

# Decentralized optimization with non-identical sampling in presence of stragglers

Tharindu Adikari and Stark Draper

## Abstract

We consider decentralized consensus optimization when workers sample data from non-identical distributions and perform variable amounts of work due to slow nodes known as stragglers. The problem of non-identical distributions and the problem of variable amount of work have been previously studied separately. In our work we analyze them together under a unified system model. We study the convergence of the optimization algorithm when combining worker outputs under two heuristic methods: (1) weighting equally, and (2) weighting by the amount of work completed by each. We prove convergence of the two methods under perfect consensus, assuming straggler statistics are independent and identical across all workers for all iterations. Our numerical results show that under approximate consensus the second method outperforms the first method for both convex and non-convex objective functions. We make use of the theory on minimum variance unbiased estimator (MVUE) to evaluate the existence of an optimal method for combining worker outputs. While we conclude that neither of the two heuristic methods are optimal, we also show that an optimal method does not exist.

## I. INTRODUCTION

The general system model for decentralized consensus optimization assumes a cluster of workers connected through a network. An important property is that there may not be a central coordinator amongst workers (unlike in master-worker systems). Fig. 1 presents an example of how 10 workers may be connected in a decentralized manner. Each worker is associated with a utility function and the goal is collaboratively and iteratively to optimize the sum of utility functions using only local computations. In each iteration, workers first perform local computations and then synchronize their results through some consensus mechanism. We refer to the two phases as ‘compute’ and ‘consensus’. Workers can arrive either at perfect or approximate consensus depending on the consensus scheme used. For example the analysis in [2] assumes perfect consensus which can be achieved using operations in the Message Passing Interface (MPI) standard. However, recently there has been significant interest in approximate averaging methods due to a number of attractive features. Such features include that these schemes are simple to implement, they support asynchronous operation, and they work well with dynamic network structures [3], [4].

A canonical application of decentralized consensus optimization is distributed machine learning. With large datasets, it is desirable to assign smaller subsets of the data to each worker, with workers collaborating to find an optimal model for the entire dataset. In this paper we study the global convergence of such a system that employs stochastic gradient descent (SGD) when data distributions at workers are non-identical *and* when workers perform variable amounts of work per iteration. We summarize the contributions of this paper with its outline as follows. In Sec. II we discuss the related work in this theme of studies. In Sec. III we outline the system model that we work with. Since workers perform variable amounts of work, the level of accuracy in worker outputs vary from one worker to the other. The worker outputs must be combined in a way that leads to the fastest convergence of the system. In Sec. IV we discuss two heuristic methods of combining the worker outputs. We numerically and theoretically analyze the

This work was supported by Huawei Technologies through a joint project with University of Toronto, and the Natural Science and Engineering Research Council (NSERC) of Canada through a Discovery Research Grant.

This paper was presented in part at the IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP), Barcelona, Spain, May 2020. ([1], <https://ieeexplore.ieee.org/document/9053329>)

The code used for the numerical experiments in this paper is available at <https://github.com/thadikari/consensus>.

The authors are with the Electrical and Computer Engineering Department, University of Toronto, Toronto, ON M5S 2E4, Canada (e-mail: tharindu.adikari@mail.utoronto.ca; stark.draper@utoronto.ca).

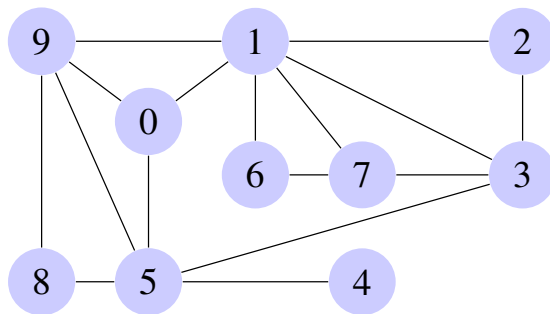


Fig. 1: Workers and the worker connections in a decentralized system. The topology is same as that in Figure 2 in [5].

performances of the two methods and prove their convergence with SGD. Subsequently, in Sec. V we discuss the existence of an optimal method of combining the worker outputs.

## II. RELATED WORK

There has been a significant amount of recent interest in the application of decentralized optimization to machine learning. A variety of system models have been considered [2], [3], [4], [5]. These studies differ in their assumptions of the network topology, the data model, the averaging method, and the underlying optimization algorithm. For example, the scheme in [2] employs exact averaging and stochastic gradient descent. In contrast, [3], [4], [5] rely on random walk-based approximate consensus and dual averaging [6], [7]. The authors of [5] consider the optimization problem in the presence of stragglers (slow workers) and introduce Anytime MiniBatch (AMB) [8] to exploit stragglers. The idea behind AMB is to allocate all workers a fixed amount of time for gradient computations in each iteration so that slow workers do not hold up the system. A time limit is imposed on the consensus phase as well to ensure the system does not stall due to random communication delays. Workers then apply the gradients obtained in the consensus phase and proceed to the next iteration.

The closest system models to ours are those of [3] and [5]. We consider a generalized model that captures important aspects of each paper. Specifically, [3] assumes non-identical data distributions at the workers but does not take into account that the gradients may be computed using different amounts of data per iteration. On the other hand, [5] considers identical data distributions but variable amount of gradient computations. In this case, gradients from workers are weighted according to the amount of work performed. In Sec. IV we discuss how these two ideas can be combined to achieve faster convergence when workers complete variable amounts of work and data distributions across workers are different.

## III. SYSTEM MODEL

Similar to [2], [3], [4], [5], we consider a distributed optimization problem defined on an undirected graph  $\mathcal{G}(\mathcal{V}, \mathcal{E})$ . Here,  $\mathcal{V}$  and  $\mathcal{E}$  stand for the set of vertices and the set of edges in the graph. A vertex represents a worker and an edge represents a bi-directional communication link. Specifically,  $(i, j) \in \mathcal{E}$  means that there exists a direct link between workers  $i$  and  $j$ . Note that communication between any two workers that are not directly connected must be relayed through others. We denote the number of workers  $|\mathcal{V}|$  by  $n$ . Also, we assume that the graph is connected, i.e., there is a path that traverses edges and connects any vertex to any other vertex in the graph.

### A. Worker data distributions

Workers share a cost function  $f : \mathbb{R}^d \times \mathcal{X} \rightarrow \mathbb{R}$  which is parameterized by a vector  $w \in \mathbb{R}^d$ . The cost for data point  $X \in \mathcal{X}$  is  $f(w, X)$ . We assume that  $f$  is differentiable,  $L$ -Lipschitz convex in  $w$  for all  $X$  and has bounded gradient, i.e.,  $\|\nabla_w f(w, X)\| \leq L$ . Also, we assume that the data distributions

across workers are non-identical. The  $i$ th worker can only sample data from distribution  $Q_i$ . We define the expected cost for the  $i$ th worker with respect to its data distribution as  $F_i(w) \stackrel{\text{def}}{=} \mathbb{E}_{X \sim Q_i}[f(w, X)]$ . Here, the random variable  $X$  is an abstraction for common optimization problems such as unsupervised and supervised learning. For example, in the supervised case  $X$  represents the (data, label) pair and  $Q_i$  is the joint distribution of the pair.

We are interested in finding a globally optimal parameter vector considering *data across all* workers. To this end we define a mixture distribution  $Q \stackrel{\text{def}}{=} \sum_{i=1}^n \gamma_i Q_i$ . The priors  $\gamma_i \geq 0$  represent the relative importance of each distribution. They are assumed known and satisfy  $\sum_{i=1}^n \gamma_i = 1$ . For example, if worker datasets are finite the normalized sizes of datasets can be taken to be  $\gamma_i$ . However, in our analysis we consider the more general situation where the size of each dataset is not necessarily proportional to its prior. In this case we define the global objective across workers to be

$$F(w) \stackrel{\text{def}}{=} \mathbb{E}_{X \sim Q}[f(w, X)] = \sum_{i=1}^n \gamma_i \mathbb{E}_{X \sim Q_i}[f(w, X)] = \sum_{i=1}^n \gamma_i F_i(w). \quad (1)$$

We want to design a system that enables all workers to converge to

$$w^* \stackrel{\text{def}}{=} \arg \min_{w \in \mathbb{R}^d} F(w).$$

Let us denote the minimizer of  $F_i(w)$  by  $w_i^*$ . Note that with no prior assumptions on the  $Q_i$ ,  $w_i^* \neq w_j^*$  for  $i \neq j$ . In other words, the  $i$ th worker has access to only  $F_i$  yet its goal is to collaboratively and iteratively find  $w^*$ . In subsequent sections we use

$$g_i(w, X) \stackrel{\text{def}}{=} \nabla_w f_i(w, X)$$

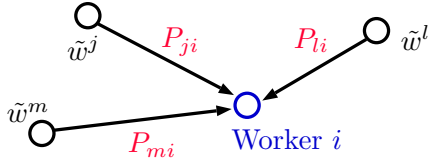
to denote the gradient computed at the  $i$ th worker using data point  $X$ . We note that the bounded gradient assumption on  $f$  also tells us that  $\|\nabla F_i(w)\| \leq L$  which can be shown by using Jensen's inequality and the convexity of norms.

## B. Variable computations

Let us now consider a case analogous to [5], i.e., when stragglers are present. In this paper we assume a *homogeneous* cluster. With this assumption straggler statistics can be considered i.i.d. across workers, although the data distributions may differ. In any iteration, the  $i$ th worker computes gradients from  $b_i$  data samples. Specifically, for a given  $w$  the worker computes  $g_i(w, X)$  using  $b_i$  realizations of  $X$  sampled from  $Q_i$ . We denote by  $\bar{g}_i(w)$  the average of the  $b_i$  realizations of  $g_i(w, X)$ . We do not include the iteration index in  $b_i$  or  $\bar{g}_i$  to avoid clutter. The  $b_i$  themselves are random variables taking values in  $\{1, 2, \dots\}$ . We assume that the  $b_i$  are i.i.d. for all  $i \in [n]$  and across iterations. This means  $b_i$  may differ from iteration to iteration even for the same worker. Our assumptions on  $b_i$  are consistent with homogeneous computing clusters whose straggler statistics are identical across compute nodes. We also note that  $b_i \geq 1$  and each worker computes at least one gradient sample. The latter assumption will prove to be useful in our analysis in Sec. IV-C.

## C. Consensus phase

In this paper we consider gradient descent as the core optimization algorithm, whereas [3] and [5] use dual averaging. After computing  $\bar{g}_i$  workers locally apply gradients and move to the consensus phase. The goal of this phase is to let workers synchronize by averaging across the local parameter vectors. We show that this strategy drives all workers to converge to  $w^*$ . Let  $k \in \{0, 1, 2, \dots\}$  denote the iterate and let  $w_i^k$  denote the parameter vector at the  $i$ th worker in iteration  $k$ . We assume that all workers are initialized to  $w_i^0 = w^0$ . As per Sec. III-B, the  $i$ th worker computes  $\bar{g}_i(w_i^k)$  in the  $k$ th iteration. In this paper we use the random walk-based approximate consensus proposed in [3]. Let  $P \in \mathbb{R}^{n \times n}$  be a doubly stochastic



$j, l, m$ : neighbours of worker  $i$

$$\tilde{w}^i \leftarrow \tilde{w}^i P_{ii} + \tilde{w}^j P_{ji} + \tilde{w}^l P_{li} + \tilde{w}^m P_{mi}$$

Fig. 2: Example of one random walk-based consensus round. The message vectors are denoted by  $\tilde{w}$ . Worker  $i$  computes the weighted sum of  $\tilde{w}^i$  and the messages coming from its neighbours. At the end of one consensus round the  $i$ th worker updates the message vector as shown.

TABLE I: Definitions of system parameters.

Parameter	Definition	Parameter	Definition
$Q_i$	data distribution of $i$ th worker	$\bar{g}_e$	$\sum \gamma_i \bar{g}_i$
$Q$	$\sum_{i=1}^n \gamma_i Q_i$	$\sigma_e^2$	$\mathbb{V}(\bar{g}_e)$
$b_i$	number of gradient samples at $i$ th worker	$\bar{g}_p$	$\sum \frac{nb_i}{b} \gamma_i \bar{g}_i$
$b$	$\sum_{i=1}^n b_i$	$B$	$\{b_1, \dots, b_n\}$
$g_i = g_i(w, X)$	$\nabla_w f_i(w, X)$	$\mu_1$	$\mathbb{E}[b_i/b]$
$\bar{g}_i = \bar{g}_i(w)$	average of $b_i$ realizations of $g_i(w, X)$	$\mu_2$	$\mathbb{E}[1/b_i]$
$\nabla_i$	$\nabla F_i(w)$	$\mu_3$	$\mathbb{E}[b_i/b^2]$
$\nabla$	$\nabla F(w)$	$c_i$	$\frac{b_i}{b} - \mu_1$
$\sum$	$\sum_{i=1}^n$	$s^2$	$\mathbb{E}[c_i^2]$

matrix whose  $i, j$ th element  $P_{i,j} > 0$  only if  $(i, j) \in \mathcal{E}$ . We denote the  $m$ th matrix power of  $P$  by  $[P]^m$ . Methods for generating this type of a matrix include those based on Metropolis-Hastings weights [9] and graph Laplacians [3]. Since  $P$  is a doubly stochastic matrix, all entries in  $[P]^m$  converge to  $\frac{1}{n}$  as  $m$  grows. This can be shown by first observing that the all-ones vector is a left *and* a right eigenvector of  $P$ , corresponding to the largest eigenvalue 1. One can then consider an eigendecomposition of  $P$  to show convergence. The rate of convergence is determined by the second largest eigenvalue of  $P$ . In random walk-based consensus the  $i$ th worker takes a weighted sum of the message vectors from itself and its neighbours. Entries in the  $i$ th row of  $P$  are taken as the weights, and as per the construction of  $P$  all non-neighbours have zero weights. This process is illustrated in Fig. 2. All workers receive the average of message vectors if the exchanging and summing operations are iteratively carried out for many rounds.

Table I summarizes the list of the system parameters and the notations used in the paper. Note that the table consists of parameters that will be introduced in next sections as well.

#### IV. HEURISTIC GRADIENT ESTIMATORS

Let  $W_k$  and  $G_k$  be the  $n$ -column matrices whose  $i$ th columns are  $w_i^k$  and  $\bar{g}_i(w_i^k)$  respectively. After computing  $\bar{g}_i(w_i^k)$ , workers can apply the gradient and update the parameter vector in any manner that leads to the fastest convergence. We discuss in Sec. IV two heuristic methods of applying gradients at workers. In Sec. IV-A we present the two methods, in Sec. IV-B we present numerical results obtained with the two methods, and finally in Sec. IV-C we present a convergence analysis of the two methods. We show in Sec. IV-C that the method that produces the unbiased estimate of  $\nabla F(w)$  with the lower variance leads to a faster convergence. We call the two methods heuristic gradient estimators to reflect that they are estimating  $\nabla F(w)$ .

##### A. Two gradient weighting methods

1) *Equal weighting*: The idea in this scheme is to locally apply gradient as  $w_i^k - tn\gamma_i \bar{g}_i(w_i^k)$  and use the result to perform  $m$  random walk-based consensus rounds. Here,  $t > 0$  is the step size. Workers then

take the output of consensus as  $w_j^{k+1}$  and proceed to the next iteration. Let  $V_1$  be the diagonal matrix whose diagonal is  $(\gamma_1, \dots, \gamma_n)$ . Then  $w_j^{k+1}$  is given by the  $j$ th column of

$$W^{k+1} = (W^k - tnV_1G^k)[P]^m.$$

Since  $P$  is a doubly stochastic matrix, the product  $n[P]^m$  converges to the all-ones matrix as  $m$  grows. In the limit of  $m$ , all columns of  $W^{k+1}$  converge to

$$w_j^{k+1} = \frac{1}{n} \sum_{i=1}^n w_i^k - t \sum_{i=1}^n \gamma_i \bar{g}_i(w_i^k). \quad (2)$$

Since all workers are initialized to  $w_i^0 = w^0$ , (2) is equivalent to perfect consensus and the parameter vectors are identical across workers for all iterations. We can equivalently write the update equation for all workers as

$$w_j^{k+1} = w_j^k - \eta_k \sum_{i=1}^n \gamma_i \bar{g}_i(w_i^k), \quad (3)$$

irrespective of  $j$ . The term ‘equal weighting’ is used in the sense that  $\gamma_i \bar{g}_i(w_i^k)$  from all workers are treated equally when computing  $w_j^{k+1}$ . This is in contrast to the scheme we describe next.

2) *Proportional weighting*: In the previous scheme, although workers compute  $\bar{g}_i(w_i^k)$  using different values for  $b_i$ , the  $b_i$  are not taken into account when combining the estimates. For example, the  $i$ th worker may compute  $\bar{g}_i(w_i^k)$  with  $b_i = 10$  whereas the  $j$ th worker may compute  $\bar{g}_j(w_j^k)$  with  $b_j = 100$ . The latter would be a less noisy estimate. Naturally, one should ask whether we can do better by taking into consideration the confidence of each gradient estimate. This gives rise to the following scheme.

Let  $b \stackrel{\text{def}}{=} \sum_{i=1}^n b_i$ . In this scheme we want to formulate an initial message that enables all workers to receive

$$w_j^{k+1} = w_j^k - \eta_k \sum_{i=1}^n \frac{nb_i}{b} \gamma_i \bar{g}_i(w_i^k) \quad (4)$$

in the limit of  $m$ . Compared to the equal weighting scheme, now  $\gamma_i \bar{g}_i(w_i^k)$  is weighted by its relative confidence  $\frac{nb_i}{b}$ . In the limit of  $m$ , all workers receive the desired parameter vector if the initial message of the  $i$ th worker is set to  $w_i^k - n^2 \frac{b_i}{b} \gamma_i \bar{g}_i(w_i^k)$ . Let  $V_2$  be the diagonal matrix whose diagonal is  $(\frac{nb_1}{b} \gamma_1, \dots, \frac{nb_n}{b} \gamma_n)$ . We can obtain the matrix form of approximate consensus by letting

$$W_{k+1} = (W_k - \eta_k n V_2 G_k)[P]^m.$$

Note that the equal weighting scheme can be recovered by setting a constant for all  $b_i$ . Although equal and proportional schemes are quite similar, the convergence properties of the latter may not be immediately apparent. This is because the  $\frac{b_i}{b}$  themselves are random variables, and at a given iteration a larger  $b_i$  pulls  $w_j^{k+1}$  towards  $w_i^*$  for any  $j$ .

Next, in Sec. IV-B we present a few numerical experiments that demonstrate empirical performance of equal and proportional weighting schemes with a real dataset. In Sec. IV-C we theoretically analyze the convergence properties of the two schemes and we prove that both schemes make all workers converge to  $w^*$ .

## B. Numerical results

We present experiments performed on 10 workers using the Fashion-MNIST dataset [10]. As shown in Fig. 3, the dataset consists of  $28 \times 28$  grayscale images of fashion items that belong to 10 classes. Two sub-datasets that are of sizes 50,000 and 10,000 are available for training and testing. To simulate different distributions we partition the original dataset into 10 groups and assign each to one of the 10

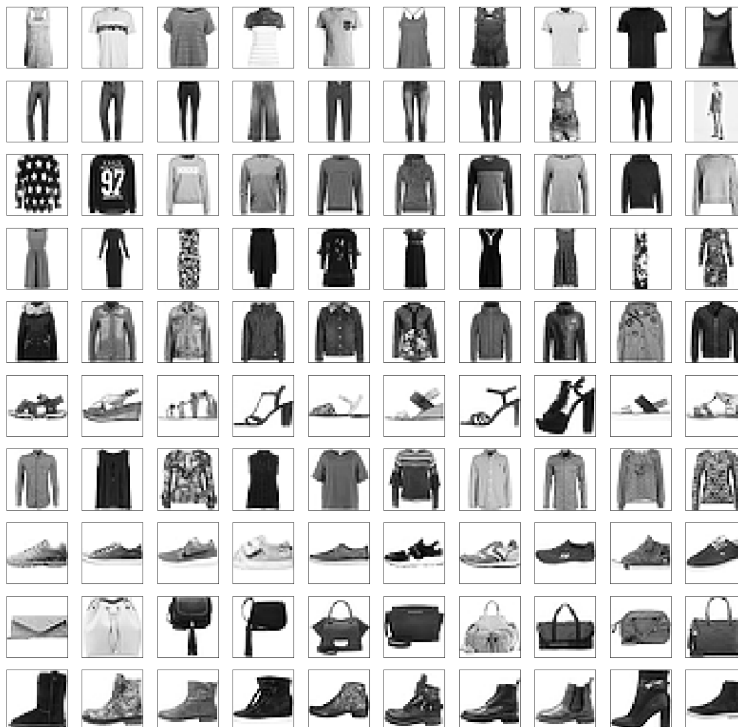


Fig. 3: A few examples from the 10 classes in the Fashion-MNIST dataset (one class per row).

workers. The partitioning is done in a way that only the  $i$ th worker receives images of class  $i$ . This is to ensure that the distributions of different workers are distinct. The workers are connected as per the topology shown in Fig. 1. The label of each worker in Fig. 1 indicates  $i$ , the class index of the images assigned to that worker.

The matrix  $P$  is defined using Metropolis-Hastings method and its second largest eigenvalue is 0.888. We select cross entropy loss of a multinomial logistic regression classifier as the cost function. The distinction between the two methods considered in this work is best illustrated when the size of the parameter vector  $w$  is large. To increase the size of  $w$  we include one hidden layer in the classifier *without* using an activation function. The resulting cost function, which we denote by  $f_{lr}$ , remains convex in  $w$  but the size of the parameter vector is now increased. We also test the non-convex function obtained by applying  $\max(0, \cdot)$  as the activation in the hidden layer. This cost function is denoted by  $f_{nn}$ .

To simulate stragglers we sample  $b_i$  from a Bernoulli distribution in an i.i.d. manner. A worker chooses  $b_i = 60$  data points with probability 0.8 and  $b_i = 1$  with probability 0.2. These choice of numbers make the distribution of  $b_i$  heavily skewed and result in quite noisy gradients at some workers. Convergence results are shown in Fig. 4. The costs of all workers are identical in perfect consensus. For approximate consensus  $m$  is set to 10. In this case the workers have different weights (and costs) as iterations progress, therefore we plot the costs of all 10 workers. Note that these are the values of  $F(w_i^k)$ , i.e., the costs with respect to the global data distribution. We observe that the proportional method outperforms equal weighting under both consensus schemes. In particular, the equal weighting plots are noticeably more noisy than those of proportional weighting.

### C. Convergence analysis

In this section we prove the convergence of the two weighting schemes by assuming perfect consensus amongst workers. Recall that with perfect consensus at the end of the consensus phase all workers possess the same gradient estimate, and apply on the same parameter vector. This means in each iteration workers apply the same gradient and obtain identical parameter vectors. To simplify notation we denote  $g(w, X)$

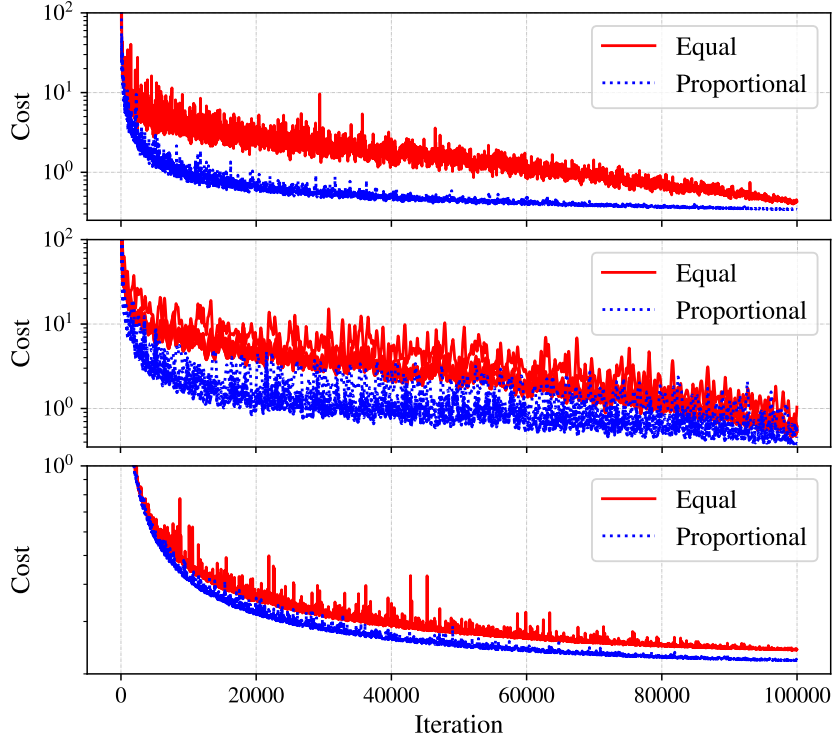


Fig. 4: Results for  $f_{lr}$  with perfect consensus (top), for  $f_{lr}$  with approximate consensus (middle), and for  $f_{nn}$  with perfect consensus (bottom). To suppress noise, the data in the middle figure have been smoothed using a Gaussian filter with standard deviation 5.

for  $X \sim Q_i$  by  $g_i$ , and  $\bar{g}_i(w)$  by  $\bar{g}_i$ . Recall that  $\bar{g}_i$  is the average of  $b_i$  instances of  $g_i$ . In the interest of reducing clutter, we use  $\nabla_i$ ,  $\nabla$  and  $\sum$  to denote  $\nabla F_i(w)$ ,  $\nabla F(w)$  and  $\sum_{i=1}^n$  respectively. Before proceeding to the main proof we first state the following few properties that will become useful.

1) Note that

$$\mathbb{E}[g_i] = \mathbb{E}_{X \sim Q_i}[g(w, X)] = \nabla \mathbb{E}_{X \sim Q_i}[f(w, X)] = \nabla_i.$$

For a given  $b_i \geq 1$ ,  $\bar{g}_i$  is the average of  $b_i$  i.i.d. realizations of  $g_i$ , which means  $\mathbb{E}[\bar{g}_i | b_i] = \nabla_i$ . Since  $\bar{g}_i$  and  $\bar{g}_j$  are independent for  $i \neq j$ , letting the set  $B \stackrel{\text{def}}{=} \{b_1, \dots, b_n\}$  we have the identities

$$\mathbb{E}[\bar{g}_i | B] = \mathbb{E}[\bar{g}_i | b_i] = \nabla_i \quad \text{and} \quad \mathbb{E}[\bar{g}_i] = \mathbb{E}[\mathbb{E}[\bar{g}_i | B]] = \mathbb{E}[\nabla_i] = \nabla_i. \quad (5)$$

2) Next, we show that  $\mathbb{E}[b_i/b] = 1/n$ . Since all the  $b_i$  have i.i.d. statistics,  $\mathbb{E}[b_i/b] = \mathbb{E}[b_j/b] = \mu_1$  for all  $i, j \in [n]$ , for some constant  $\mu_1$ . We have

$$1 = \mathbb{E} \left[ \sum b_i/b \right] = \sum \mathbb{E}[b_i/b] = n\mu_1. \quad (6)$$

Solving for  $\mu_1$  yields the desired result. Similarly, note that  $\mathbb{E}[1/b_i]$  and  $\mathbb{E}[b_i/b^2]$  are constants independent of  $i$ , thus we denote them by  $\mu_2$  and  $\mu_3$  respectively.

3) We denote the variance of a random vector by  $\mathbb{V}$  where  $\mathbb{V}(\cdot) = \mathbb{E}[\|\cdot\|^2] - \|\mathbb{E}[\cdot]\|^2$ . The law of total variance states that if  $Z$  and  $Y$  are random variables on the same probability space, and the variance of  $Y$  is finite, then

$$\mathbb{V}(Y) = \mathbb{E}[\mathbb{V}(Y | Z)] + \mathbb{V}(\mathbb{E}[Y | Z]).$$

Using the law of total variance and (5), we upper bound  $\mathbb{V}(\bar{g}_i)$  as

$$\mathbb{V}(\bar{g}_i) = \mathbb{E}[\mathbb{V}(\bar{g}_i | b_i)] + \mathbb{V}(\mathbb{E}[\bar{g}_i | b_i]) \leq \mathbb{E}[\sigma^2/b_i] + \mathbb{V}(\nabla_i) \leq \sigma^2 \mu_2. \quad (7)$$

The first inequality is due to the assumption  $\mathbb{V}(g_i) \leq \sigma^2$  and the second inequality follows by observing that  $\mathbb{V}(\nabla_i) = 0$ . Also, note that from (1) we have

$$\nabla = \sum_{i=1}^n \gamma_i \nabla_i. \quad (8)$$

The proofs presented next rely on the the following theorem ([11] pg.192) that characterizes the convergence of stochastic gradient descent.

**Theorem 1.** *Let  $h : \mathbb{R}^d \rightarrow \mathbb{R}$  be a  $L$ -Lipschitz convex function and  $w^* = \arg \min_w h(w)$ . For all  $k \in \{0, 1, 2, \dots\}$  let  $v^k$  be a random vector such that  $\mathbb{E}[v^k] = \nabla_w h(w^k)$ , and  $\mathbb{V}(v^k) \leq \sigma^2$ . For a given  $w^0$  let the sequence  $(w^1, w^2, \dots)$  be such that  $w^{k+1} = w^k - tv^k$ . Then for a learning rate  $t \leq 1/L$  and  $\bar{w}^k = (w^1 + \dots + w^k)/k$  we have  $\mathbb{E}[h(\bar{w}^k)] \leq h(w^*) + \frac{\|w^0 - w^*\|_2}{2tk} + \frac{t\sigma^2}{2}$ .*

Note that convergence of SGD requires only that  $v^k$  be an unbiased estimator of the true gradient and that it have finite variance. Next we make use of this observation to prove convergence of the equal and proportional weighting schemes, by showing that they are in fact unbiased estimators of  $\nabla$ .

1) *Proof for equal weighting:* From (3), all workers possess the same parameter vector and apply the gradient  $\bar{g}_e \stackrel{\text{def}}{=} \sum \gamma_i \bar{g}_i$ . To prove convergence it suffices to show that  $\mathbb{E}[\bar{g}_e] = \nabla$  and that  $\sigma_e^2 \stackrel{\text{def}}{=} \mathbb{V}(\bar{g}_e)$  is bounded. First note that

$$\mathbb{E}[\bar{g}_e] = \mathbb{E} \left[ \sum \gamma_i \bar{g}_i \right] = \mathbb{E} \left[ \sum \gamma_i \mathbb{E}[\bar{g}_i | b_i] \right] = \sum \gamma_i \nabla_i = \nabla,$$

which is implied by (5) and (8). Second, using (7) we have

$$\sigma_e^2 = \mathbb{V} \left( \sum \gamma_i \bar{g}_i \right) = \sum \mathbb{V}(\gamma_i \bar{g}_i) \leq \mu_2 \sigma^2 \sum \gamma_i^2. \quad (9)$$

The last equality follows because the  $\bar{g}_i$  are independent. This is the tightest possible bound without making additional assumptions on the distribution of  $b_i$ . Now, according to Theorem 1, all workers converge to  $w^*$ .

2) *Proof for proportional weighting:* All workers update the parameter vectors according to (4) using  $\bar{g}_p \stackrel{\text{def}}{=} \sum \frac{nb_i}{b} \gamma_i \bar{g}_i$  as the gradient. From (5), (6) and (8) we have

$$\mathbb{E}[\bar{g}_p] = \mathbb{E} \left[ \mathbb{E} \left[ \frac{n}{b} \sum b_i \gamma_i \bar{g}_i \middle| B \right] \right] = \mathbb{E} \left[ \frac{n}{b} \sum b_i \gamma_i \mathbb{E}[\bar{g}_i | B] \right],$$

which gives

$$\mathbb{E}[\bar{g}_p] = \mathbb{E} \left[ \frac{n}{b} \sum b_i \gamma_i \nabla_i \right] = \sum \gamma_i \nabla_i n \mathbb{E}[b_i/b] = \nabla.$$

This result proves  $\bar{g}_p$  is an unbiased estimator of  $\nabla$ . To bound  $\mathbb{V}(\bar{g}_p)$  we apply the law of total variance as

$$\sigma_p^2 \stackrel{\text{def}}{=} \mathbb{V}(\bar{g}_p) = \underbrace{\mathbb{E} \left[ \mathbb{V} \left( \frac{n}{b} \sum b_i \gamma_i \bar{g}_i \middle| B \right) \right]}_{C_1} + \underbrace{\mathbb{V} \left( \mathbb{E} \left[ \frac{n}{b} \sum b_i \gamma_i \bar{g}_i \middle| B \right] \right)}_{C_2}.$$

Since  $b_i \bar{g}_i$  is the sum of  $b_i$  i.i.d. realizations of  $g_i$ ,  $\mathbb{V}(b_i \bar{g}_i | B) \leq b_i \sigma^2$ . This gives us

$$C_1 \leq n^2 \sigma^2 \mathbb{E} \left[ \sum \frac{b_i}{b^2} \gamma_i^2 \right] = n^2 \mu_3 \sigma^2 \sum \gamma_i^2,$$

which is the tightest possible bound with the given assumptions on  $b_i$ .

Next, we bound  $C_2$ . From (5) we have

$$C_2 = \mathbb{V} \left( \frac{n}{b} \sum \gamma_i b_i \mathbb{E}[\bar{g}_i | B] \right) = \mathbb{V} \left( n \sum \gamma_i \frac{b_i}{b} \nabla_i \right). \quad (10)$$



Noting from (6) that

$$\mathbb{E} \left[ n \sum \gamma_i \frac{b_i}{b} \nabla_i \right] = n \sum \gamma_i \mathbb{E} \left[ \frac{b_i}{b} \right] \nabla_i = \sum \gamma_i \nabla_i = \nabla,$$

we expand the variance term in (10) to write

$$C_2 = \mathbb{E} \left[ \left\| n \sum \gamma_i \frac{b_i}{b} \nabla_i - \nabla \right\|^2 \right] = \mathbb{E} \left[ \left\| \sum \frac{b_i}{b} (n\gamma_i \nabla_i - \nabla) \right\|^2 \right] = \mathbb{E} \left[ \left\| \sum \frac{b_i}{b} \Delta_i \right\|^2 \right].$$

Here we use  $\Delta_i$  to denote the difference  $(n\gamma_i \nabla_i - \nabla)$ . Note that from (8) we have  $\sum \Delta_i = 0$ . Letting  $c_i \stackrel{\text{def}}{=} \frac{b_i}{b} - \mu_1$  we get  $\mathbb{E}[c_i] = 0$  and

$$\mathbb{V} \left( \frac{b_i}{b} \right) = \mathbb{E}[c_i^2] \stackrel{\text{def}}{=} s^2 \quad (11)$$

for all  $i \in [n]$  for some constant  $s$ . Next we use the identity

$$\sum c_i \Delta_i = \sum \left( \frac{b_i}{b} - \mu_1 \right) \Delta_i = \sum \frac{b_i}{b} \Delta_i - \mu_1 \sum \Delta_i = \sum \frac{b_i}{b} \Delta_i \quad (12)$$

to obtain

$$C_2 = \mathbb{E} \left[ \left\| \sum c_i \Delta_i \right\|^2 \right] \leq \mathbb{E} \left[ \sum c_i^2 \sum \|\Delta_i\|^2 \right] = \mathbb{E} \left[ \sum c_i^2 \right] \sum \|\Delta_i\|^2 = n s^2 \sum \|\Delta_i\|^2 \leq n^3 s^2 D.$$

The first equality is due to (12), and the first inequality is due to triangle inequality and Cauchy-Schwarz inequality. The constant  $s$  is as defined in (11). Note that  $\sum \|\Delta_i\|^2$  is proportional to the mean squared error of true gradients. The last inequality is obtained by assuming that there exists a constant  $D \geq 0$  such that

$$\sum \|\Delta_i\|^2 \leq n^2 D. \quad (13)$$

Later we show that in fact such a constant exists. Combining the bounds for  $C_1$  and  $C_2$  give us

$$\sigma_p^2 \leq n^2 \mu_3 \sigma^2 \sum \gamma_i^2 + n^3 s^2 D. \quad (14)$$

This result, along with  $\mathbb{E}[\bar{g}_p] = \nabla$  prove the convergence of the proportional method.

Now we show that a constant  $D \geq 0$  exists that satisfies (13). The  $L$ -Lipschitz continuous assumption on  $f$  implies  $\|\nabla_w f(w, X)\| \leq L$ . We have

$$\|\nabla_i\| = \|\mathbb{E}_{X \sim Q_i}[\nabla_w f(w, X)]\| \leq \mathbb{E}_{X \sim Q_i}[\|\nabla_w f(w, X)\|] \leq \mathbb{E}[L] = L,$$

where we get the first inequality by asserting convexity of norms and applying the Jensen's inequality. Using this result, an easy albeit loose candidate for  $D$  can be obtained as follows. We have

$$\sum \|\Delta_i\|^2 = \sum \|n\gamma_i \nabla_i - \nabla\|^2 = n^2 \sum \gamma_i^2 \|\nabla_i\|^2 - n \|\nabla\|^2 \leq n^2 L^2 \sum \gamma_i^2,$$

which lets us take  $L^2 \sum \gamma_i^2$  as  $D$ . This shows the existence of a finite  $D$  that satisfies (13), and concludes the proof of the convergence of the proportional method.

3) *Convergence rates comparison:* In this section we compare the two upper bounds obtained for  $\sigma_c^2$  and  $\sigma_p^2$ . We start by showing that  $n^2\mu_3 \leq \mu_2$ . We have  $\mathbb{E}[b_i|b] = \mathbb{E}[b_j|b] = \frac{1}{n}\sum\mathbb{E}[b_i|b] = \frac{1}{n}\mathbb{E}[\sum b_i|b] = \frac{b}{n}$ . Since the reciprocal function is convex, by Jensen's inequality  $\mathbb{E}\left[\frac{1}{b_i}|b\right] \geq \frac{1}{\mathbb{E}[b_i|b]} = \frac{n}{b}$ . Now we have

$$\mathbb{E}\left[\frac{n}{b} - \frac{1}{b_i}\right] = \mathbb{E}\left[\frac{n}{b} - \mathbb{E}\left[\frac{1}{b_i}|b\right]\right] \leq \mathbb{E}\left[\frac{n}{b} - \frac{n}{b}\right] = 0,$$

which gives  $\mathbb{E}\left[\frac{n}{b}\right] \leq \mathbb{E}\left[\frac{1}{b_i}\right]$ . Therefore,

$$n^2\mu_3 = n^2\mathbb{E}\left[\frac{b_i}{b^2}\right] = n^2\mathbb{E}\left[\mathbb{E}\left[\frac{b_i}{b^2}|b\right]\right] = \mathbb{E}\left[\frac{n^2}{b^2}\mathbb{E}[b_i|b]\right],$$

which gives

$$n^2\mu_3 = \mathbb{E}\left[\frac{n^2}{b^2}\mathbb{E}[b_i|b]\right] = \mathbb{E}\left[\frac{n^2}{b^2}\frac{b}{n}\right] = \mathbb{E}\left[\frac{n}{b}\right] \leq \mathbb{E}\left[\frac{1}{b_i}\right] = \mu_2. \quad (15)$$

As suggested by the SGD convergence properties summarized in Theorem 1, the proportional method converges faster if  $\sigma_p^2 \leq \sigma_c^2$ . The two upper bounds we obtained for  $\sigma_p^2$  and  $\sigma_c^2$  are the tightest possible bounds without making additional assumptions on  $b_i$ . Therefore, in general we expect the proportional method to converge faster if the upper bound for  $\sigma_p^2$  is less than that for  $\sigma_c^2$ . Substituting from (9) and (14) to the upper bounds we have

$$n^2\mu_3\sigma^2\sum\gamma_i^2 + n^3s^2D \leq \mu_2\sigma^2\sum\gamma_i^2.$$

Rearranging the terms gives the condition for the proportional method to converge faster than equal weighting

$$D/\sigma^2 \leq (\mu_2 - n^2\mu_3)\sum\gamma_i^2/(n^3s^2). \quad (16)$$

The right side is non-negative due to the result in (15).

This condition is insightful. The right side depends only on the straggler statistics. Recall that  $\sigma^2$  and  $D$  are due to our assumptions  $\mathbb{V}(g_i) \leq \sigma^2$  and  $\sum\|\Delta_i\|^2 \leq n^2D$ . While  $D$  measures how different the true gradients across workers are (global variation),  $\sigma^2$  measures the variance of the gradient distribution within one worker (local variance). The inequality binds together attributes of three different phenomena: stragglers, measurement noise and data distributions. If the other factors remain the same, a large  $\sigma^2$  is likely to satisfy the condition, making the proportional method converge faster.

In Fig. 5 we attempt to visualize the condition in (16) using an example. In this example we assume a 2-dimensional gradient. The  $x$  and  $y$  axes represent the two components of the gradient. Let  $\gamma_i = 1/4$ , and we assume that  $w$  is fixed. In the left figure, the four shaded clusters represent distributions of the random variables  $\nabla_w f(w, X); X \sim Q_i$  for  $i \in [4]$ . At the centre of each cluster is the expected gradient  $\nabla F_i(w) = \nabla_i = \mathbb{E}_{X \sim Q_i}[\nabla_w f(w, X)]$ . The variance of  $\nabla_w f(w, X)$  is  $\sigma^2$  in all clusters. As per (8)  $\nabla = \frac{1}{4}\sum_{i=0}^4 \nabla_i$  and per (13)  $\nabla$  has a distance  $\sqrt{D}$  to all cluster centres. The right figure is same as the left, except that its  $\sigma^2$  is smaller. In both figures  $D$  is the same. The figures illustrate that  $\sigma^2$  is a measurement local to the workers (a single shaded region), and  $D$  is a measure of the global variation among workers (all four shaded regions). Going back to (16), we observe that a smaller  $\sigma^2$  is likely to violate the inequality, making the equal method converge faster. As per the figure on the right, a smaller  $\sigma^2$  also implies a smaller noise radius. Even  $b_i = 1$  is enough to estimate a cluster centre  $\nabla_i$  with a relatively high accuracy, and the error will be low compared to  $D$ . However, if  $\sigma^2$  is large, we need a large  $b_i$  to keep the measurement error low. In such cases the proportional method performs better by leveraging the gradients with higher confidence.

Note that  $D$  is an implicit measure of the divergence of worker distributions  $Q_i$ , which is observed through  $\nabla_i$ s. Authors of [5] consider the case when distributions of workers are the same. This means we

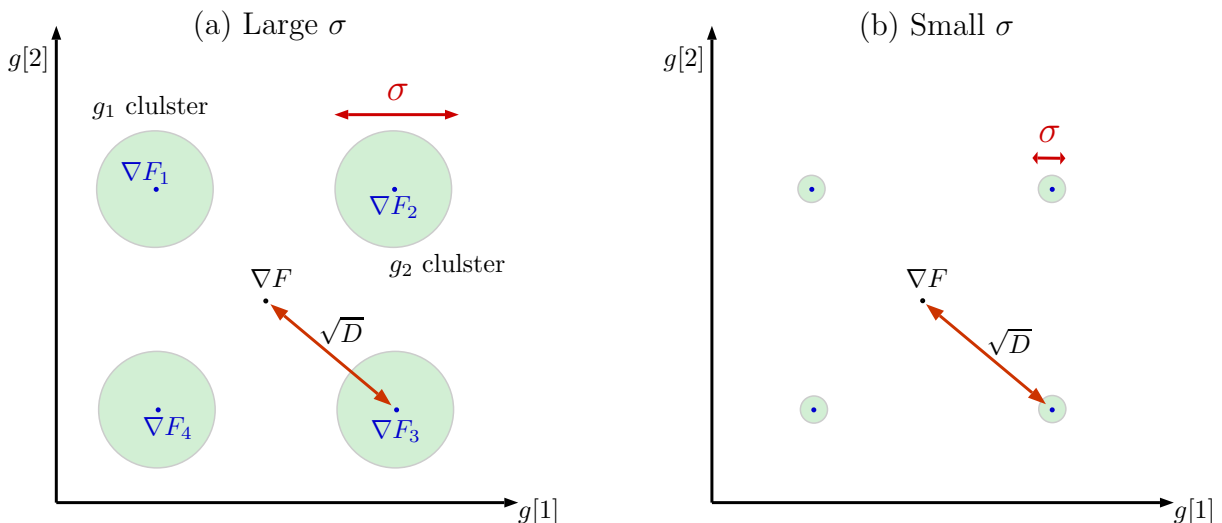


Fig. 5: Visualizing (16), the condition for the proportional method to converge faster than equal weighting.

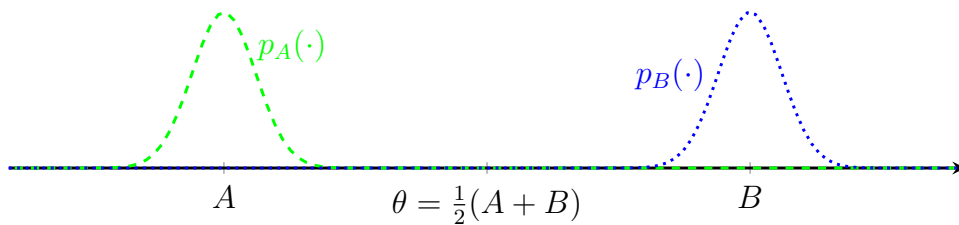


Fig. 6: Samples are drawn from  $p_A$  and  $p_B$  with the goal of estimating  $\theta$ .

have  $\nabla_i = \nabla$ , and we are allowed to set  $D = 0$  to get the tightest bound in (13). If  $D = 0$ , (15) implies that the condition in (16) is always true. In this case the proportional scheme is guaranteed to perform better, and is inline with the findings in [5]. However, when  $D \neq 0$  the proportional scheme outperforms equal weighting only if  $D/\sigma^2$  is small enough. This means no matter how noisy the  $\bar{g}_i$  (measured through  $\sigma^2$ ) from some workers are, they cannot be weighted down if the gradient distributions are significantly different (observed through  $D$ ). While  $D/\sigma^2$  is not measurable for real world datasets, the proportional weighting method performs better in the experiments we present in Sec. IV-B. Therefore, the condition (16) seems to be satisfied in those experiments.

## V. OPTIMAL GRADIENT ESTIMATOR

In Sec. IV we analyzed two methods of combining the gradient estimates of workers. We showed that  $\bar{g}_e$  and  $\bar{g}_p$  corresponding to the equal and proportional methods are unbiased estimators of  $\nabla$  and have finite variances. As per Theorem 1, the unbiased estimator with the smallest variance wins the gradient descent race. We would like to know if there exist any other estimators that achieve a smaller variance. In this section we make use of the theory on ‘minimum variance unbiased estimator’ (MVUE) to answer the question.

To better understand the problem at hand, we start with the toy problem illustrated in Fig. 6, which is obtained by assuming two workers and a scalar parameter vector. We assume the parameter vector is fixed and  $\gamma_1 = \gamma_2 = \frac{1}{2}$ . Let  $\nabla_1 = A$  and  $\nabla_2 = B$  be the true gradients at the first and second workers respectively. Our goal is to obtain the minimum variance estimate of  $\theta = \frac{1}{2}(A + B)$ , the average gradient of the two workers. We assume the gradients sampled at the two workers are normally distributed with variance  $\sigma^2$ , and means  $A$  and  $B$ . This distribution is due to the randomness of data available at each worker. The probability distributions for the samples drawn at the workers are denoted by  $p_A$  and  $p_B$ . Let

the first worker obtain  $a_1, \dots, a_M$  i.i.d. samples and the second worker to obtain  $b_1, \dots, b_N$  i.i.d. samples. Now we apply Theorem 2 on the Cramer-Rao Lower Bound (CRLB) to this problem.

**Theorem 2** (Cramer-Rao Lower Bound - Scalar parameter [12]). *For an observation vector  $\mathbf{x}$  it is assumed that the PDF  $p(\mathbf{x}; \theta)$  satisfies the ‘regularity’ condition  $\mathbb{E}\left[\frac{\partial p(\mathbf{x}; \theta)}{\partial \theta}\right] = 0$  for all  $\theta$ , where the expectation is taken with respect to  $p(\mathbf{x}; \theta)$ . Then, the variance of any unbiased estimator  $\hat{\theta}$  must satisfy  $\mathbb{V}(\hat{\theta}) \geq -\mathbb{E}\left[\frac{\partial^2 p(\mathbf{x}; \theta)}{\partial \theta^2}\right]^{-1}$ , where the derivative is evaluated at the true value of  $\theta$  and the expectation is taken with respect to  $p(\mathbf{x}; \theta)$ . Furthermore, an unbiased estimator may be found that attains the bound for all  $\theta$  if and only if  $\frac{\partial p(\mathbf{x}; \theta)}{\partial \theta} = I(\theta)(g(\mathbf{x}) - \theta)$  for some functions  $g$  and  $I$ . That estimator which is the MVUE is  $\hat{\theta} = g(\mathbf{x})$ , and the minimum variance is  $1/I(\theta)$ .*

We look at the following two cases and apply CRLB.

#### A. Constant observation counts

Let  $M$  and  $N$  be constants and let the observation vector be defined as

$$\mathbf{x} = [a_1, \dots, a_M, b_1, \dots, b_N].$$

Since all samples are independent we have the joint PDF

$$p = p(\mathbf{x}; \theta) = \prod_{m=1}^M p_A(a_m) \prod_{n=1}^N p_B(b_n) = \prod_{m=1}^M \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(a_m - A)^2}{2\sigma^2}} \prod_{n=1}^N \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(b_n - B)^2}{2\sigma^2}}.$$

By taking partial derivative of  $\ln p$  with respect to  $\theta$  we have

$$\frac{\partial \ln p}{\partial \theta} = \frac{1}{\sigma^2} \sum_{m=1}^M (a_m - A) \frac{\partial A}{\partial \theta} + \frac{1}{\sigma^2} \sum_{n=1}^N (b_n - B) \frac{\partial B}{\partial \theta}. \quad (17)$$

Since

$$\mathbb{E}\left[\frac{\partial \ln p}{\partial \theta}\right] = \frac{1}{\sigma^2} \sum_{m=1}^M (\mathbb{E}[a_m] - A) \frac{\partial A}{\partial \theta} + \frac{1}{\sigma^2} \sum_{n=1}^N (\mathbb{E}[b_n] - B) \frac{\partial B}{\partial \theta} = 0,$$

the PDF satisfies the regularity condition and we can apply CRLB. Also note that  $\theta = \frac{1}{2}(A + B)$  we have  $\frac{\partial \theta}{\partial A} = \frac{\partial \theta}{\partial B} = \frac{1}{2}$  which gives  $\frac{\partial A}{\partial \theta} = \frac{\partial B}{\partial \theta} = 2$ . Differentiating (17) with respect to  $\theta$  gives

$$\frac{\partial^2 \ln p}{\partial \theta^2} = \frac{1}{\sigma^2} \sum_{m=1}^M \left( (a_m - A) \frac{\partial^2 A}{\partial \theta^2} - \left( \frac{\partial A}{\partial \theta} \right)^2 \right) + \frac{1}{\sigma^2} \sum_{n=1}^N \left( (b_n - B) \frac{\partial^2 B}{\partial \theta^2} - \left( \frac{\partial B}{\partial \theta} \right)^2 \right)$$

and by taking expectation and substituting  $\frac{\partial A}{\partial \theta} = \frac{\partial B}{\partial \theta} = 2$  we have  $\mathbb{E}\left[\frac{\partial^2 \ln p}{\partial \theta^2}\right] = -\frac{4(M+N)}{\sigma^2}$ . This means for any unbiased estimator  $\hat{\theta}$ , by CRLB  $\mathbb{V}(\hat{\theta}) \geq \frac{\sigma^2}{4(M+N)}$ . To find an estimator that achieves the minimum variance we rewrite (17) as follows:

$$\begin{aligned} \frac{\partial \ln p}{\partial \theta} &= \frac{2}{\sigma^2} \sum_{m=1}^M (a_m - A) + \frac{2}{\sigma^2} \sum_{n=1}^N (b_n - B) \\ &= \frac{2}{\sigma^2} \left( \sum_{m=1}^M a_m + \sum_{n=1}^N b_n - (MA + NB) \right). \end{aligned} \quad (18)$$

Note that (18) cannot be written in  $I(\theta)(g(\mathbf{x}) - \theta)$  form for  $\theta = \frac{1}{2}(A + B)$  and we conclude that the MVUE does not exist.

Knowing that an estimator achieving CRLB does not exist, we can compare the variances of the two estimators discussed in Sec. IV with the lower bound. Gradient estimator for the equal method is given

by  $\bar{g}_e = \frac{1}{2} \left( \frac{1}{M} \sum_{m=1}^M a_m + \frac{1}{N} \sum_{n=1}^N b_n \right)$ . We can verify  $\bar{g}_e$  is in fact an unbiased estimator by noting that  $\mathbb{E}[\bar{g}_e] = \frac{1}{2}(A + B) = \theta$ . For variance we have

$$\begin{aligned} \mathbb{V}(\bar{g}_e) &= \frac{1}{4} \left( \frac{1}{M^2} \sum_{m=1}^M \mathbb{V}(a_m) + \frac{1}{N^2} \sum_{n=1}^N \mathbb{V}(b_n) \right) \\ &= \frac{\sigma^2}{4} \left( \frac{1}{M} + \frac{1}{N} \right) \\ &= \frac{\sigma^2}{4(M+N)} \frac{(M+N)^2}{MN} \\ &= \frac{\sigma^2}{4(M+N)} \left( \frac{M}{N} + \frac{N}{M} + 2 \right) \\ &\geq \frac{\sigma^2}{(M+N)}, \end{aligned}$$

where we used the fact that the observations are independent, and  $\frac{M}{N} + \frac{N}{M} \geq 2$  for  $M, N \geq 1$ . This shows that the variance of equal estimator  $\mathbb{V}(\bar{g}_e)$  is at least 4 times larger than the lower bound. For the proportional method the gradient estimator is

$$\bar{g}_p = \frac{M}{M+N} \left( \frac{1}{M} \sum_{m=1}^M a_m \right) + \frac{N}{M+N} \left( \frac{1}{N} \sum_{n=1}^N b_n \right) = \frac{1}{M+N} \left( \sum_{m=1}^M a_m + \sum_{n=1}^N b_n \right).$$

In this case  $\bar{g}_e$  is not an unbiased estimator of  $\theta$  since  $\mathbb{E}[\bar{g}_p] = \frac{MA+NB}{M+N} \neq \theta$ . However for variance we have

$$\mathbb{V}(\bar{g}_p) = \frac{1}{(M+N)^2} (M\sigma^2 + N\sigma^2) = \frac{\sigma^2}{M+N},$$

which is again larger than the CRLB. Next we assume  $M$  and  $N$  are i.i.d. RVs and redo the analysis.

### B. Random observation counts

Let  $M$  and  $N$  be i.i.d. random variables. In this case the observation vector itself is of random length. Let  $q(\cdot)$  be the PDF of  $M$  and  $N$ . We assume  $\mathbb{E}[M] = \mathbb{E}[N] = \mu$ , and  $\mathbb{E}[1/M] = \mathbb{E}[1/N] = \mu_2$ . We assume that  $q$  is *not* a function of  $\theta$ . This assumption is consistent with our distributed optimization setup where the  $b_i$  are independent of  $\nabla_i$ . Now we have  $p = p(\mathbf{x} | M, N; \theta) p_{M,N}(M, N)$ . The first term is same as what we derived in Sec. V-A. By observing  $M$  and  $N$  are independent we write

$$\ln p = \ln p(\mathbf{x} | M, N; \theta) + \ln q(M) + \ln q(N).$$

Noting that  $\theta$  is not a parameter in  $q$ , we differentiate with respect to  $\theta$  to obtain (17) once again. Therefore,  $\mathbb{E} \left[ \frac{\partial \ln p}{\partial \theta} \right] = 0$  for this case as well. The CRLB in this case is given by

$$\mathbb{E} \left[ \frac{\partial^2 \ln p}{\partial \theta^2} \right] = \mathbb{E} \left[ -\frac{4(M+N)}{\sigma^2} \right] = -\frac{4}{\sigma^2} \mathbb{E}[M+N] = -\frac{8\mu}{\sigma^2},$$

which means for any unbiased estimator  $\hat{\theta}$ ,  $\mathbb{V}(\hat{\theta}) \geq \frac{\sigma^2}{8\mu}$ . Regarding the existence of an estimator achieving the lower bound, the same argument as before applies. We cannot write (18) in  $I(\theta)(g(\mathbf{x}) - \theta)$  form, and MVUE does not exist.

Let us now compare the lower bound with the variances of  $\bar{g}_e$  and  $\bar{g}_p$ . The estimator  $\bar{g}_e$  remains unbiased and we have

$$\mathbb{V}(\bar{g}_e) = \frac{1}{4} \left( \mathbb{V} \left( \frac{1}{M} \sum_{m=1}^M a_m \right) + \mathbb{V} \left( \frac{1}{N} \sum_{n=1}^N b_n \right) \right)$$

and for the first (and second) variance term

$$\begin{aligned}
\mathbb{V}\left(\frac{1}{M}\sum_{m=1}^M a_m\right) &= \mathbb{E}\left[\mathbb{V}\left(\frac{1}{M}\sum_{m=1}^M a_m \middle| M\right)\right] + \mathbb{V}\left(\mathbb{E}\left[\frac{1}{M}\sum_{m=1}^M a_m \middle| M\right]\right) \\
&= \mathbb{E}\left[\frac{1}{M^2}\mathbb{V}\left(\sum_{m=1}^M a_m \middle| M\right)\right] + \mathbb{V}(\mathbb{E}[A|M]) \\
&= \mathbb{E}\left[\frac{1}{M^2}\sigma^2 M\right] + 0 \\
&= \sigma^2 \mu_2.
\end{aligned}$$

Substituting the result we get  $\mathbb{V}(\bar{g}_e) = \frac{\sigma^2 \mu_2}{2}$ . The reciprocal of a positive real number is a convex function and by Jensen's inequality we can show that  $\frac{1}{\mu} = \frac{1}{\mathbb{E}[M]} \leq \mathbb{E}\left[\frac{1}{M}\right] = \mu_2$ . We again conclude that  $\mathbb{V}(\bar{g}_e)$  is at least 4 times larger than the CRLB. In contrast to the Sec. V-A proportional estimator now is unbiased since

$$\begin{aligned}
\mathbb{E}[\bar{g}_p] &= \mathbb{E}\left[\frac{M}{M+N}\left(\frac{1}{M}\sum_{m=1}^M a_m\right) + \frac{N}{M+N}\left(\frac{1}{N}\sum_{n=1}^N b_n\right)\right] \\
&= \mathbb{E}\left[\frac{M}{M+N}\mathbb{E}\left[\frac{1}{M}\sum_{m=1}^M a_m \middle| M, N\right]\right] + \mathbb{E}\left[\frac{N}{M+N}\mathbb{E}\left[\frac{1}{N}\sum_{n=1}^N b_n \middle| M, N\right]\right] \\
&= \mathbb{E}\left[\frac{M}{M+N}\right]A + \mathbb{E}\left[\frac{N}{M+N}\right]B \\
&= \frac{1}{2}(A+B).
\end{aligned}$$

The last equality is because  $\mathbb{E}\left[\frac{M}{M+N}\right] = \mathbb{E}\left[\frac{N}{M+N}\right] = \frac{1}{2}$  for i.i.d.  $M, N$ . For the variance we can show that  $\mathbb{V}(\bar{g}_p) \geq \sigma^2 \mathbb{E}\left[\frac{1}{M+N}\right]$  which is greater than the CRLB. Therefore, we can conclude that neither of the two methods achieve the minimum variance, and, more importantly no other estimator does.

## VI. CONCLUSIONS AND FUTURE WORK

In this paper we study the convergence of decentralized optimization when data distributions at workers are non-identical *and* when workers perform variable amounts of work. We numerically and theoretically analyze two heuristic methods of combining worker outputs. We make use of the theory on minimum variance unbiased estimator (MVUE) to evaluate the existence of an optimal method for combining worker outputs. While we conclude that neither of the two heuristic methods are optimal, we also show that an optimal method does *not* exist. A few possible next steps to this theme of works is as follows.

First is improving the theoretical analysis presented in Sec. IV-C. We provide a convergence proof for the proportional method when consensus is perfect. We would like to extend the convergence proof and generalize the condition in (16) to approximate consensus. Such an analysis will largely benefit from the work in [3] as there are many similarities between the algorithms. Also, in our analysis we assume that  $b_i \geq 1$ . A natural next step will be to allow  $b_i = 0$  and generalize our convergence results accordingly. This inclusion means that some workers may not produce a gradient estimate at all, therefore, may skip over some iterations. In other words, only a subset of workers participate in each iteration [13].

Second next step is regarding the interplay between worker connectivity graph and worker distributions. In our work we consider that data distributions at workers are different. However, in reality it is safe to assume that the distributions are similar to some extent, but with a few workers having 'atypical' distributions scattered around. It will be interesting to understand how the position on the graph of these atypical workers impact the convergence. For example, assume three workers whose network connectivity

makes a chain. We have two workers with similar distributions and one worker with a different distribution. It will be interesting to see whether a faster convergence is obtained by placing the outlier distribution in middle of the chain or at one end of the chain. We leave the consideration of this type of a system model as future work.

#### ACKNOWLEDGMENT

The authors would like to thank Haider Al-Lawati at University of Toronto, Jason Lam and Zhenhua Hu at Huawei Technologies Canada for technical discussions, and Compute Canada ([www.computeCanada.ca](http://www.computeCanada.ca)) for providing computing resources for the experiments.

#### REFERENCES

- [1] T. Adikari and S. Draper, “Decentralized optimization with non-identical sampling in presence of stragglers,” in *Proc. Int. Conf. Acoust. Speech, Signal Processing*. Barcelona: IEEE, May 2020.
- [2] O. Dekel, R. Gilad-Bachrach, O. Shamir, and L. Xiao, “Optimal distributed online prediction using mini-batches,” *J. of Machine Learning Research*, pp. 165–202, Jan 2012.
- [3] J. C. Duchi, A. Agarwal, and M. J. Wainwright, “Dual averaging for distributed optimization: Convergence analysis and network scaling,” *IEEE Trans. Automat. Contr.*, pp. 592–606, Jun 2011.
- [4] K. I. Tsianos and M. G. Rabbat, “Efficient distributed online prediction and stochastic optimization with approximate distributed averaging,” *IEEE Trans. Signal Inf. Proc. over Networks*, pp. 489–506, Oct 2016.
- [5] N. Ferdinand, H. Al-Lawati, S. Draper, and M. Nokleby, “Anytime minibatch: Exploiting stragglers in online distributed optimization,” in *Int. Conf. Learning Representations*, New Orleans, May 2019.
- [6] Y. Nesterov, “Primal-dual subgradient methods for convex problems,” *Mathematical programming*, pp. 221–259, Aug 2009.
- [7] L. Xiao, “Dual averaging methods for regularized stochastic learning and online optimization,” *J. of Machine Learning Research*, pp. 2543–2596, Oct 2010.
- [8] N. Ferdinand, B. Gharachorloo, and S. C. Draper, “Anytime exploitation of stragglers in synchronous stochastic gradient descent,” in *Int. Conf. Machine Learning and Applications*, Cancun, Dec 2017.
- [9] L. Xiao, S. Boyd, and S.-J. Kim, “Distributed average consensus with least-mean-square deviation,” *J. of parallel and distributed computing*, pp. 33–46, Jan 2007.
- [10] H. Xiao, K. Rasul, and R. Vollgraf, “Fashion-MNIST: a novel image dataset for benchmarking machine learning algorithms,” *arXiv preprint arXiv:1708.07747*, Aug 2017.
- [11] S. Shalev-Shwartz and S. Ben-David, *Understanding machine learning: From theory to algorithms*. Cambridge University Press, 2014.
- [12] S. M. Kay, *Fundamentals of statistical signal processing*. Prentice Hall PTR, 1993.
- [13] J. Konečný, H. B. McMahan, F. X. Yu, P. Richtárik, A. T. Suresh, and D. Bacon, “Federated learning: Strategies for improving communication efficiency,” *arXiv preprint arXiv:1610.05492*, 2016.