
OPTIMAL DATA SPLITTING IN DISTRIBUTED OPTIMIZATION FOR MACHINE LEARNING

TECHNICAL REPORT

Gleb Molodtsov
MIPT, Russia
molodtsov.gl@phystech.edu

Daniil Medyakov
MIPT, Russia
mediakov.do@phystech.edu

Alexander Beznosikov
MIPT, Russia
anbeznosikov@gmail.com

ABSTRACT

The distributed optimization problem has become increasingly relevant recently. We consider this problem in the context of varying capacities of devices between which data is shared. The objective of this study is to achieve an optimal ratio of distributed data between the server and local machines. Optimal gradient descent and its application to distributed optimization under similarity are employed to address this problem. However, most distributed approaches suffer from a significant bottleneck - the cost of communications. The paper proposes a solution that takes into account this cost. The running times of the system are compared between uniform and optimal distributions. The superior theoretical performance of our solutions is experimentally validated.

1 Introduction

1.1 Distributed optimization

To achieve state-of-the-art performance in modern machine learning and minimization tasks, researchers and practitioners face various challenges. Training modern machine learning models remains an extremely difficult task. In order to improve the generalization of deployed models, machine learning engineers are compelled to rely on increasingly large training datasets. These datasets are typically collected in a distributed manner and stored across a network of edge devices, as is the case with federated learning. In the distributed training of complex models, the communication overhead often becomes the bottleneck of the training system. Motivated by the need for more efficient distributed learning methods, we consider optimization problems of the following form:

$$\min_{x \in \mathbb{R}^d} r(x) = \frac{1}{n} \sum_{i=1}^n f_i(x), \quad (1)$$

Here, $x \in \mathbb{R}^d$ collects the parameters of a statistical model to be trained, n is the number of devices, and f_i represents the convex loss-function of agent i , which is unknown to the other agents. Several solution methods have been proposed to solve 1. The prototype approach involves interleaving

edge devices calculations (nodes $i = 1, \dots, n$) with communications to and from the master node ($i = 1$). The master node maintains and updates the authoritative copy of the optimization variables, eventually producing the final solution estimate.

1.2 Distributed optimization under similarity

Since communication cost often becomes the bottleneck in distributed computing, significant research has focused on developing communication-efficient distributed algorithms. Acceleration, based on Nesterov's concept, has been extensively studied for reducing the communication burden. For L -smooth and μ -strongly convex functions r in 1, first-order methods guarantee linear convergence with computation and communication complexities proportional to $\sqrt{\kappa}$. Here, $\kappa := \frac{L}{\mu}$ represents the condition number of r . However, for ill-conditioned functions with a large κ , the polynomial dependence on κ may be unsatisfactory. This is often the case for many empirical risk minimization (ERM) problems where the optimal regularization parameter for test predictive performance is very small.

To further improve communication complexity, we can exploit the additional structure typically found in ERM problems, known as function similarity. Specifically, for all x in a suitable domain of interest and all $i \neq j = 1, \dots, n$, the difference between the Hessian matrices of local losses, denoted by $\|\nabla^2 f_i(x) - \nabla^2 f_j(x)\|$, is bounded by δ , where $\delta > 0$ measures the degree of similarity.

To achieve lower communication and local gradient complexity, we can refer to Algorithm 1 in [2]. For this purpose, the function r needs to be transformed into the following form:

$$r(x) = \underbrace{f_1(x)}_{q(x)} + \underbrace{\frac{1}{n} \sum_{i=1}^n [f_i(x) - f_1(x)]}_{p(x)}, \quad (2)$$

Here, r is assumed to be convex and decomposed as the sum of a smooth, potentially non-convex function p and a smooth convex function q . First-order information of p and q can be accessed separately. We are particularly interested in scenarios where evaluating the gradients of these two functions incurs heterogeneous costs, i.e., the cost of computing ∇p is significantly higher than that of computing ∇q .

1.3 Contributions

For this problem, as shown above, there already exist algorithms that achieve lower estimates, but they do not take into account in any way the fact that devices in the network may have different capacities, i.e., process a unit of information in different times. In this paper, we ask the question:

Can we find such a distribution of data among the devices in the network to reduce the actual running time of the optimal algorithm [2]?

We pay special attention to the limiting cases and obtain results in them. The case where communications are too expensive is not of practical interest as the whole idea of distributed learning is lost, but the case of inexpensive communications (not so expensive that the communication takes longer than processing all data by just one device) is of great interest. Actively using the algorithm and estimates from [2], we derive formulas for data distribution among devices in the network, both in general and in the limiting cases. We also conducted experiments confirming that with the obtained distribution it takes less time to solve the selected problem.

2 Problem Statement

Let's consider the accelerated extragradient algorithm (Algorithm 1 in [2]). Let us calculate how many operations this algorithm performs per iteration. In line 5, there is one communication, one

local computation, one central node calculation, and additional central node computations. In line 6, there is one communication, one local computation, and one central node calculation. Let us make the following notations: τ_i - the time of one local computation on i -th device, K - the number of iterations, τ_{comm} - the time for one communication, k_{some} - additional computations of the central node, n - the number of nodes in the network. Then we can write the general running time of the algorithm as:

$$T_{sum} = 2 \cdot \max(\tau_1, \tau_2, \dots, \tau_n) \cdot K + 2 \cdot K \cdot \tau_{comm} + \tau_1 \cdot k_{some} \quad (3)$$

Our task is to minimize the time T_{sum} . Let's represent the time τ_i as $\tau_i = \tau_i^{loc} \cdot b_i$, where τ_i^{loc} is the time spent by the i -th device to process a unit of information submitted to its input, and b_i is the size of dataset submitted to the i -th device. b_i must satisfy the following constraints: $\sum_{i=1}^n b_i = N$, where N is the size of the whole dataset, $\delta = \frac{L}{\sqrt{b_i}}$ or $\delta = \frac{L}{b_i}$ (this estimate is given in [2]). We obtained the following optimization problem:

$$\min_{\sum_{i=1}^n b_i = N; \delta = \frac{L}{b_1^\gamma}} [2 \cdot \max(\tau_1^{loc} \cdot b_1, \tau_2^{loc} \cdot b_2, \dots, \tau_n^{loc} \cdot b_n) \cdot K + 2 \cdot K \cdot \tau_{comm} + \tau_1 \cdot k_{some}], \gamma \in \{\frac{1}{2}, 1\} \quad (4)$$

3 Problem solution 4, case $\delta = \frac{L}{\sqrt{b_1}}$

3.1 The primary problem of minimization

In [2] the estimates of K and k_{some} are found, namely:

$$2 \cdot K = \mathcal{O}(\max\{1, \sqrt{\frac{L_p}{\mu}}\} \log(\frac{1}{\varepsilon})), \tau_1 \cdot k_{some} = \mathcal{O}(\max\{1, \sqrt{\frac{L_q}{L_p}}, \sqrt{\frac{L_p}{\mu}}, \sqrt{\frac{L_q}{\mu}}\} \log(\frac{1}{\varepsilon})).$$

The value of τ_{comm} will be determined later.

Thus, our minimization problem is reduced to:

$$\begin{aligned} \min_{\sum_{i=1}^n b_i = N; \delta = \frac{L}{b_1^\gamma}} & [(\max(\tau_1^{loc} \cdot b_1, \tau_2^{loc} \cdot b_2, \dots, \tau_n^{loc} \cdot b_n) + \tau_{comm}) * \mathcal{O}(\max\{1, \sqrt{\frac{L_p}{\mu}} \log(\frac{1}{\varepsilon})\}) \\ & + \mathcal{O}(\max\{1, \sqrt{\frac{L_q}{L_p}}, \sqrt{\frac{L_p}{\mu}}, \sqrt{\frac{L_q}{\mu}}\} \log(\frac{1}{\varepsilon}))], \gamma \in \{\frac{1}{2}, 1\} \end{aligned} \quad (5)$$

3.2 Auxiliary problem

Consider an auxiliary problem:

$$\min_{\sum_{i=2}^n b_i = N} [\max(\tau_2^{loc} \cdot b_2, \tau_3^{loc} \cdot b_3, \dots, \tau_n^{loc} \cdot b_n)] \quad (6)$$

Lemma 1. The solution of problem 6 is \vec{b} satisfying $\tau_2^{loc} \cdot b_2 = \tau_3^{loc} \cdot b_3 = \dots = \tau_n^{loc} \cdot b_n$

Proof. Without loss of generality, let us assume fixed values for $\tau_2^{loc} \leq \tau_3^{loc} \leq \dots \leq \tau_n^{loc}$.

Then let us arbitrarily choose $b_2 \geq b_3 \geq \dots \geq b_n$.

This is indeed the case, otherwise we would have a situation where $\exists i \neq j : i, j \in \{2, \dots, n\} : \max(\tau_i^{loc} \cdot b_i, \tau_j^{loc} \cdot b_j) > \max(\tau_i^{loc} \cdot b_j, \tau_j^{loc} \cdot b_i)$, and therefore the distribution would be suboptimal.

Our goal is to minimize the function $g(b) = \max(\tau_2^{loc} \cdot b_2, \tau_3^{loc} \cdot b_3, \dots, \tau_n^{loc} \cdot b_n)$.

Suppose that there exists a distribution such that $\exists i \in \{2, \dots, n\} : g(\vec{b}^0) = \tau_i^{loc} \cdot b_i^0$ is the minimum, and $\forall j : j \geq 2, j \neq i \hookrightarrow \tau_i^{loc} \cdot b_i^0 > \tau_j^{loc} \cdot b_j^0$.

It follows that $b_i^0 > \frac{\tau_j^{loc}}{\tau_i^{loc}} b_j^0 > \frac{\tau_{j_1}^{loc}}{\tau_i^{loc}} b_{j_1}^0 > \dots > \frac{\tau_{j_k}^{loc}}{\tau_i^{loc}} b_{j_k}^0$.

Then, considering $\sum_{i=2}^n b_i = N \hookrightarrow b_i^0 + \frac{\tau_j^{loc}}{\tau_i^{loc}} b_j^0 + \frac{\tau_{j_1}^{loc}}{\tau_i^{loc}} b_{j_1}^0 + \dots + \frac{\tau_{j_k}^{loc}}{\tau_i^{loc}} b_{j_k}^0 > N$, we obtain

$$b_i^0 > N(1 + \tau_i^{loc} \sum_{\substack{j=2 \\ j \neq i}}^n \frac{1}{\tau_j^{loc}})^{-1}.$$

Next, let us consider $b_i = N(1 + \tau_i^{loc} \sum_{\substack{j=2 \\ j \neq i}}^n \frac{1}{\tau_j^{loc}})^{-1}$, $b_j = \frac{\tau_i^{loc}}{\tau_j^{loc}} \cdot b_i \ \forall j \in \{2, \dots, n\}$. This distribution

yields a minimum of $g(\vec{b}) = \tau_i^{loc} \cdot b_i = \tau_j^{loc} \cdot b_j \ \forall j \in \{2, \dots, n\}$, and $g(\vec{b}) < g(\vec{b}^0)$. This contradicts the assumption of minimality.

Thus, for the distribution that minimizes the function $g(b) = \max(\tau_2^{loc} \cdot b_2, \tau_3^{loc} \cdot b_3, \dots, \tau_n^{loc} \cdot b_n)$, it holds that $\tau_2^{loc} \cdot b_2 = \tau_3^{loc} \cdot b_3 = \dots = \tau_n^{loc} \cdot b_n$.

□

Let us return to the problem 5. In addition to the minimum expression already studied in the problem 6, there are additional terms in the problem 5. Note that $L_q = L$, $L_p = \delta = \frac{L}{\sqrt{b_1}}$ or $L_p = \delta = \frac{L}{b_1}$ (this estimate is given in [2]). They depend on the value of b_1 , but do not depend on $b_i, i \in \overline{2, n}$. From this and Lemma 1, it follows that in the original problem 5, the data sharing between the 2nd, 3rd, and subsequent devices should be proportional. Thus, the problem 5 is reduced to a new problem with additional constraints:

$$\begin{aligned} \min_{\substack{\sum_{i=1}^n b_i = N; \delta = \frac{L}{b_1 \gamma}; \\ \tau_2^{loc} \cdot b_2 = \dots = \tau_n^{loc} \cdot b_n}} & [(\max(\tau_1^{loc} \cdot b_1, \tau_2^{loc} \cdot b_2, \dots, \tau_n^{loc} \cdot b_n) + \tau_{comm}) \cdot \mathcal{O}(\max\{1, \sqrt{\frac{L_p}{\mu}} \log(\frac{1}{\varepsilon})\})] \quad (7) \\ & + \mathcal{O}(\max\{1, \sqrt{\frac{L_q}{L_p}}, \sqrt{\frac{L_p}{\mu}}, \sqrt{\frac{L_q}{\mu}}\} \log(\frac{1}{\varepsilon})), \gamma \in \{\frac{1}{2}, 1\} \end{aligned}$$

3.3 Define the final minimization problem

It follows from Lemma 1 that $b_i \tau_i^{loc} = \text{const} \ \forall i \in \overline{2, n}$ due to the symmetry of the problem. Therefore,

$$N - b_1 = \sum_{i=2}^n b_i = \sum_{i=2}^n \frac{\tau_2^{loc} \cdot b_2}{\tau_i^{loc}} = \tau_2^{loc} \cdot b_2 \cdot \sum_{i=2}^n \frac{1}{\tau_i^{loc}} \Rightarrow b_2 = \frac{N - b_1}{\tau_2^{loc}} \left(\sum_{i=2}^n \frac{1}{\tau_i^{loc}} \right)^{-1}$$

3.3.1 $\delta = \frac{L}{b_1}$

In the first case the following relations are fulfilled:

$$L_p = \delta, L_q = L, \mu \leq \delta \leq L \Rightarrow \begin{cases} 2 \cdot K = \mathcal{O}(\sqrt{\frac{L_p}{\mu}} \log(\frac{1}{\varepsilon})) = \mathcal{O}(\sqrt{\frac{L}{\mu b_1}} \log(\frac{1}{\varepsilon})) \\ \tau_1 \cdot k_{some} = \mathcal{O}(\sqrt{\frac{L}{\mu}} \log(\frac{1}{\varepsilon})) \end{cases}.$$

Then the problem will take the following form:

$$\min_{\sum_{i=1}^n b_i = N} [(\max\{\tau_1^{loc} \cdot b_1, \tau_2^{loc} \cdot b_2\} + \tau_{comm}) \cdot \mathcal{O}(\sqrt{\frac{L}{\mu b_1}} \log(\frac{1}{\varepsilon})) + \tau_1^{loc} \cdot b_1 \cdot \mathcal{O}(\sqrt{\frac{L}{\mu}} \log(\frac{1}{\varepsilon}))]$$

Consider the final form of the minimization problem:

$$\min_{0 < b_1 \leq N} [(\max\{\tau_1^{loc} \cdot b_1; (N - b_1) \cdot (\sum_{i=2}^n \frac{1}{\tau_i^{loc}})^{-1}\} + \tau_{comm}) \cdot \mathcal{O}(\sqrt{\frac{L}{\mu b_1}} \log(\frac{1}{\varepsilon})) + \tau_1^{loc} \cdot b_1 \cdot \mathcal{O}(\sqrt{\frac{L}{\mu}} \log(\frac{1}{\varepsilon}))] \quad (8)$$

Let us investigate the problem further. To do this, find the point at which the expressions under the maximum coincide.

$$b_1^0 \cdot (\tau_1^{loc} + (\sum_{i=2}^n \frac{1}{\tau_i^{loc}})^{-1}) = N(\sum_{i=2}^n \frac{1}{\tau_i^{loc}})^{-1} \Rightarrow b_1^0 = \frac{N(\sum_{i=2}^n \frac{1}{\tau_i^{loc}})^{-1}}{\tau_1^{loc} + (\sum_{i=2}^n \frac{1}{\tau_i^{loc}})^{-1}}$$

Thus, we obtained two half-intervals, on each of which we can formulate a different minimization problem:

$$\begin{cases} (a) & 0 < b_1 \leq b_1^0 \Rightarrow \max\{\tau_1^{loc} \cdot b_1; (N - b_1) \cdot (\sum_{i=2}^n \frac{1}{\tau_i^{loc}})^{-1}\} = (N - b_1) \cdot (\sum_{i=2}^n \frac{1}{\tau_i^{loc}})^{-1} \\ (b) & b_1^0 < b_1 \leq N \Rightarrow \max\{\tau_1^{loc} \cdot b_1; (N - b_1) \cdot (\sum_{i=2}^n \frac{1}{\tau_i^{loc}})^{-1}\} = \tau_1^{loc} \cdot b_1 \end{cases} \quad (9)$$

$$(a) : \mathcal{F}_1(b_1) = [N(\sum_{i=2}^n \frac{1}{\tau_i^{loc}})^{-1} + \tau_{comm}] \cdot c_1 \sqrt{\frac{L}{\mu}} \log(\frac{1}{\varepsilon}) b_1^{-\frac{1}{2}} - c_1 \sqrt{\frac{L}{\mu}} \log(\frac{1}{\varepsilon}) (\sum_{i=2}^n \frac{1}{\tau_i^{loc}})^{-1} b_1^{\frac{1}{2}} + \tau_1^{loc} \cdot c_2 \sqrt{\frac{L}{\mu}} \log(\frac{1}{\varepsilon}) b_1$$

$$(b) : \mathcal{F}_2(b_1) = \tau_{comm} \cdot c_1 \sqrt{\frac{L}{\mu}} \log(\frac{1}{\varepsilon}) b_1^{-\frac{1}{2}} + c_1 \sqrt{\frac{L}{\mu}} \log(\frac{1}{\varepsilon}) \tau_1^{loc} b_1^{\frac{1}{2}} + \tau_1^{loc} \cdot c_2 \sqrt{\frac{L}{\mu}} \log(\frac{1}{\varepsilon}) b_1$$

$$(a) : \mathcal{F}'_1(b_1) = -\frac{1}{2} c_1 b_1^{-\frac{3}{2}} [N(\sum_{i=2}^n \frac{1}{\tau_i^{loc}})^{-1} + \tau_{comm}] \cdot \sqrt{\frac{L}{\mu}} \log(\frac{1}{\varepsilon}) - \frac{1}{2} c_1 b_1^{-\frac{1}{2}} \sqrt{\frac{L}{\mu}} \log(\frac{1}{\varepsilon}) (\sum_{i=2}^n \frac{1}{\tau_i^{loc}})^{-1} + \tau_1^{loc} \cdot c_2 \sqrt{\frac{L}{\mu}} \log(\frac{1}{\varepsilon})$$

$$(b) : \mathcal{F}'_2(b_1) = -\frac{1}{2} c_1 b_1^{-\frac{3}{2}} \tau_{comm} \cdot \sqrt{\frac{L}{\mu}} \log(\frac{1}{\varepsilon}) + \frac{1}{2} c_1 b_1^{-\frac{1}{2}} \sqrt{\frac{L}{\mu}} \log(\frac{1}{\varepsilon}) \tau_1^{loc} + \tau_1^{loc} \cdot c_2 \sqrt{\frac{L}{\mu}} \log(\frac{1}{\varepsilon})$$

3.3.2 $\delta = \frac{L}{\sqrt{b_1}}$

Let us perform similar transformations for this case. Our relations turn into:

$$L_p = \delta, L_q = L, \mu \leq \delta \leq L \Rightarrow \begin{cases} 2 \cdot K = \mathcal{O}(\sqrt{\frac{L_p}{\mu}} \log(\frac{1}{\varepsilon})) = \mathcal{O}(\sqrt{\frac{L}{\mu \sqrt{b_1}}} \log(\frac{1}{\varepsilon})) \\ \tau_1 \cdot k_{some} = \mathcal{O}(\sqrt{\frac{L}{\mu}} \log(\frac{1}{\varepsilon})) \end{cases}$$

And the main problem will take the following form:

$$\min_{\sum_{i=1}^n b_i = N} [(\max\{\tau_1^{loc} \cdot b_1, \tau_2^{loc} \cdot b_2\} + \tau_{comm}) \cdot \mathcal{O}(\sqrt{\frac{L}{\mu \sqrt{b_1}}} \log(\frac{1}{\varepsilon})) + \tau_1^{loc} \cdot b_1 \cdot \mathcal{O}(\sqrt{\frac{L}{\mu}} \log(\frac{1}{\varepsilon}))]$$

Then write down the final minimization problem:

$$\min_{0 < b_1 \leq N} [(\max\{\tau_1^{loc} \cdot b_1; (N - b_1) \cdot (\sum_{i=2}^n \frac{1}{\tau_i^{loc}})^{-1}\} + \tau_{comm}) \cdot \mathcal{O}(\sqrt{\frac{L}{\mu \sqrt{b_1}}} \log(\frac{1}{\varepsilon})) + \tau_1^{loc} \cdot b_1 \cdot \mathcal{O}(\sqrt{\frac{L}{\mu}} \log(\frac{1}{\varepsilon}))] \quad (10)$$

Let us investigate the problem further. To do this, find the point at which the expressions under the maximum coincide.

$$b_1^0 \cdot (\tau_1^{loc} + (\sum_{i=2}^n \frac{1}{\tau_i^{loc}})^{-1}) = N(\sum_{i=2}^n \frac{1}{\tau_i^{loc}})^{-1} \Rightarrow b_1^0 = \frac{N(\sum_{i=2}^n \frac{1}{\tau_i^{loc}})^{-1}}{\tau_1^{loc} + (\sum_{i=2}^n \frac{1}{\tau_i^{loc}})^{-1}}$$

Thus, we obtained two half-intervals, on each of which we can formulate a different minimization problem:

$$\begin{cases} (a) \ 0 < b_1 \leq b_1^0 \Rightarrow \max\{\tau_1^{loc} \cdot b_1; (N - b_1) \cdot (\sum_{i=2}^n \frac{1}{\tau_i^{loc}})^{-1}\} = (N - b_1) \cdot (\sum_{i=2}^n \frac{1}{\tau_i^{loc}})^{-1} \\ (b) \ b_1^0 < b_1 \leq N \Rightarrow \max\{\tau_1^{loc} \cdot b_1; (N - b_1) \cdot (\sum_{i=2}^n \frac{1}{\tau_i^{loc}})^{-1}\} = \tau_1^{loc} \cdot b_1 \end{cases} \quad (11)$$

$$\begin{aligned} (a) : \mathcal{F}_1(b_1) &= [N(\sum_{i=2}^n \frac{1}{\tau_i^{loc}})^{-1} + \tau_{comm}] \cdot c_1 \sqrt{\frac{L}{\mu}} \log(\frac{1}{\varepsilon}) b_1^{-\frac{1}{4}} - c_1 \sqrt{\frac{L}{\mu}} \log(\frac{1}{\varepsilon}) (\sum_{i=2}^n \frac{1}{\tau_i^{loc}})^{-1} b_1^{\frac{3}{4}} + \\ &\tau_1^{loc} \cdot c_2 \sqrt{\frac{L}{\mu}} \log(\frac{1}{\varepsilon}) b_1 \\ (b) : \mathcal{F}_2(b_1) &= \tau_{comm} \cdot c_1 \sqrt{\frac{L}{\mu}} \log(\frac{1}{\varepsilon}) b_1^{-\frac{1}{4}} + c_1 \sqrt{\frac{L}{\mu}} \log(\frac{1}{\varepsilon}) \tau_1^{loc} b_1^{\frac{3}{4}} + \tau_1^{loc} \cdot c_2 \sqrt{\frac{L}{\mu}} \log(\frac{1}{\varepsilon}) b_1 \end{aligned}$$

Let's find the derivatives of the functions:

$$\begin{aligned} (a) : \mathcal{F}'_1(b_1) &= -\frac{1}{4} c_1 b_1^{-\frac{5}{4}} [N(\sum_{i=2}^n \frac{1}{\tau_i^{loc}})^{-1} + \tau_{comm}] \cdot \sqrt{\frac{L}{\mu}} \log(\frac{1}{\varepsilon}) - \frac{3}{4} c_1 b_1^{-\frac{1}{4}} \sqrt{\frac{L}{\mu}} \log(\frac{1}{\varepsilon}) (\sum_{i=2}^n \frac{1}{\tau_i^{loc}})^{-1} + \\ &\tau_1^{loc} \cdot c_2 \sqrt{\frac{L}{\mu}} \log(\frac{1}{\varepsilon}) \\ (b) : \mathcal{F}'_2(b_1) &= -\frac{1}{4} c_1 b_1^{-\frac{5}{4}} \tau_{comm} \cdot \sqrt{\frac{L}{\mu}} \log(\frac{1}{\varepsilon}) + \frac{3}{4} c_1 b_1^{-\frac{1}{4}} \sqrt{\frac{L}{\mu}} \log(\frac{1}{\varepsilon}) \tau_1^{loc} + \tau_1^{loc} \cdot c_2 \sqrt{\frac{L}{\mu}} \log(\frac{1}{\varepsilon}) \end{aligned}$$

It is impossible to write out the solution of these equations in analytical form, because of their degrees, so we will consider the limiting cases.

3.4 Final solution in limiting cases

3.4.1 $\delta = \frac{L}{b_1}$

To solve the resulting cubic equation, we can use the Cardano formula. Consider the equation $ax^{-\frac{1}{2}} + bx^{-\frac{3}{2}} + c = 0$,

where in cases (a) : $0 < b_1 \leq b_1^0$ and (b) : $b_1^0 < b_1 \leq N$ we assume:

$$\begin{cases} (a) : a = \frac{1}{2} c_1 \sqrt{\frac{L}{\mu}} \log(\frac{1}{\varepsilon}) (\sum_{i=2}^n \frac{1}{\tau_i^{loc}})^{-1}; \ b = -\frac{1}{2} c_1 [N(\sum_{i=2}^n \frac{1}{\tau_i^{loc}})^{-1} + \tau_{comm}] \cdot \sqrt{\frac{L}{\mu}} \log(\frac{1}{\varepsilon}); \ c = \tau_1^{loc} \cdot c_2 \sqrt{\frac{L}{\mu}} \log(\frac{1}{\varepsilon}) \\ (b) : a = \frac{1}{2} c_1 \sqrt{\frac{L}{\mu}} \log(\frac{1}{\varepsilon}) \tau_1^{loc}; \ b = -\frac{1}{2} c_1 \tau_{comm} \cdot \sqrt{\frac{L}{\mu}} \log(\frac{1}{\varepsilon}); \ c = \tau_1^{loc} \cdot c_2 \sqrt{\frac{L}{\mu}} \log(\frac{1}{\varepsilon}) \end{cases}$$

Then on the condition that

$$N \geq \frac{a^2}{3c^2} + \frac{\sqrt[3]{2a^6 + 3\sqrt{3}\sqrt{4a^3b^3c^6 + 27b^4c^8 + 18a^3bc^2 + 27b^2c^4}}}{3\sqrt[3]{2}c^2} - \frac{\sqrt[3]{2}(-a^4 - 6abc^2)}{3c^2 \sqrt[3]{2a^6 + 3\sqrt{3}\sqrt{4a^3b^3c^6 + 27b^4c^8 + 18a^3bc^2 + 27b^2c^4}}},$$

We get a solution:

$$x = \frac{a^2}{3c^2} + \frac{\sqrt[3]{2a^6 + 3\sqrt{3}\sqrt{4a^3b^3c^6 + 27b^4c^8 + 18a^3bc^2 + 27b^2c^4}}}{3\sqrt[3]{2}c^2} - \frac{\sqrt[3]{2}(-a^4 - 6abc^2)}{3c^2\sqrt[3]{2a^6 + 3\sqrt{3}\sqrt{4a^3b^3c^6 + 27b^4c^8 + 18a^3bc^2 + 27b^2c^4}}}.$$

3.4.2 $\delta = \frac{L}{\sqrt{b_1}}$

Let us find an analytical solution to the problem in two particular cases:

1. $\forall i \hookrightarrow \tau_{comm} \ll \tau_i^{loc}$
2. $\forall i \hookrightarrow \tau_{comm} \gg \tau_i^{loc}, \forall i \neq j \hookrightarrow \tau_i^{loc} = \tau_j^{loc}$

Establish $\alpha = c_1 \cdot \sqrt{\frac{L}{\mu}} \cdot \log(\frac{1}{\varepsilon}), \beta = c_2 \cdot \sqrt{\frac{L}{\mu}} \cdot \log(\frac{1}{\varepsilon})$

Consider case 1.

a) $0 < b_1 \leq b_1^0$

$$\mathcal{F}_1(b_1) = [N(\sum_{i=2}^n \frac{1}{\tau_i^{loc}})^{-1} + \tau_{comm}] \cdot \alpha b_1^{-\frac{1}{4}} - \alpha(\sum_{i=2}^n \frac{1}{\tau_i^{loc}})^{-1} b_1^{\frac{3}{4}} + \tau_1^{loc} \cdot \beta b_1$$

Consider

$$\tau_1^{loc} \leq \tau_2^{loc} \leq \dots \leq \tau_n^{loc} \quad (12)$$

$$\begin{aligned} (\sum_{i=2}^n \frac{1}{\tau_i^{loc}})^{-1} &= \frac{1}{\frac{1}{\tau_1^{loc}} + \dots + \frac{1}{\tau_n^{loc}}} = \frac{\tau_2^{loc} \dots \tau_n^{loc}}{\tau_3^{loc} \dots \tau_n^{loc} + \tau_2^{loc} \cdot \tau_4^{loc} \dots \tau_n^{loc} + \dots + \tau_2^{loc} \cdot \dots \cdot \tau_{n-1}^{loc}} \\ &\geq \frac{\tau_2^{loc}}{12} >> \tau_{comm} \end{aligned} \quad (13)$$

Then taking into account 13:

$$\mathcal{F}_1(b_1) = \alpha(\sum_{i=2}^n \frac{1}{\tau_i^{loc}})^{-1} \cdot b_1^{-\frac{1}{4}} (N - b_1) + \tau_1^{loc} \beta \cdot b_1$$

$$\mathcal{F}'_1(b_1) = \alpha(\sum_{i=2}^n \frac{1}{\tau_i^{loc}})^{-1} \cdot (-\frac{1}{4} b_1^{-\frac{5}{4}} N - \frac{3}{4} b_1^{-\frac{1}{4}}) + \tau_1^{loc} \beta$$

The analytical solution in this case is not given.

b) $b_1^0 \leq b_1 \leq N$

$$\mathcal{F}_2(b_1) = \tau_{comm} \cdot \alpha b_1^{-\frac{1}{4}} + \alpha \tau_1^{loc} b_1^{\frac{3}{4}} + \tau_1^{loc} \cdot \beta b_1 = \alpha \cdot b_1^{-\frac{1}{4}} (\tau_{comm} + \tau_1^{loc} b_1) + \beta \cdot \tau_1^{loc} \cdot b_1$$

$$\tau_1^{loc} b_1 \geq_{b_1 \geq b_1^0, 12} \frac{\tau_1^{loc} N \frac{\tau_2^{loc}}{n-1}}{\tau_1^{loc} + \frac{\tau_n^{loc}}{n-1}} \geq \frac{\tau_1^{loc} \tau_2^{loc} N}{(n-1)(\tau_1^{loc} + \tau_n^{loc})} \geq \frac{\tau_1^{loc} \tau_2^{loc} N}{2(n-1)\tau_n^{loc}} >> \tau_{comm} \frac{N}{2(n-1)} >> \tau_{comm} \quad (14)$$

Then taking into account 14:

$$\mathcal{F}_2(b_1) = \alpha \cdot \tau_1^{loc} \cdot b_1^{\frac{3}{4}} + \beta \tau_1^{loc} \cdot b_1$$

$$\mathcal{F}'_2(b_1) = \frac{3}{4} \alpha \cdot \tau_1^{loc} \cdot b_1^{-\frac{1}{4}} + \beta \cdot \tau_1^{loc} > 0$$

Since the derivative of the function is positive, the function is increasing, and therefore the

$$\text{minimum will be taken at } b_1 = b_1^0 = \frac{N(\sum_{i=2}^n \frac{1}{\tau_i^{loc}})^{-1}}{\tau_1^{loc} + (\sum_{i=2}^n \frac{1}{\tau_i^{loc}})^{-1}}$$

Thus, in the case of small τ_{comm} we obtained the following result:

$$b_{1\min} \leq b_1^0 = \frac{N(\sum_{i=2}^n \frac{1}{\tau_i^{loc}})^{-1}}{\tau_1^{loc} + (\sum_{i=2}^n \frac{1}{\tau_i^{loc}})^{-1}}.$$

Consider case 2.

Establish: $\tau := \tau_i^{loc} \forall i \in 1, \dots, n$.

$$\text{Then } \mathcal{F} = (\max\{\tau b_1; (N - b_1) \frac{\tau}{n-1}\} + \tau_{comm}) \cdot \frac{\alpha}{\sqrt[4]{b_1}} + \tau \beta b_1 \quad (15)$$

Consider the case $\tau_{comm} = N^2 \tau$. N can be considered large, so $\tau_{comm} \gg N\tau$. Then:

$$\max\{\tau b_1; (N - b_1) \frac{\tau}{n-1}\} < \tau N \ll \tau_{comm} \Rightarrow \mathcal{F} \approx \frac{\alpha \tau_{comm}}{\sqrt[4]{b_1}} + \beta \tau b_1$$

$$\mathcal{F}'(b_1) = -\frac{\alpha \tau_{comm}}{4b_1 \sqrt[4]{b_1}} + \beta \tau = 0 \Rightarrow b_{1\min}^{\frac{5}{4}} = \frac{\tau_{comm} \alpha}{4\beta \tau} \Rightarrow b_{1\min} = (\frac{\tau_{comm} \alpha}{4\beta \tau})^{\frac{4}{5}}$$

Assuming that the found value b_1 lies on the interval $(0, N)$, that is, at $0 < (\frac{\tau_{comm} \alpha}{4\beta \tau})^{\frac{4}{5}} < N$, it will be the point of minimum function \mathcal{F} . Then:

$$\mathcal{F}(b_{1\min}) = (\alpha \tau_{comm})^{\frac{4}{5}} \cdot (4\beta \tau)^{\frac{1}{5}} + (\beta \tau)^{\frac{1}{5}} \cdot (\frac{\alpha \tau_{comm}}{4})^{\frac{4}{5}} = (\alpha \tau_{comm})^{\frac{4}{5}} \cdot (4\beta \tau)^{\frac{1}{5}} (4^{\frac{1}{5}} + 4^{-\frac{4}{5}}).$$

Otherwise, the minimum will be reached at the right boundary, since at zero we can say that the function is increasing.

Summarizing all of the above in this case, it is worth noting that for very large values of N the second special case generalizes to the following condition:

$$\forall i \hookrightarrow \tau_{comm} = \mathcal{O}(N^k \tau_i^{loc}), k > 1, \forall i \neq j \hookrightarrow \tau_i^{loc} = \tau_j^{loc} \quad (16)$$

$$\min \mathcal{F}(b_1) = \begin{cases} (\alpha \tau_{comm})^{\frac{4}{5}} \cdot (4\beta \tau)^{\frac{1}{5}} (4^{\frac{1}{5}} + 4^{-\frac{4}{5}}), & 0 < (\frac{\tau_{comm} \alpha}{4\beta \tau})^{\frac{4}{5}} < N \\ \frac{\alpha \tau_{comm}}{N} + \beta \tau N, & (\frac{\tau_{comm} \alpha}{4\beta \tau})^{\frac{4}{5}} \geq N \end{cases} \quad (17)$$

3.5 Practical solution

In order to determine the minimum of these functions on the respective half-intervals, we will examine points where the derivatives of $\mathcal{F}'_1(b_1)$ and $\mathcal{F}'_2(b_1)$ approach zero. It should be noted that, given the nature of these functions, their derivatives can only be zero once on the desired half-interval. Hence, by employing basic methods, we can locate the zeros of the derivatives. Subsequently, we need to compare the values of the corresponding function at these points with the value at the extreme point of the interval. One of these points will provide the minimum solution, thereby serving as the ultimate solution to the problem 8 and 10.

4 Experiments

4.1 Description of experiments

For experimental verification of the theoretical results we consider the problem "Ridge Regression":

$$\min_{\omega} [\frac{1}{2N} X\omega - y^2 + \frac{\lambda}{2} \omega^2], q(\omega) = \frac{1}{2N} X\omega - y^2, p(\omega) = \frac{\lambda}{2} \omega^2 \quad (18)$$

The file a9a.txt with number of lines $N = 97683$ was chosen as dataset. The first step was the implementation of algorithm 1 of [2].

The iterative OGM-G method of [1] was applied to find the solution of line 5 of Algorithm 1. After calculating the required number of iterations to achieve a certain accuracy, the running time was measured.

From the obtained convergence graph we found the number of iterations to achieve the given accuracy, the values of constants c_1, c_2 , and, respectively, α, β . With their help, we were able to distribute the data from the dataset to the devices according to the above formulas.

Next, we ran the algorithm and measured the running time on the resulting distribution of data across devices and uniform distribution. The cases of large and small communications were considered.

In the end, two problems were considered:

1. $\delta = \frac{L}{\sqrt{b_1}}$
2. $\delta = \frac{L}{b_1}$

For problem 1, following cases were considered:

1. small communications (3.4.2)
2. large communications (3.4.2)
3. search for optimal allocation using Python optimization tools (3.5)

For problem 2, following cases were considered:

1. search for the optimal solution using the Cardano formula (3.4.1)
2. search for optimal allocation using Python optimization tools (3.5)

For all cases, acceleration was found and graphs were plotted 1

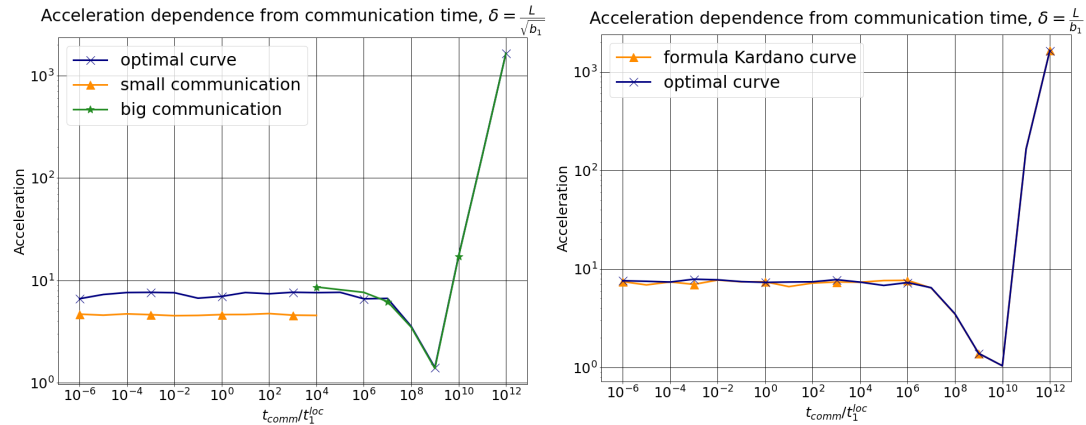


Figure 1: Final results

4.2 Analyze

Let us analyze the obtained graphs. The formula for the case of large communications and the Cardano formula practically coincided with the optimal solution search. The case of small communications showed worse results. This is explained by the fact that the formula was obtained in rough approximation. But if we take into account the constants α, β , we can get a better result, which is shown below.

$$F = \left(\max \left\{ \tau_1^{loc} \cdot b_1; (N - b_1) \cdot \left(\sum_{i=2}^n \frac{1}{\tau_i^{loc}} \right)^{-1} \right\} + \tau_{comm} \right) \cdot \frac{\alpha}{\sqrt[4]{b_1}} + \tau_1^{loc} b_1 \cdot \beta,$$

It has already been evaluated that $b_1 \leq b_1^0 \Rightarrow F = (N - b_1) \cdot \left(\sum_{i=2}^n \frac{1}{\tau_i^{loc}} \right)^{-1} \cdot \frac{\alpha}{\sqrt[4]{b_1}} + \tau_1^{loc} b_1 \cdot \beta$

$$F = N \left(\sum_{i=2}^n \frac{1}{\tau_i^{loc}} \right)^{-1} \alpha b^{-\frac{1}{4}} - \left(\sum_{i=2}^n \frac{1}{\tau_i^{loc}} \right)^{-1} \alpha b_1^{\frac{3}{4}} + \tau_1^{loc} \beta b_1$$

Consider that $\alpha \sim 10^6, \beta \sim 10^9 \Rightarrow$

$$F \cong 10^6 N \left(\sum_{i=2}^n \frac{1}{\tau_i^{loc}} \right)^{-1} \cdot b^{-\frac{1}{4}} - 10^6 \left(\sum_{i=2}^n \frac{1}{\tau_i^{loc}} \right)^{-1} \cdot b_1^{\frac{3}{4}} + 10^9 \tau_1^{loc} b_1$$

$$\frac{1}{4} 10^6 N \left(\sum_{i=2}^n \frac{1}{\tau_i^{loc}} \right)^{-1} \cdot b_1^{-\frac{5}{4}} \leq 10^5 \tau_1^{loc} \Rightarrow b_1 \leq \frac{4 \cdot 10^3 \tau_1^{loc}}{N \left(\sum_{i=2}^n (\tau_i^{loc})^{-1} \right)^{-1}}$$

5 Conclusion

In this paper we presented a new way to partition the data for the distributed optimization problem. Our solution is based on separating convex and non-convex subproblems and applying the accelerated extragradient algorithm from [2] as well as the OGM-G algorithm from [1]. Our method works well on star topology networks with various communication costs between the server and the local machines. The theoretical results have been confirmed experimentally. This indicates that our method gives acceleration on tasks of this type.

References

- [1] Donghwan Kim and Jeffrey A Fessler. Optimizing the efficiency of first-order methods for decreasing the gradient of smooth convex functions. *Journal of optimization theory and applications*, 188(1):192–219, 2021.
- [2] Dmitry Kovalev, Aleksandr Beznosikov, Ekaterina Borodich, Alexander Gasnikov, and Gesualdo Scutari. Optimal gradient sliding and its application to optimal distributed optimization under similarity. *Advances in Neural Information Processing Systems*, 35:33494–33507, 2022.