar{w}_i^k + \tilde{w}_i^k - w_0^{k-1} \\ &= -\frac{1}{\lambda} \sum_j x_{ij} \left(\lambda \ln^{-s_i^k}\left(\frac{1}{\lambda}\right)\right)^{x_{ij}^T \bar{w}_i^k} \exp(-x_{ij}^T \tilde{w}_i^k) + \left(\ln\left(\frac{1}{\lambda}\right) + s_i^k \ln \ln\left(\frac{1}{\lambda}\right)\right) \bar{w}_i^k + \tilde{w}_i^k - w_0^{k-1}. \end{aligned} \quad (37)$$

Then we define the set of support vectors as $S_i^k = \{x_{ij} | x_{ij}^T \bar{w}_i^k = 1\}$. Recall that we assume $r^{k-1} = w_0^{k-1} - \ln(\frac{1}{\lambda}) \bar{w}_0^{k-1}$ is at order $O((k-1)\ln \ln(\frac{1}{\lambda}))$. We can obtain

$$\begin{aligned} A_i &= -\frac{1}{\lambda} \left(\lambda \ln^{-s_i^k}\left(\frac{1}{\lambda}\right)\right)^1 \sum_{x_{ij} \in S_i^k} x_{ij} \exp(-x_{ij}^T \tilde{w}_i^k) - \frac{1}{\lambda} \sum_{x_{ij} \notin S_i^k} x_{ij} \left(\lambda \ln^{-s_i^k}\left(\frac{1}{\lambda}\right)\right)^{x_{ij}^T \bar{w}_i^k} \exp(-x_{ij}^T \tilde{w}_i^k) \\ &\quad + \ln\left(\frac{1}{\lambda}\right) (\bar{w}_i^k - \bar{w}_0^{k-1}) - r^{k-1} + s_i^k \ln \ln\left(\frac{1}{\lambda}\right) \bar{w}_i^k + \tilde{w}_i^k \\ &= -\ln^{-s_i^k}\left(\frac{1}{\lambda}\right) \sum_{x_{ij} \in S_i^k} x_{ij} \exp(-x_{ij}^T \tilde{w}_i^k) - \sum_{x_{ij} \notin S_i^k} x_{ij} \lambda^{x_{ij}^T \bar{w}_i^k - 1} \left(\ln\left(\frac{1}{\lambda}\right)\right)^{-s_i^k x_{ij}^T \bar{w}_i^k} \exp(-x_{ij}^T \tilde{w}_i^k) \\ &\quad + \ln\left(\frac{1}{\lambda}\right) (\bar{w}_i^k - \bar{w}_0^{k-1}) - r^{k-1} + s_i^k \ln \ln\left(\frac{1}{\lambda}\right) \bar{w}_i^k + \tilde{w}_i^k. \end{aligned} \quad (38)$$

By the triangle inequality, we have

$$\begin{aligned} \|A_i\| &\leq \underbrace{\left\| \ln\left(\frac{1}{\lambda}\right) (\bar{w}_i^k - \bar{w}_0^{k-1}) - \ln^{-s_i^k}\left(\frac{1}{\lambda}\right) \sum_{x_{ij} \in S_i^k} x_{ij} \exp(-x_{ij}^T \tilde{w}_i^k) \right\|}_{B_1} \\ &\quad + \underbrace{\left\| \sum_{x_{ij} \notin S_i^k} x_{ij} \lambda^{x_{ij}^T \bar{w}_i^k - 1} \left(\ln\left(\frac{1}{\lambda}\right)\right)^{-s_i^k x_{ij}^T \bar{w}_i^k} \exp(-x_{ij}^T \tilde{w}_i^k) \right\|}_{B_2} \\ &\quad + \underbrace{\|r^{k-1}\|}_{O((k-1)\ln \ln(\frac{1}{\lambda}))} + \underbrace{\ln \ln\left(\frac{1}{\lambda}\right) \|\bar{w}_i^k\|}_{O(1)} + \underbrace{\|\tilde{w}_i^k\|}_{O(1)}. \end{aligned} \quad (39)$$

We just need to show B_1 and B_2 approach to 0 then $\|A_i\|$ can approach to $O(k \ln \ln(\frac{1}{\lambda}))$.

We divide it into two cases.

1. When $\bar{w}_i^k = P(\bar{w}_0^{k-1}) \neq \bar{w}_0^{k-1}$, meaning \bar{w}_0^{k-1} is not in the convex set C_i . In this case we choose $s_i^k = -1$ then

$$\begin{aligned} B_1 &= \left\| \ln\left(\frac{1}{\lambda}\right)(\bar{w}_i^k - \bar{w}_0^{k-1}) - \ln\left(\frac{1}{\lambda}\right) \sum_{x_{ij} \in S_i^k} x_{ij} \exp(-x_{ij}^T \bar{w}_i^k) \right\| \\ &= \ln\left(\frac{1}{\lambda}\right) \left\| (\bar{w}_i^k - \bar{w}_0^{k-1}) - \sum_{x_{ij} \in S_i^k} x_{ij} \exp(-x_{ij}^T \bar{w}_i^k) \right\|. \end{aligned} \quad (40)$$

We now want to choose \tilde{w}_i^k to make B_1 as 0. Since \bar{w}_i^k is the solution of SVM problem (13), by the KKT condition of SVM problem, it can be written as

$$\bar{w}_i^k = \bar{w}_0^{k-1} + \sum_{x_{ij} \in S_i^k} \beta_{ij} x_{ij} \quad (41)$$

where β_{ij} is the dual variable corresponding to x_{ij} in the set of support vectors. Thus we want to choose \tilde{w}_i^k as

$$\sum_{x_{ij} \in S_i^k} \exp(-x_{ij}^T \tilde{w}_i^k) x_{ij} = \sum_{x_{ij} \in S_i^k} \beta_{ij} x_{ij}. \quad (42)$$

We can prove such a \tilde{w}_i^k almost surely exists in Lemma 5.

For the term B_2 , since $\lim_{\lambda \rightarrow 0} \lambda^{c-1} \ln^c\left(\frac{1}{\lambda}\right) \rightarrow 0$ for any constant $c > 1$, and $x_{ij}^T \bar{w}_i^k - 1 > 0$ for any x_{ij} being not a support vector, then we can see

$$B_2 = \left\| \sum_{x_{ij} \notin S_i^k} x_{ij} \lambda^{x_{ij}^T \bar{w}_i^k - 1} \left(\ln\left(\frac{1}{\lambda}\right)\right)^{x_{ij}^T \bar{w}_i^k} \exp(-x_{ij}^T \bar{w}_i^k) \right\| \xrightarrow{\lambda \rightarrow 0} 0. \quad (43)$$

Here we choose \tilde{w}_i^k and s_i^k to make $B_1 = 0$ and $B_2 \rightarrow 0$.

2. When $\bar{w}_i^k = P(\bar{w}_0^{k-1}) = \bar{w}_0^{k-1}$, meaning \bar{w}_0^{k-1} is already in the convex set C_i . Then $\bar{w}_i^k - \bar{w}_0^{k-1} = 0$. In this case we choose $\tilde{w}_i^k = 0$ and $s_i^k = +1$. We can have

$$B_1 = \ln^{-1}\left(\frac{1}{\lambda}\right) \left\| \sum_{x_{ij} \in S_i^k} x_{ij} \right\| \xrightarrow{\lambda \rightarrow 0} 0, \quad (44)$$

since $\ln^{-1}\left(\frac{1}{\lambda}\right) \xrightarrow{\lambda \rightarrow 0} 0$ and $\left\| \sum_{x_{ij} \in S_i^k} x_{ij} \right\|$ is $O(1)$.

And since $x_{ij}^T \bar{w}_i^k - 1 > 0$ for any x_{ij} being not a support vector, we have

$$B_2 = \left\| \sum_{x_{ij} \notin S_i^k} x_{ij} \lambda^{x_{ij}^T \bar{w}_i^k - 1} \left(\ln\left(\frac{1}{\lambda}\right)\right)^{-x_{ij}^T \bar{w}_i^k} \right\| \xrightarrow{\lambda \rightarrow 0} 0, \quad (45)$$

where $\lambda^{x_{ij}^T \bar{w}_i^k - 1} \xrightarrow{\lambda \rightarrow 0} 0$ and $\left(\ln\left(\frac{1}{\lambda}\right)\right)^{-x_{ij}^T \bar{w}_i^k} \xrightarrow{\lambda \rightarrow 0} 0$. Thus we choose \tilde{w}_i^k and s_i^k to make $B_1 \rightarrow 0$ and $B_2 \rightarrow 0$.

Plugging 39 back into 35, we can obtain

$$\begin{aligned} \|r_i^k\| &\leq \|A_i^k\| + \ln \ln\left(\frac{1}{\lambda}\right) \|\bar{w}_i^k\| + \|\bar{w}_i^k\| \\ &\leq \underbrace{B_1 + B_2}_{\rightarrow 0} + 2 \ln \ln\left(\frac{1}{\lambda}\right) \|\bar{w}_i^k\| + 2 \|\tilde{w}_i^k\| + \|r^{k-1}\| \\ &\leq 2 \ln \ln\left(\frac{1}{\lambda}\right) \|\bar{w}_i^k\| + 2 \|\tilde{w}_i^k\| + \|r^{k-1}\|. \end{aligned} \quad (46)$$

By the assumption $\|r^{k-1}\| = O((k-1) \ln \ln(\frac{1}{\lambda}))$ and $\|\bar{w}_i^k\| = O(1)$, $\|\tilde{w}_i^k\| = O(1)$, we have $\|r_i^k\| = O(k \ln \ln(\frac{1}{\lambda}))$.

From 31, we finally obtain

$$\|r^k\| \leq \frac{1}{M} \|r_i^k\| = O(k \ln \ln(\frac{1}{\lambda})), \quad (47)$$

which confirms our assumption. Then we have $\lim_{\lambda \rightarrow 0} \frac{w_0^k}{\|w_0^k\|} = \frac{\bar{w}_0^k}{\|\bar{w}_0^k\|}$ for any k at order $O\left(\frac{\ln(1/\lambda)}{\ln \ln(1/\lambda)}\right)$.

D.2 PROOFS OF AUXILIARY LEMMAS

Lemma 5. *For the sequence $\{\bar{w}_0^k\}$ generated by sequential SVM problems 13 and aggregations, and for almost all datasets sampled from M continuous distributions, the unique dual solution $\beta_i^k \in \mathbb{R}^{|S_i| \times 1}$ satisfying the KKT conditions of SVM problem 13 has non-zero elements. Then there exists \tilde{w}_i^k satisfying $X_{S_i} \tilde{w}_i^k = -\ln \beta_i^k$.*

For almost all datasets, a hyperplane can be determined by d points. Thus there are at most d support vectors and the set of support vectors is linearly independent.

Proof. By the KKT condition of SVM problem, we can write the solution as

$$\bar{w}_i^k = \bar{w}_0^{k-1} + \sum_{x_{ij} \in S_i} \beta_{ij}^k x_{ij} = \bar{w}_0^{k-1} + X_{S_i}^T \beta_i^k. \quad (48)$$

where $X_{S_i} \in \mathbb{R}^{|S_i| \times d}$ is the data matrix with all the support vectors, and $\beta_i^k \in \mathbb{R}^{|S_i| \times 1}$ is the dual variable vector. Thus we can obtain

$$\beta_i^k = (X_{S_i} X_{S_i}^T)^{-1} X_{S_i} (\bar{w}_i^k - \bar{w}_0^{k-1}) = (X_{S_i} X_{S_i}^T)^{-1} \mathbf{1}_{S_i} - (X_{S_i} X_{S_i}^T)^{-1} X_{S_i} \bar{w}_0^{k-1}, \quad (49)$$

where $X_{S_i} X_{S_i}^T$ is invertible since X_{S_i} has full row rank $|S_i|$, and the second equality is from $X_{S_i} \bar{w}_i^k = \mathbf{1}_{S_i}$ with $\mathbf{1}_{S_i} \in \mathbb{R}^{|S_i| \times 1}$ being all one vector. Plugging β_i^k back, we have

$$\bar{w}_i^k = \left[I - X_{S_i}^T (X_{S_i} X_{S_i}^T)^{-1} X_{S_i} \right] \bar{w}_0^{k-1} + X_{S_i}^T (X_{S_i} X_{S_i}^T)^{-1} \mathbf{1}_{S_i}. \quad (50)$$

After averaging, the global model is

$$\bar{w}_0^k = \left[I - \frac{1}{M} \sum_{i=1}^M X_{S_i}^T (X_{S_i} X_{S_i}^T)^{-1} X_{S_i} \right] \bar{w}_0^{k-1} + \frac{1}{M} \sum_{i=1}^M X_{S_i}^T (X_{S_i} X_{S_i}^T)^{-1} \mathbf{1}_{S_i}. \quad (51)$$

It implies \bar{w}_0^k is a rational function in the components of X_1, X_2, \dots, X_M , and also β_i^k is also a rational function in the components of data matrices. So its entries can be expressed as $\beta_{ij}^k = p_{ij}^k(X_1, X_2, \dots, X_M) / q_{ij}^k(X_1, X_2, \dots, X_M)$ for some polynomials p_{ij}^k, q_{ij}^k . Note that $\beta_{ij}^k = 0$ only if $p_{ij}^k(X_1, X_2, \dots, X_M) = 0$, and the components of X_1, X_2, \dots, X_M must constitute a root of polynomial p_{ij}^k . However, the root of any polynomial has measure zero, unless the polynomial is the zero polynomial, i.e., $p_{ij}^k(X_1, X_2, \dots, X_M) = 0$ for any X_1, X_2, \dots, X_M .

Next we need to show p_{ij}^k cannot be zero polynomials. To do this, we just need to construct a specific X_1, X_2, \dots, X_M where the p_{ij}^k is not zero polynomial. Denote $e_i \in \mathbb{R}^d$ as the i -th standard unit vector, and v_1, v_2, \dots, v_M be the number of support vectors at M compute nodes. We construct the datasets as

$$X_i = r_i [e_1, e_2, \dots, e_{v_i}]^T, \text{ for all } i. \quad (52)$$

where r_i are positive constants that will be chosen later. For these datasets, the set of support vector is dataset itself, i.e., $X_{S_i} = X_i$. We can calculate

$$X_i X_i^T = r_i^2 I_{v_i}, \quad X_i^T X_i = r_i^2 \begin{bmatrix} I_{v_i} & \mathbf{0} \\ \mathbf{0} & \mathbf{0}_{(d-v_i) \times (d-v_i)} \end{bmatrix}, \quad X_i^T \mathbf{1}_{S_i} = r_i \begin{bmatrix} \mathbf{1}_{v_i} \\ \mathbf{0}_{d-v_i} \end{bmatrix} \quad (53)$$

1188 because the maximum value of a_j is $\frac{M-1}{M}$ and the maximum value of b_j is $\frac{1}{M} \sum_{i=1}^M \frac{1}{r_i^2}$.

1189 Thus we require

$$1191 \sum_{i=1}^M \frac{1}{r_i} < \frac{1}{1 - \left(\frac{M-1}{M}\right)^{k-1}}. \quad (60)$$

1192 Since $\left(\frac{M-1}{M}\right)^{k-1} \rightarrow 0$ when $k \rightarrow \infty$, we only require the left-hand side is less than the lower bound
1193 of right-hand side:

$$1194 \sum_{i=1}^M \frac{1}{r_i} < 1. \quad (61)$$

1195 Therefore we can choose $r_i = M + 1$ to make it happen.

1196 Then we can obtain $\beta_{ij}^k > 0$ holds for any support vector x_{ij} and any round k . And the \tilde{w}_i^k simply
1197 satisfies $X_{S_i} \tilde{w}_i^k = -\ln \beta_i^k$. \square

1205 E LEMMA AND PROOFS IN SECTION 4

1206 Here we provide a lemma of Modified Local-GD similar to Lemma 2 of vanilla Local-GD.

1207 **Lemma 6.** *For almost all datasets sampled from a continuous distribution satisfying Assumption 1,*
1208 *we train the global model w_0 from Modified Local-GD in Algorithm 3 and \bar{w}_0 from Modified PPM. The*
1209 *parameter is chosen as $\alpha^k = 1 - \frac{1}{k+1}$. With initialization $w_0^0 = \bar{w}_0^0 = 0$, we have $w_0^k \rightarrow \ln\left(\frac{1}{\lambda}\right) \bar{w}_0^k$, and*
1210 *the residual $\|w_0^k - \ln\left(\frac{1}{\lambda}\right) \bar{w}_0^k\| = O(k \ln \ln \frac{1}{\lambda})$, as $\lambda \rightarrow 0$. It implies that at any round $k = o\left(\frac{\ln(1/\lambda)}{\ln \ln(1/\lambda)}\right)$,*
1211 *w_0^k converges in direction to \bar{w}_0^k :*

$$1212 \lim_{\lambda \rightarrow 0} \frac{w_0^k}{\|w_0^k\|} = \frac{\bar{w}_0^k}{\|\bar{w}_0^k\|}. \quad (62)$$

1213 *Proof.* With initialization $w_0^0 = \bar{w}_0^0 = 0$, the Modified Local-GD is just a scaling of vanilla Local-GD:

$$1214 w_0^{k+1} = \frac{k}{k+1} \frac{1}{M} \sum_{i=1}^M w_i^{k+1}. \quad (63)$$

1215 Also, the Modified PPM is a scaling of vanilla PPM: $\bar{w}_0^{k+1} = \frac{k}{k+1} \frac{1}{M} \sum_{i=1}^M \bar{w}_i^{k+1}$.

1216 When $k \geq 1$, we can know the residual between Modified Local-GD and Modified PPM is

$$1217 \begin{aligned} 1218 \|r^k\| &= \left\| w_0^k - \ln\left(\frac{1}{\lambda}\right) \bar{w}_0^k \right\| = \frac{k}{k+1} \frac{1}{M} \left\| \sum_{i=1}^M w_i^k - \ln\left(\frac{1}{\lambda}\right) \bar{w}_i^k \right\| \\ 1219 &\leq \frac{1}{M} \sum_{i=1}^M \left\| w_i^k - \ln\left(\frac{1}{\lambda}\right) \bar{w}_i^k \right\| = \frac{1}{M} \sum_{i=1}^M \|r_i^k\|. \end{aligned} \quad (64)$$

1220 Then we can follow the same process in the proof of Lemma 2 to obtain

$$1221 \|r^k\| \leq \frac{1}{M} \|r_i^k\| = O\left(k \ln \ln \left(\frac{1}{\lambda}\right)\right), \quad (65)$$

1222 As a result we have $\lim_{\lambda \rightarrow 0} \frac{w_0^k}{\|w_0^k\|} = \frac{\bar{w}_0^k}{\|\bar{w}_0^k\|}$.

1223 \square