

DET-CGD: COMPRESSED GRADIENT DESCENT WITH MATRIX STEPSIZES FOR NON-CONVEX OPTIMIZATION

Hanmin Li Avetik Karagulyan Peter Richtárik

King Abdullah University of Science and Technology

{hanmin.li, avetik.karagulyan, peter.richtarik}@kaust.edu.sa

ABSTRACT

This paper introduces a new method for minimizing matrix-smooth non-convex objectives through the use of novel Compressed Gradient Descent (CGD) algorithms enhanced with a matrix-valued stepsize. The proposed algorithms are theoretically analyzed first in the single-node and subsequently in the distributed settings. Our theoretical results reveal that the matrix stepsize in CGD can capture the objective’s structure and lead to faster convergence compared to a scalar stepsize. As a byproduct of our general results, we emphasize the importance of selecting the compression mechanism and the matrix stepsize in a layer-wise manner, taking advantage of model structure. Moreover, we provide theoretical guarantees for free compression, by designing specific layer-wise compressors for the non-convex matrix smooth objectives. Our findings are supported with empirical evidence.

1 INTRODUCTION

The minimization of smooth and non-convex functions is a fundamental problem in various domains of applied mathematics. Most machine learning algorithms rely on solving optimization problems for training and inference, often with structural constraints or non-convex objectives to accurately capture the learning and prediction problems in high-dimensional or non-linear spaces. However, non-convex problems are typically NP-hard to solve, leading to the popular approach of relaxing them to convex problems and using traditional methods. Direct approaches to non-convex optimization have shown success but their convergence and properties are not well understood, making them challenging for large scale optimization. While its convex alternative has been extensively studied and is generally an easier problem, the non-convex setting is of greater practical interest often being the computational bottleneck in many applications.

In this paper, we consider the general minimization problem:

$$\min_{x \in \mathbb{R}^d} f(x), \tag{1}$$

where $f : \mathbb{R}^d \rightarrow \mathbb{R}$ is a differentiable function. In order for this problem to have a finite solution we will assume throughout the paper that f is bounded from below.

Assumption 1. *There exists $f^{\text{inf}} \in \mathbb{R}$ such that $f(x) \geq f^{\text{inf}}$ for all $x \in \mathbb{R}^d$.*

The stochastic gradient descent (SGD) algorithm (Moulines & Bach, 2011; Bubeck et al., 2015; Gower et al., 2019) is one of the most common algorithms to solve this problem. In its most general form, it can be written as

$$x^{k+1} = x^k - \gamma g(x^k), \tag{2}$$

where $g(x^k)$ is a stochastic estimator of $\nabla f(x^k)$ and $\gamma > 0$ is a positive scalar stepsize. A particular case of interest is the compressed gradient descent (CGD) algorithm (Khirirat et al., 2018), where the estimator g is taken as a compressed alternative of the initial gradient:

$$g(x^k) = \mathcal{C}(\nabla f(x^k)), \tag{3}$$

and the compressor \mathcal{C} is chosen to be a "sparser" estimator that aims to reduce the communication overhead in distributed or federated settings. This is crucial, as highlighted in the seminal paper by

Konečný et al. (2016), which showed that the bottleneck of distributed optimization algorithms is the communication complexity. In order to deal with the limited resources of current devices, there are various compression objectives that are practical to achieve. These include also compressing the model broadcasted from server to clients for local training, and reducing the computational burden of local training. These objectives are mostly complementary, but compressing gradients has the potential for the greatest practical impact due to slower upload speeds of client connections and the benefits of averaging Kairouz et al. (2021). In this paper we will focus on this latter problem.

An important subclass of compressors are the sketches. Sketches are linear operators defined on \mathbb{R}^d , i.e., $\mathcal{C}(y) = \mathbf{S}y$ for every $y \in \mathbb{R}^d$, where \mathbf{S} is a random matrix. A standard example of such a compressor is the Rand- k compressor, which randomly chooses k entries of its argument and scales them with a scalar multiplier to make the estimator unbiased. Instead of communicating all d coordinates of the gradient, one communicates only a subset of size k , thus reducing the number of communicated bits by a factor of d/k . Formally, Rand- k is defined as follows: $\mathbf{S} = \frac{d}{k} \sum_{j=1}^k e_{i_j} e_{i_j}^\top$, where e_{i_j} is the i_j -th standard basis vector in \mathbb{R}^d . We refer the reader to (Safaryan et al., 2022) for an overview on compressions.

Besides the assumption that function f is bounded from below, we also assume that it is \mathbf{L} matrix smooth, as we are trying to take advantage of the entire information contained in the smoothness matrix \mathbf{L} and the stepsize matrix \mathbf{D} .

Assumption 2 (Matrix smoothness). *There exists $\mathbf{L} \in \mathbb{S}_+^d$ such that*

$$f(x) \leq f(y) + \langle \nabla f(y), x - y \rangle + \frac{1}{2} \langle \mathbf{L}(x - y), x - y \rangle \quad (4)$$

holds for all $x, y \in \mathbb{R}^d$.

The assumption of matrix smoothness, which is a generalization of scalar smoothness, has been shown to be a more powerful tool for improving supervised model training. In Safaryan et al. (2021), the authors proposed using smoothness matrices and suggested a novel communication sparsification strategy to reduce communication complexity in distributed optimization for convex objectives. The technique was adapted to three distributed optimization algorithms in the convex setting, resulting in significant communication complexity savings and consistently outperforming the baselines. The results of this study demonstrate the efficacy of the matrix smoothness assumption in improving distributed optimization algorithms.

The case of block-diagonal smoothness matrices is particularly relevant in various applications, such as neural networks (NN). In this setting, each block corresponds to a layer of the network, and we characterize the smoothness with respect to nodes in the i -th layer by a corresponding matrix \mathbf{L}_i . Unlike in the scalar setting, we favor the similarity of certain entries of the argument over the others. This is because the information carried by the layers becomes more complex, while the nodes in the same layers are similar. This phenomenon has been observed visually in various studies, such as those by Yosinski et al. (2015) and Zintgraf et al. (2017).

We study two matrix stepsized CGD-type algorithms and analyze their convergence properties for non-convex matrix-smooth functions. As mentioned earlier, we put special emphasis on the block-diagonal case. We design our sketches and stepsizes in a way that leverages this structure, and we show that in certain cases, we can achieve compression without losing in the overall communication complexity.

1.1 RELATED WORK

Many successful convex optimization techniques have been adapted for use in the non-convex setting. Here is a non-exhaustive list: adaptivity (Dvinskikh et al., 2019; Zhang et al., 2020), variance reduction (J Reddi et al., 2016; Li et al., 2021), and acceleration (Guminov et al., 2019). A paper of particular importance for our work is that of Khaled & Richtárik (2020), which proposes a unified scheme for analyzing stochastic gradient descent in the non-convex regime. A comprehensive overview of non-convex optimization can be found in (Jain et al., 2017; Danilova et al., 2022).

A classical example of a matrix stepsized method is Newton’s method. This method has been popular in the optimization community for a long time (Gragg & Tapia, 1974; Miel, 1980; Yamamoto, 1987).

However, computing the stepsize as the inverse Hessian of the current iteration results in significant computational complexity. Instead, quasi-Newton methods use an easily computable estimator to replace the inverse of the Hessian (Broyden, 1965; Dennis & Moré, 1977; Al-Baali & Khalfan, 2007; Al-Baali et al., 2014). An example is the Newton-Star algorithm (Islamov et al., 2021), which we discuss in Section 2.

Gower & Richtárik (2015) analyzed sketched gradient descent by making the compressors unbiased with a sketch-and-project trick. They provided an analysis of the resulting algorithm for the linear feasibility problem. Later, Hanzely et al. (2018) proposed a variance-reduced version of this method. Sketches are also of independent interest. In particular, Song et al. (2023) described a way of designing the distribution of sketch matrices, while Lee et al. (2019); Qin et al. (2023) used sketches in solving empirical risk minimization problems.

Leveraging the layer-wise structure of neural networks has been widely studied for optimizing the training loss function. For example, (Zheng et al., 2019) propose SGD with different scalar stepsizes for each layer, (Yu et al., 2017; Ginsburg et al., 2019) propose layer-wise normalization for Stochastic Normalized Gradient Descent, and (Dutta et al., 2020; Wang et al., 2022) propose layer-wise compression in the distributed setting.

DCGD, proposed by Khirirat et al. (2018), has since been improved in various ways, such as in (Horvath et al., 2019; Li et al., 2020). There is also a large body of literature on other federated learning algorithms with unbiased compressors (Alistarh et al., 2017; Mishchenko et al., 2019; Gorbunov et al., 2021; Mishchenko et al., 2022; Maranjyan et al., 2022; Horváth et al., 2023).

1.2 CONTRIBUTIONS

Our paper contributes in the following ways:

- We propose two novel matrix stepsize sketch CGD algorithms in Section 2, which, to the best of our knowledge, are the first attempts to analyze a fixed matrix stepsize for non-convex optimization. We present a unified theorem in Section 3 that guarantees stationarity for minimizing matrix-smooth non-convex functions. The results show that taking our algorithms improve on their scalar alternatives. The complexities are summarized in Table 1 for some particular cases.
- We design our algorithms’ sketches and stepsize to take advantage of the layer-wise structure of neural networks, assuming that the smoothness matrix is block-diagonal. In Section 4, we prove that our algorithms achieve better convergence than classical methods.
- Assuming that the server-to-client communication is less expensive Konečný et al. (2016); Kairouz et al. (2021), we propose distributed versions of our algorithms in Section 5, following the standard FL scheme, and prove weighted stationarity guarantees. Our theorem recovers the result for DCGD in the scalar case and improves it in general.
- We validate our theoretical results with experiments. The plots and framework are provided in the Appendix.

1.3 PRELIMINARIES

The usual Euclidean norm on \mathbb{R}^d is defined as $\|\cdot\|$. We use bold capital letters to denote matrices. By \mathbf{I}_d we denote the $d \times d$ identity matrix, and by \mathbf{O}_d we denote the $d \times d$ zero matrix. Let \mathbb{S}_{++}^d (resp. \mathbb{S}_+^d) be the set of $d \times d$ symmetric positive definite (resp. semi-definite) matrices. Given $\mathbf{Q} \in \mathbb{S}_{++}^d$ and $x \in \mathbb{R}^d$, we write $\|x\|_{\mathbf{Q}} := \sqrt{\langle \mathbf{Q}x, x \rangle}$, where $\langle \cdot, \cdot \rangle$ is the standard Euclidean inner product on \mathbb{R}^d . For a matrix $\mathbf{A} \in \mathbb{S}_{++}^d$, we define by $\lambda_{\max}(\mathbf{A})$ (resp. $\lambda_{\min}(\mathbf{A})$) the largest (resp. smallest) eigenvalue of the matrix \mathbf{A} . Let $\mathbf{A}_i \in \mathbb{R}^{d_i \times d_i}$ and $d = d_1 + \dots + d_\ell$. Then the matrix $\mathbf{A} = \text{Diag}(\mathbf{A}_1, \dots, \mathbf{A}_\ell)$ is defined as a block diagonal $d \times d$ matrix where the i -th block is equal to \mathbf{A}_i . We will use $\text{diag}(\mathbf{A}) \in \mathbb{R}^{d \times d}$ to denote the diagonal of any matrix $\mathbf{A} \in \mathbb{R}^{d \times d}$. Given a function $f : \mathbb{R}^d \rightarrow \mathbb{R}$, its gradient and its Hessian at point $x \in \mathbb{R}^d$ are respectively denoted as $\nabla f(x)$ and $\nabla^2 f(x)$. A random vector $x \in \mathbb{R}^d$ is an ε -stationary point if $\mathbb{E} \left[\|\nabla f(x)\|^2 \right] \leq \varepsilon^2$, where the expectation is over the randomness of the algorithm.

2 THE ALGORITHMS

Below we define our two main algorithms:

$$x^{k+1} = x^k - \mathbf{D}\mathbf{S}^k\nabla f(x^k), \quad (\text{det-CGD1})$$

and

$$x^{k+1} = x^k - \mathbf{T}^k\mathbf{D}\nabla f(x^k). \quad (\text{det-CGD2})$$

Here, $\mathbf{D} \in \mathbb{S}_{++}^d$ is the fixed stepsize matrix. The sequences of random matrices \mathbf{S}^k and \mathbf{T}^k satisfy the following assumption.

Assumption 3. *We will assume that the random sketches that appear in our algorithms are i.i.d., unbiased, symmetric and positive semi-definite for each algorithm. That is*

$$\begin{aligned} \mathbf{S}^k, \mathbf{T}^k &\in \mathbb{S}_+^d, \quad \mathbf{S}^k \stackrel{iid}{\sim} \mathcal{S} \quad \text{and} \quad \mathbf{T}^k \stackrel{iid}{\sim} \mathcal{T} \\ \mathbb{E}[\mathbf{S}^k] &= \mathbb{E}[\mathbf{T}^k] = \mathbf{I}_d, \quad \text{for every } k \in \mathbb{N}, \end{aligned}$$

where \mathcal{S} and \mathcal{T} are probability distributions over \mathbb{S}_+^d .

A simple instance of [det-CGD1](#) and [det-CGD2](#) is the vanilla GD. Indeed, if $\mathbf{S}^k = \mathbf{T}^k = \mathbf{I}_d$ and $\mathbf{D} = \gamma\mathbf{I}_d$, then $x^{k+1} = x^k - \gamma\nabla f(x^k)$. In general, one may view these algorithms as Newton-type methods. In particular, our setting includes the Newton Star (NS) algorithm by [Islamov et al. \(2021\)](#):

$$x^{k+1} = x^k - (\nabla^2 f(x^{\text{inf}}))^{-1} \nabla f(x^k). \quad (\text{NS})$$

The authors prove that in the convex case it converges to the unique solution x^{inf} locally quadratically, provided certain assumptions are met. However, it is not a practical method as it requires knowledge of the Hessian at the optimal point. This method, nevertheless, hints that constant matrix stepsize can yield fast convergence guarantees. Our results allow us to choose the \mathbf{D} depending on the smoothness matrix \mathbf{L} . The latter can be seen as a uniform upper bound on the Hessian.

The difference between [det-CGD1](#) and [det-CGD2](#) is the update rule. In particular, the order of the sketch and the stepsize is interchanged. When the sketch \mathbf{S} and the stepsize \mathbf{D} are commutative w.r.t. matrix product, the algorithms become equivalent. In general, a simple calculation shows that if we take

$$\mathbf{T}^k = \mathbf{D}\mathbf{S}^k\mathbf{D}^{-1}, \quad (5)$$

then [det-CGD1](#) and [det-CGD2](#) are the same. Defining \mathbf{T}^k according to (5), we recover the unbiasedness condition:

$$\mathbb{E}[\mathbf{T}^k] = \mathbf{D}\mathbb{E}[\mathbf{S}^k]\mathbf{D}^{-1} = \mathbf{I}_d. \quad (6)$$

However, in general $\mathbf{D}\mathbb{E}[\mathbf{S}^k]\mathbf{D}^{-1}$ is not necessarily symmetric, which contradicts to Assumption 3. Thus, [det-CGD1](#) and [det-CGD2](#) are not equivalent for our purposes.

3 MAIN RESULTS

Before we state the main result, we present a stepsize condition for [det-CGD1](#) and [det-CGD2](#), respectively:

$$\mathbb{E}[\mathbf{S}^k\mathbf{D}\mathbf{L}\mathbf{D}\mathbf{S}^k] \preceq \mathbf{D}, \quad (7)$$

and

$$\mathbb{E}[\mathbf{D}\mathbf{T}^k\mathbf{L}\mathbf{T}^k\mathbf{D}] \preceq \mathbf{D}. \quad (8)$$

In the case of vanilla GD (7) and (8) become $\gamma < L^{-1}$, which is the standard condition for convergence. Below is the main convergence theorem for both algorithms in the single-node regime.

Theorem 1. *Suppose that Assumptions 1-3 are satisfied. Then, for each $k \geq 0$*

$$\frac{1}{K} \sum_{k=0}^{K-1} \mathbb{E} \left[\|\nabla f(x^k)\|_{\mathbf{D}}^2 \right] \leq \frac{2(f(x^0) - f^{\text{inf}})}{K}, \quad (9)$$

if one of the below conditions is true:

- i) The vectors x^k are the iterates of *det-CGD1* and \mathbf{D} satisfies (7);
- ii) The vectors x^k are the iterates of *det-CGD2* and \mathbf{D} satisfies (8).

It is important to note that Theorem 1 yields the same convergence rate for any $\mathbf{D} \in \mathbb{S}_{++}^d$, despite the fact that the matrix norms on the left-hand side cannot be compared for different weight matrices. To ensure comparability of the right-hand side of (9), it is necessary to normalize the weight matrix \mathbf{D} that is used to measure the gradient norm. We propose using determinant normalization, which involves dividing both sides of (9) by $\det(\mathbf{D})^{1/d}$, yielding the following:

$$\frac{1}{K} \sum_{k=0}^{K-1} \mathbb{E} \left[\left\| \nabla f(x^k) \right\|_{\frac{\mathbf{D}}{\det(\mathbf{D})^{1/d}}}^2 \right] \leq \frac{2(f(x^0) - f^{\text{inf}})}{\det(\mathbf{D})^{1/d} K}. \quad (10)$$

This normalization is meaningful because adjusting the weight matrix to $\frac{\mathbf{D}}{\det(\mathbf{D})^{1/d}}$ allows its determinant to be 1, making the norm on the left-hand side comparable to the standard Euclidean norm. It is important to note that the volume of the normalized ellipsoid $\{x \in \mathbb{R}^d : \|x\|_{\mathbf{D}/\det(\mathbf{D})^{1/d}}^2 \leq 1\}$ does not depend on the choice of $\mathbf{D} \in \mathbb{S}_{++}^d$. Therefore, the results of (9) are comparable across different \mathbf{D} in the sense that the right-hand side of (9) measures the volume of the ellipsoid containing the gradient.

3.1 OPTIMAL MATRIX STEPSIZE

In this section, we describe how to choose the optimal stepsize that minimizes the iteration complexity. The problem is easier for *det-CGD2*. We notice that (8) can be explicitly solved. Specifically, it is equivalent to

$$\mathbf{D} \preceq (\mathbb{E} [\mathbf{T}^k \mathbf{L} \mathbf{T}^k])^{-1}. \quad (11)$$

We want to emphasize that the RHS matrix is invertible despite the sketches not being so. Indeed. The map $h : \mathbf{T} \rightarrow \mathbf{T} \mathbf{L} \mathbf{T}$ is convex on \mathbb{S}_+^d . Therefore, Jensen's inequality implies

$$\mathbb{E} [\mathbf{T}^k \mathbf{L} \mathbf{T}^k] \succeq \mathbb{E} [\mathbf{T}^k] \mathbf{L} \mathbb{E} [\mathbf{T}^k] = \mathbf{L} \succ \mathbf{O}_d.$$

This explicit condition on \mathbf{D} can assist in determining the optimal stepsize. Since both \mathbf{D} and $(\mathbb{E} [\mathbf{T}^k \mathbf{L} \mathbf{T}^k])^{-1}$ are positive definite, then the right-hand side of (10) is minimized exactly when

$$\mathbf{D} = (\mathbb{E} [\mathbf{T}^k \mathbf{L} \mathbf{T}^k])^{-1}. \quad (12)$$

Note that the explicit solution of \mathbf{D} needs to be calculated only once, at the beginning of the algorithm. It is then fixed for all iterations. The situation is different for *det-CGD1*. According to (10), the optimal \mathbf{D} is defined as the solution of the following constrained optimization problem:

$$\begin{aligned} & \text{minimize} && \log \det(\mathbf{D}^{-1}) \\ & \text{subject to} && \mathbb{E} [\mathbf{S}^k \mathbf{D} \mathbf{L} \mathbf{D} \mathbf{S}^k] \preceq \mathbf{D} \\ & && \mathbf{D} \in \mathbb{S}_{++}^d. \end{aligned} \quad (13)$$

Proposition 1. *The optimization problem (13) with respect to stepsize matrix $\mathbf{D} \in \mathbb{S}_{++}^d$, is a convex optimization problem with a convex constraint.*

The proof of this proposition can be found in the Appendix. It is based on the reformulation of the constraint to its equivalent quadratic form inequality. Using the trace trick, we can prove that for every vector chosen in the quadratic form, it is convex. Since the intersection of convex sets is convex, we conclude the proof.

One could consider using the CVXPY (Diamond & Boyd, 2016) package to solve (13), provided that it is first transformed into a Disciplined Convex Programming (DCP) form (Grant et al., 2006). Nevertheless, (7) is not recognized as a DCP constraint in the general case. To make CVXPY applicable, additional steps tailored to the problem at hand must be taken.

Table 1: Summary of communication complexities of **det-CGD1** and **det-CGD2** with different sketches and stepsize matrices. The \mathbf{D}_i here for **det-CGD1** is \mathbf{W}_i with the optimal scaling determined using Theorem 2, for **det-CGD2** it is the optimal stepsize matrix defined in (12). The constant $2(f(x^0) - f^{\text{inf}})/\varepsilon^2$ is hidden, ℓ is the number of layers, k_i is the mini-batch size for the i -th layer if we use the rand- k sketch. The notation $\tilde{\mathbf{L}}_{i,k}$ is defined as $\frac{d-k}{d-1} \text{diag}(\mathbf{L}_i) + \frac{k-1}{d-1} \mathbf{L}_i$.

No.	The method	$(\mathbf{S}_i^k, \mathbf{D}_i)$	$l \geq 1, d_i, k_i, \sum_{i=1}^{\ell} k_i = k$, layer structure	$l = 1, k_i = k$, general structure
1.	det-CGD1	$(\mathbf{I}_d, \gamma \mathbf{L}_i^{-1})$	$d \cdot \det(\mathbf{L})^{1/d}$	$d \cdot \det(\mathbf{L})^{1/d}$
2.	det-CGD1	$(\mathbf{I}_d, \gamma \text{diag}^{-1}(\mathbf{L}_i))$	$d \cdot \det(\text{diag}(\mathbf{L}))^{1/d}$	$d \cdot \det(\text{diag}(\mathbf{L}))^{1/d}$
3.	det-CGD1	$(\mathbf{I}_d, \gamma \mathbf{I}_{d_i})$	$d \cdot \left(\prod_{i=1}^{\ell} \lambda_{\max}^{d_i}(\mathbf{L}_i)\right)^{1/d}$	$d \cdot \lambda_{\max}(\mathbf{L})$
4.	det-CGD1	$(\text{rand-1}, \gamma \mathbf{I}_{d_i})$	$\ell \cdot \left(\prod_{i=1}^{\ell} d_i^{d_i} (\max_j (\mathbf{L}_i)_{jj})^{d_i}\right)^{1/d}$	$d \cdot \max_j (\mathbf{L}_{jj})$
5.	det-CGD1	$(\text{rand-1}, \gamma \mathbf{L}_i^{-1})$	$\ell \cdot \left(\frac{\prod_{i=1}^{\ell} d_i^{d_i} \lambda_{\max}^{d_i}(\mathbf{L}_i^{\frac{1}{2}} \text{diag}(\mathbf{L}_i^{-1}) \mathbf{L}_i^{\frac{1}{2}})}{\prod_{i=1}^{\ell} \det(\mathbf{L}_i^{-1})}\right)^{1/d}$	$\frac{d \lambda_{\max}(\mathbf{L}^{\frac{1}{2}} \text{diag}(\mathbf{L}^{-1}) \mathbf{L}^{\frac{1}{2}})}{\det(\mathbf{L}^{-1})^{1/d}}$
6.	det-CGD1	$(\text{rand-1}, \gamma \mathbf{L}_i^{-1/2})$	$\ell \cdot \left(\frac{\prod_{i=1}^{\ell} d_i^{d_i} \lambda_{\max}^{d_i}(\mathbf{L}_i^{1/2})}{\prod_{i=1}^{\ell} \det(\mathbf{L}_i^{-1/2})}\right)^{1/d}$	$d \cdot \lambda_{\max}^{1/2}(\mathbf{L}) \det(\mathbf{L})^{1/(2d)}$
7.	det-CGD1	$(\text{rand-1}, \gamma \text{diag}^{-1}(\mathbf{L}_i))$	$\ell \cdot \left(\frac{\prod_{i=1}^{\ell} d_i^{d_i}}{\prod_{j=1}^d (\mathbf{L}_{jj}^{-1})}\right)^{1/d}$	$d \cdot \det(\text{diag}(\mathbf{L}))^{1/d}$
8.	det-CGD1	$(\text{rand-}k_i, \gamma \text{diag}^{-1}(\mathbf{L}_i))$	$k \cdot \left(\prod_{i=1}^{\ell} \left(\frac{d_i}{k_i}\right)^{d_i} \det(\text{diag}(\mathbf{L}))\right)^{1/d}$	$d \cdot \det(\text{diag}(\mathbf{L}))^{1/d}$
9.	det-CGD2	$(\mathbf{I}_d, \mathbf{L}_i^{-1})$	$d \cdot \det(\mathbf{L})^{1/d}$	$d \cdot \det(\mathbf{L})^{1/d}$
10.	det-CGD2	$(\text{rand-1}, \frac{\text{diag}^{-1}(\mathbf{L}_i)}{d_i})$	$\ell \cdot \left(\prod_{i=1}^{\ell} d_i^{d_i}\right)^{1/d} \det(\text{diag} \mathbf{L})^{1/d}$	$d \cdot \det(\text{diag}(\mathbf{L}))^{1/d}$
11.	det-CGD2	$(\text{rand-}k, \frac{k_i}{d_i} \tilde{\mathbf{L}}_{i,k_i}^{-1})$	$k \cdot \left(\prod_{i=1}^{\ell} \left(\frac{d_i}{k_i}\right)^{d_i}\right) \left(\prod_{i=1}^{\ell} \det(\tilde{\mathbf{L}}_{i,k_i})\right)^{1/d}$	$d \cdot \det(\tilde{\mathbf{L}}_{1,k})$
12.	det-CGD2	$(\text{Bern-}q_i, q_i \mathbf{L}_i^{-1})$	$\left(\sum_{i=1}^{\ell} q_i d_i\right) \cdot \prod_{i=1}^{\ell} \left(\frac{d_i}{q_i}\right)^{d_i} \det(\mathbf{L})^{1/d}$	$d \cdot \det(\mathbf{L})^{1/d}$
13.	GD	$(\mathbf{I}_d, \lambda_{\max}^{-1}(\mathbf{L}) \mathbf{I}_d)$	N/A	$d \cdot \lambda_{\max}(\mathbf{L})$

4 LEVERAGING THE LAYER-WISE STRUCTURE

In this section we focus on the block-diagonal case of \mathbf{L} for both **det-CGD1** and **det-CGD2**. In particular, we propose hyper-parameters of **det-CGD1** designed specifically for training NNs. Let us assume that $\mathbf{L} = \text{Diag}(\mathbf{L}_1, \dots, \mathbf{L}_{\ell})$, where $\mathbf{L}_i \in \mathbb{S}_{++}^{d_i}$. This setting is a generalization of the classical smoothness condition, as in the latter case $\mathbf{L}_i = \mathbf{L} \mathbf{I}_{d_i}$ for all $i = 1, \dots, \ell$. Respectively, we choose both the sketches and the stepsize to be block diagonal: $\mathbf{D} = \text{Diag}(\mathbf{D}_1, \dots, \mathbf{D}_{\ell})$ and $\mathbf{S}^k = \text{Diag}(\mathbf{S}_1^k, \dots, \mathbf{S}_{\ell}^k)$, where $\mathbf{D}_i, \mathbf{S}_i^k \in \mathbb{S}_{++}^{d_i}$.

Let us notice that the left hand side of the inequality constraint in (13) has quadratic dependence on \mathbf{D} , while the right hand side is linear. Thus, for every matrix $\mathbf{W} \in \mathbb{S}_{++}^d$, there exists $\gamma > 0$ such that

$$\gamma^2 \lambda_{\max}(\mathbb{E}[\mathbf{S}^k \mathbf{W} \mathbf{L} \mathbf{W} \mathbf{S}^k]) \leq \gamma \lambda_{\min}(\mathbf{W}).$$

Therefore, for $\gamma \mathbf{W}$ we deduce

$$\mathbb{E}[\mathbf{S}^k (\gamma \mathbf{W}) \mathbf{L} (\gamma \mathbf{W}) \mathbf{S}^k] \preceq \gamma^2 \lambda_{\max}(\mathbb{E}[\mathbf{S}^k \mathbf{W} \mathbf{L} \mathbf{W} \mathbf{S}^k]) \mathbf{I}_d \preceq \gamma \lambda_{\min}(\mathbf{W}) \mathbf{I}_d \preceq \gamma \mathbf{W}. \quad (14)$$

The following theorem is based on this simple fact applied to the corresponding blocks of the matrices $\mathbf{D}, \mathbf{L}, \mathbf{S}^k$ for **det-CGD1**.

Theorem 2. Let $f : \mathbb{R}^d \rightarrow \mathbb{R}$ satisfy Assumptions 1 and 2, with \mathbf{L} admitting the layer-separable structure $\mathbf{L} = \text{Diag}(\mathbf{L}_1, \dots, \mathbf{L}_{\ell})$, where $\mathbf{L}_1, \dots, \mathbf{L}_{\ell} \in \mathbb{S}_{++}^{d_i}$. Choose random matrices $\mathbf{S}_1^k, \dots, \mathbf{S}_{\ell}^k \in \mathbb{S}_{++}^{d_i}$ to satisfy Assumption 3 for all $i \in [\ell]$, and let $\mathbf{S}^k := \text{Diag}(\mathbf{S}_1^k, \dots, \mathbf{S}_{\ell}^k)$. Furthermore, choose matrices $\mathbf{W}_1, \dots, \mathbf{W}_{\ell} \in \mathbb{S}_{++}^{d_i}$ and scalars $\gamma_1, \dots, \gamma_{\ell} > 0$ such that

$$\gamma_i \leq \lambda_{\max}^{-1} \left(\mathbb{E} \left[\mathbf{W}_i^{-1/2} \mathbf{S}_i^k \mathbf{W}_i \mathbf{L}_i \mathbf{W}_i \mathbf{S}_i^k \mathbf{W}_i^{-1/2} \right] \right) \quad \forall i \in [\ell]. \quad (15)$$

Letting $\mathbf{W} := \text{Diag}(\mathbf{W}_1, \dots, \mathbf{W}_{\ell})$, $\Gamma := \text{Diag}(\gamma_1 \mathbf{I}_{d_1}, \dots, \gamma_{\ell} \mathbf{I}_{d_{\ell}})$ and $\mathbf{D} := \Gamma \mathbf{W}$, we get

$$\frac{1}{K} \sum_{k=0}^{K-1} \mathbb{E} \left[\left\| \nabla f(x^k) \right\|_{\frac{\Gamma \mathbf{W}}{\det(\Gamma \mathbf{W})^{1/d}}}^2 \right] \leq \frac{2(f(x^0) - f^{\text{inf}})}{\det(\Gamma \mathbf{W})^{1/d} K}. \quad (16)$$

In particular, if the scalars $\{\gamma_i\}$ are chosen to be equal to their maximum allowed values from (15), then the convergence factor of (16) is equal to

$$\det(\Gamma\mathbf{W})^{-\frac{1}{d}} = \left[\prod_{i=1}^{\ell} \lambda_{\max}^{d_i} \left(\mathbb{E} \left[\mathbf{W}_i^{-\frac{1}{2}} \mathbf{S}_i^k \mathbf{W}_i \mathbf{L}_i \mathbf{W}_i \mathbf{S}_i^k \mathbf{W}_i^{-\frac{1}{2}} \right] \right) \right]^{\frac{1}{d}} \det(\mathbf{W}^{-1})^{\frac{1}{d}}.$$

Table 1 contains the (expected) communication complexities of **det-CGD1**, **det-CGD2** and GD for several choices of \mathbf{W} , \mathbf{D} and \mathbf{S}^k . Here are a few comments about the table. We deduce that taking a matrix stepsize without compression (row 1) we improve GD (row 13). A careful analysis reveals that the result in row 5 is always worse than row 7 in terms of both communication and iteration complexity. However, the results in row 6 and row 7 are not comparable in general, meaning that neither of them is universally better. More discussion on this table can be found in the Appendix.

Compression for free. Now, let us focus on row 12, which corresponds to a sampling scheme where the i -th layer is independently selected with probability q_i . Mathematically, it goes as follows:

$$\mathbf{T}_i^k = \frac{\eta_i}{q_i} \mathbf{I}_{d_i}, \quad \text{where } \eta_i \sim \text{Bernoulli}(q_i). \quad (17)$$

Jensen’s inequality implies that

$$\left(\sum_{i=1}^{\ell} q_i d_i \right) \cdot \prod_{i=1}^{\ell} \left(\frac{1}{q_i} \right)^{\frac{d_i}{d}} \geq d. \quad (18)$$

The equality is attained when $q_i = q$ for all $i \in [\ell]$. The expected bits transferred per iteration of this algorithm is then equal to $k_{\text{exp}} = qd$ and the communication complexity equals $d \det(\mathbf{L})^{1/d}$. Comparing with the results for **det-CGD2** with $\text{rand-}k_{\text{exp}}$ on row 11 and using the fact that $\det(\mathbf{L}) \leq \det(\text{diag}(\mathbf{L}))$, we deduce that the Bernoulli scheme is better than the uniform sampling scheme. Notice also, the communication complexity matches the one for the uncompressed **det-CGD2** displayed on row 9. This, in particular means that using the Bern- q sketches we can compress the gradients for free. The latter means that we reduce the number of bits broadcasted at each iteration without losing in the total communication complexity. In particular, when all the layers have the same width d_i , the number of broadcasted bits for each iteration is reduced by a factor of q .

5 DISTRIBUTED SETTING

In this section we describe the distributed versions of our algorithms and present convergence guarantees for them. Let us consider an objective function that is sum decomposable:

$$f(x) := \frac{1}{n} \sum_{i=1}^n f_i(x),$$

where each $f_i : \mathbb{R}^d \rightarrow \mathbb{R}$ is a differentiable function. We assume that f satisfies Assumption 1 and the component functions satisfy the below condition.

Assumption 4. *Each component function f_i is L_i -smooth and is bounded from below: $f_i(x) \geq f_i^{\text{inf}}$ for all $x \in \mathbb{R}^d$.*

This assumption also implies that f is of matrix smoothness with $\bar{\mathbf{L}} \in \mathbb{S}_{++}^d$, where $\bar{\mathbf{L}} = \frac{1}{n} \sum_{i=1}^n \mathbf{L}_i$. Following the standard FL framework (Konečný et al., 2016; McMahan et al., 2017; Khirirat et al., 2018), we assume that the i -th component function f_i is stored on the i -th client. At each iteration, the clients in parallel compute and compress the local gradient ∇f_i and communicate it to the central server. The server, then aggregates the compressed gradients, computes the next iterate, and in parallel broadcasts it to the clients. See the below pseudo-codes for the details.

Theorem 3. *Let $f_i : \mathbb{R}^d \rightarrow \mathbb{R}$ satisfy Assumption 4 and let f satisfy Assumption 1 and Assumption 2 with smoothness matrix \mathbf{L} . If the stepsize satisfies*

$$DLD \preceq D, \quad (19)$$

then the following convergence bound is true for the iterates of Algorithm 1:

$$\min_{0 \leq k \leq K-1} \mathbb{E} \left[\left\| \nabla f(x^k) \right\|_{\frac{\mathbf{D}}{\det(\mathbf{D})^{1/d}}}^2 \right] \leq \frac{2(1 + \frac{\lambda_{\mathbf{D}}}{n})^K (f(x^0) - f^{\text{inf}})}{\det(\mathbf{D})^{1/d} K} + \frac{2\lambda_{\mathbf{D}}\Delta^{\text{inf}}}{\det(\mathbf{D})^{1/d} n}, \quad (20)$$

where $\Delta^{\text{inf}} := f^{\text{inf}} - \frac{1}{n} \sum_{i=1}^n f_i^{\text{inf}}$ and

$$\lambda_{\mathbf{D}} := \max_i \left\{ \lambda_{\max} \left(\mathbb{E} \left[\mathbf{L}_i^{\frac{1}{2}} (\mathbf{S}_i^k - \mathbf{I}_d) \mathbf{D} \mathbf{L} \mathbf{D} (\mathbf{S}_i^k - \mathbf{I}_d) \mathbf{L}_i^{\frac{1}{2}} \right] \right) \right\}.$$

Algorithm 1 Distributed **det-CGD1**

1: **Input:** Starting point x^0 , stepsize matrix \mathbf{D} , number of iterations K
2: **for** $k = 0, 1, 2, \dots, K - 1$ **do**
3: The devices in parallel:
4: sample $\mathbf{S}_i^k \sim \mathcal{S}$;
5: compute $\mathbf{S}_i^k \nabla f_i(x^k)$;
6: broadcast $\mathbf{S}_i^k \nabla f_i(x^k)$.
7: **The server:**
8: combines $g^k = \frac{\mathbf{D}}{n} \sum_{i=1}^n \mathbf{S}_i^k \nabla f_i(x^k)$;
9: computes $x^{k+1} = x^k - g^k$;
10: broadcasts x^{k+1} .
11: **end for**
12: **Return:** x^K

Algorithm 2 Distributed **det-CGD2**

1: **Input:** Starting point x^0 , stepsize matrix \mathbf{D} , number of iterations K
2: **for** $k = 0, 1, 2, \dots, K - 1$ **do**
3: The devices in parallel:
4: sample $\mathbf{T}_i^k \sim \mathcal{T}$;
5: compute $\mathbf{T}_i^k \mathbf{D} \nabla f_i(x^k)$;
6: broadcast $\mathbf{T}_i^k \mathbf{D} \nabla f_i(x^k)$.
7: **The server:**
8: combines $g^k = \frac{1}{n} \sum_i \mathbf{T}_i^k \mathbf{D} \nabla f_i(x^k)$;
9: computes $x^{k+1} = x^k - g^k$;
10: broadcasts x^{k+1} .
11: **end for**
12: **Return:** x^K

The same result is true for Algorithm 2 with a different constant $\lambda_{\mathbf{D}}$. The proof of Theorem 3 and its analogue for Algorithm 2 are presented in the Appendix. The analysis is largely inspired by (Khaled & Richtárik, 2020, Theorem 1). Now, let us examine the right-hand side of (20). We start by observing that the first term has exponential dependence in K . However, the term inside the brackets, $1 + \lambda_{\mathbf{D}}/n$, depends on the stepsize \mathbf{D} . Furthermore, it has a second-order dependence on \mathbf{D} , implying that $\lambda_{\alpha \mathbf{D}} = \alpha^2 \lambda_{\mathbf{D}}$, as opposed to $\det(\alpha \mathbf{D})^{1/d}$, which is linear in α . Therefore, we can choose a small enough coefficient α to ensure that $\lambda_{\mathbf{D}}$ is of order n/K . This means that for a fixed number of iterations K , we choose the matrix stepsize to be "small enough" to guarantee that the numerator of the first term is bounded. The following corollary summarizes these arguments, and its proof can be found in the Appendix.

Corollary 1. We reach an ε -stationarity, that is the right-hand side of (20) is upper bounded by ε^2 , if the following conditions are satisfied:

$$\mathbf{D} \mathbf{L} \mathbf{D} \preceq \mathbf{D}, \quad \lambda_{\mathbf{D}} \leq \min \left\{ \frac{n}{K}, \frac{n\varepsilon^2}{4\Delta^{\text{inf}} \det(\mathbf{D})^{1/d}} \right\}, \quad K \geq \frac{12(f(x^0) - f^{\text{inf}})}{\det(\mathbf{D})^{1/d} \varepsilon^2}. \quad (21)$$

Proposition 3 in the Appendix proves that these conditions with respect to \mathbf{D} are convex. In order to minimize the iteration complexity for getting ε^2 error, one needs to solve the following optimization problem

$$\text{minimize} \quad \log \det(\mathbf{D}^{-1})$$

$$\text{subject to} \quad \mathbf{D} \text{ satisfies (21).}$$

Choosing the optimal stepsize for Algorithm 1 is analogous to solving (13). One can formulate the distributed counterpart of Theorem 2 and attempt to solve it for different sketches. Furthermore, this leads to a convex matrix minimization problem involving \mathbf{D} . We provide a formal proof of this property in the Appendix. Similar to the single-node case, computational methods can be employed using the CVXPY package. However, some additional effort is required to transform (21) into the disciplined convex programming (DCP) format.

The second term in (20) corresponds to the convergence neighborhood of the algorithm. It does not depend on the number of iteration, thus it remains unchanged, after we choose the stepsize. Nevertheless, it depends on the number of clients n . In general, the term Δ^{inf}/n can be unbounded, when $n \rightarrow +\infty$. However, per Corollary 1, we require $\lambda_{\mathbf{D}}$ to be upper-bounded by n/K . Thus,

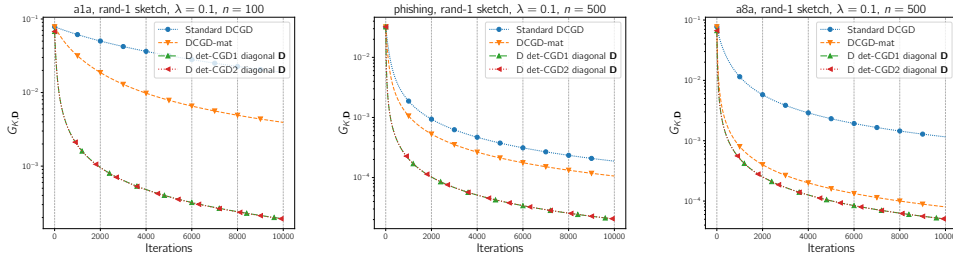


Figure 1: Comparison of standard DCGD, DCGD with matrix smoothness, D-det-CGD1 and D-det-CGD2 with optimal diagonal stepsizes under rand-1 sketch. The stepsize for standard DCGD is determined using (Khaled & Richtárik, 2020, Proposition 4), the stepsize for DCGD with matrix smoothness along with D_1 , D_2 is determined using Corollary 1, the error level is set to be $\varepsilon^2 = 0.0001$. Here $G_{K,D} := \frac{1}{K} \left(\sum_{k=0}^{K-1} \|\nabla f(x^k)\|_{D/\det(D)^{1/d}}^2 \right)$.

the neighborhood term will indeed converge to zero when $K \rightarrow +\infty$, if we choose the stepsize accordingly.

We compare our results with the existing results for DCGD. In particular we use the technique from Khaled & Richtárik (2020) for the scalar smooth DCGD with scalar stepsizes with the results from (Khaled & Richtárik, 2020, Corollary 1). See the Appendix for the details on the analysis of Khaled & Richtárik (2020). Finally, we back up our theoretical findings with experiments. See Figure 1 for a simple experiment confirming that Algorithms 1 and 2 have better iteration and communication complexity compared to scalar stepped DCGD. The graphs of the two proposed algorithms coincide, as the diagonal stepsize and the diagonal sketch commute, resulting in the same method. For more details on the experiments we refer the reader to the corresponding section in the Appendix.

6 CONCLUSION

In this paper, we enhance compressed gradient descent method with matrix-valued stepsize for general non-convex objectives. Convergence guarantees are provided for the algorithms both in the single node case and the distributed setting. By considering the layer-wise structure of models such as neural networks, we are able to design compression mechanisms that achieve compression for free. This is the first time matrix stepsize is used and analyzed together with compression in the non-convex case. Our theoretical findings are supported with abundant numerical experiments.

6.1 LIMITATIONS

It is worth noting that every point in \mathbb{R}^d can be enclosed within some volume 1 ellipsoid. Therefore, having the average D -norm of the gradient bounded by a small number does not guarantee that the average Euclidean norm is small. However, for a fixed D , the standard Euclidean norm is equivalent to the weighted D -norm. This is due to

$$\frac{\lambda_{\min}(D)}{\det(D)^{1/d}} \cdot \|\nabla f(x)\|^2 \leq \|\nabla f(x)\|_{D/(\det(D))^{1/d}}^2 \leq \frac{\lambda_{\max}(D)}{\det(D)^{1/d}} \cdot \|\nabla f(x)\|^2. \quad (22)$$

This relation is further validated by our experiments described in the Appendix.

REFERENCES

- Mehiddin Al-Baali and H Khalfan. An overview of some practical quasi-newton methods for unconstrained optimization. *Sultan Qaboos University Journal for Science [SQUJS]*, 12(2): 199–209, 2007.
- Mehiddin Al-Baali, Emilio Spedicato, and Francesca Maggioni. Broyden’s quasi-Newton methods for a nonlinear system of equations and unconstrained optimization: a review and open problems. *Optimization Methods and Software*, 29(5):937–954, 2014.
- Dan Alistarh, Demjan Grubic, Jerry Li, Ryota Tomioka, and Milan Vojnovic. QSGD: Communication-efficient SGD via gradient quantization and encoding. *Advances in neural information processing systems*, 30, 2017.
- Charles G Broyden. A class of methods for solving nonlinear simultaneous equations. *Mathematics of computation*, 19(92):577–593, 1965.
- Sébastien Bubeck et al. Convex optimization: Algorithms and complexity. *Foundations and Trends® in Machine Learning*, 8(3-4):231–357, 2015.
- Chih-Chung Chang and Chih-Jen Lin. Libsvm: a library for support vector machines. *ACM transactions on intelligent systems and technology (TIST)*, 2(3):1–27, 2011.
- Marina Danilova, Pavel Dvurechensky, Alexander Gasnikov, Eduard Gorbunov, Sergey Guminov, Dmitry Kamzolov, and Innokentiy Shibaev. Recent theoretical advances in non-convex optimization. In *High-Dimensional Optimization and Probability: With a View Towards Data Science*, pp. 79–163. Springer, 2022.
- John E Dennis, Jr and Jorge J Moré. Quasi-Newton methods, motivation and theory. *SIAM review*, 19(1):46–89, 1977.
- Steven Diamond and Stephen Boyd. CVXPY: A Python-embedded modeling language for convex optimization. *The Journal of Machine Learning Research*, 17(1):2909–2913, 2016.
- Aritra Dutta, El Houcine Bergou, Ahmed M Abdelmoniem, Chen-Yu Ho, Atal Narayan Sahu, Marco Canini, and Panos Kalnis. On the discrepancy between the theoretical analysis and practical implementations of compressed communication for distributed deep learning. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 34, pp. 3817–3824, 2020.
- Darina Dvinskikh, Aleksandr Ogaltsov, Alexander Gasnikov, Pavel Dvurechensky, Alexander Tyurin, and Vladimir Spokoiny. Adaptive gradient descent for convex and non-convex stochastic optimization. *arXiv preprint arXiv:1911.08380*, 2019.
- Boris Ginsburg, Patrice Castonguay, Oleksii Hrinchuk, Oleksii Kuchaiev, Vitaly Lavruchin, Ryan Leary, Jason Li, Huyen Nguyen, Yang Zhang, and Jonathan M Cohen. Stochastic gradient methods with layer-wise adaptive moments for training of deep networks. *arXiv preprint arXiv:1905.11286*, 2019.
- Eduard Gorbunov, Konstantin P Burlachenko, Zhize Li, and Peter Richtárik. Marina: Faster non-convex distributed learning with compression. In *International Conference on Machine Learning*, pp. 3788–3798. PMLR, 2021.
- Robert M Gower and Peter Richtárik. Randomized iterative methods for linear systems. *SIAM Journal on Matrix Analysis and Applications*, 36(4):1660–1690, 2015.
- Robert Mansel Gower, Nicolas Loizou, Xun Qian, Alibek Sailanbayev, Egor Shulgin, and Peter Richtárik. Sgd: General analysis and improved rates. In *International Conference on Machine Learning*, pp. 5200–5209. PMLR, 2019.
- William B Gragg and Richard A Tapia. Optimal error bounds for the Newton–Kantorovich theorem. *SIAM Journal on Numerical Analysis*, 11(1):10–13, 1974.
- Michael Grant, Stephen Boyd, and Yinyu Ye. Disciplined convex programming. *Global optimization: From theory to implementation*, pp. 155–210, 2006.

- SV Guminov, Yu E Nesterov, PE Dvurechensky, and AV Gasnikov. Accelerated primal-dual gradient descent with linesearch for convex, nonconvex, and nonsmooth optimization problems. In *Doklady Mathematics*, volume 99, pp. 125–128. Springer, 2019.
- Filip Hanzely, Konstantin Mishchenko, and Peter Richtárik. SEGA: Variance reduction via gradient sketching. *Advances in Neural Information Processing Systems*, 31, 2018.
- Samuel Horvath, Chen-Yu Ho, Ludovit Horvath, Atal Narayan Sahu, Marco Canini, and Peter Richtárik. Natural compression for distributed deep learning. *arXiv preprint arXiv:1905.10988*, 2019.
- Samuel Horváth, Dmitry Kovalev, Konstantin Mishchenko, Peter Richtárik, and Sebastian Stich. Stochastic distributed learning with gradient quantization and double-variance reduction. *Optimization Methods and Software*, 38(1):91–106, 2023.
- Rustem Islamov, Xun Qian, and Peter Richtárik. Distributed second order methods with fast rates and compressed communication. In *International conference on machine learning*, pp. 4617–4628. PMLR, 2021.
- Sashank J Reddi, Suvrit Sra, Barnabas Poczos, and Alexander J Smola. Proximal stochastic methods for nonsmooth nonconvex finite-sum optimization. *Advances in neural information processing systems*, 29, 2016.
- Prateek Jain, Purushottam Kar, et al. Non-convex optimization for machine learning. *Foundations and Trends® in Machine Learning*, 10(3-4):142–363, 2017.
- Peter Kairouz, H Brendan McMahan, Brendan Avent, Aurélien Bellet, Mehdi Bennis, Arjun Nitin Bhagoji, Kallista Bonawitz, Zachary Charles, Graham Cormode, Rachel Cummings, et al. Advances and open problems in federated learning. *Foundations and Trends® in Machine Learning*, 14(1–2):1–210, 2021.
- Ahmed Khaled and Peter Richtárik. Better theory for sgd in the nonconvex world. *arXiv preprint arXiv:2002.03329*, 2020.
- Sarit Khirirat, Hamid Reza Feyzmahdavian, and Mikael Johansson. Distributed learning with compressed gradients. *arXiv preprint arXiv:1806.06573*, 2018.
- Jakub Konečný, H Brendan McMahan, Felix X Yu, Peter Richtárik, Ananda Theertha Suresh, and Dave Bacon. Federated learning: Strategies for improving communication efficiency. *arXiv preprint arXiv:1610.05492*, 2016.
- Yin Tat Lee, Zhao Song, and Qiuyi Zhang. Solving empirical risk minimization in the current matrix multiplication time. In *Conference on Learning Theory*, pp. 2140–2157. PMLR, 2019.
- Zhize Li, Dmitry Kovalev, Xun Qian, and Peter Richtárik. Acceleration for compressed gradient descent in distributed and federated optimization. *arXiv preprint arXiv:2002.11364*, 2020.
- Zhize Li, Hongyan Bao, Xiangliang Zhang, and Peter Richtárik. PAGE: A simple and optimal probabilistic gradient estimator for nonconvex optimization. In *International conference on machine learning*, pp. 6286–6295. PMLR, 2021.
- Artavazd Maranjyan, Mher Safaryan, and Peter Richtárik. GradSkip: Communication-Accelerated Local Gradient Methods with Better Computational Complexity. *arXiv preprint arXiv:2210.16402*, 2022.
- H. Brendan McMahan, Eider Moore, Daniel Ramage, Seth Hampson, and Blaise Aguera y Arcas. Communication-efficient learning of deep networks from decentralized data. In *Proceedings of the 20th International Conference on Artificial Intelligence and Statistics (AISTATS)*, 2017. URL <http://arxiv.org/abs/1602.05629>.
- George J Miel. Majorizing sequences and error bounds for iterative methods. *Mathematics of Computation*, 34(149):185–202, 1980.

- Konstantin Mishchenko, Eduard Gorbunov, Martin Takáč, and Peter Richtárik. Distributed learning with compressed gradient differences. *arXiv preprint arXiv:1901.09269*, 2019.
- Konstantin Mishchenko, Grigory Malinovsky, Sebastian Stich, and Peter Richtarik. ProxSkip: Yes! Local gradient steps provably lead to communication acceleration! Finally! In Kamalika Chaudhuri, Stefanie Jegelka, Le Song, Csaba Szepesvari, Gang Niu, and Sivan Sabato (eds.), *Proceedings of the 39th International Conference on Machine Learning*, volume 162 of *Proceedings of Machine Learning Research*, pp. 15750–15769. PMLR, 17–23 Jul 2022.
- Eric Moulines and Francis Bach. Non-asymptotic analysis of stochastic approximation algorithms for machine learning. *Advances in neural information processing systems*, 24, 2011.
- Lianke Qin, Zhao Song, Lichen Zhang, and Danyang Zhuo. An online and unified algorithm for projection matrix vector multiplication with application to empirical risk minimization. In *International Conference on Artificial Intelligence and Statistics*, pp. 101–156. PMLR, 2023.
- Peter Richtárik, Igor Sokolov, and Ilyas Fatkhullin. EF21: A new, simpler, theoretically better, and practically faster error feedback. *Advances in Neural Information Processing Systems*, 34: 4384–4396, 2021.
- Mher Safaryan, Filip Hanzely, and Peter Richtárik. Smoothness matrices beat smoothness constants: Better communication compression techniques for distributed optimization. *Advances in Neural Information Processing Systems*, 34:25688–25702, 2021.
- Mher Safaryan, Egor Shulgin, and Peter Richtárik. Uncertainty principle for communication compression in distributed and federated learning and the search for an optimal compressor. *Information and Inference: A Journal of the IMA*, 11(2):557–580, 2022.
- Zhao Song, Yitan Wang, Zheng Yu, and Lichen Zhang. Sketching for first order method: Efficient algorithm for low-bandwidth channel and vulnerability. In *International Conference on Machine Learning*, pp. 32365–32417. PMLR, 2023.
- Sebastian U Stich. Unified optimal analysis of the (stochastic) gradient method. *arXiv preprint arXiv:1907.04232*, 2019.
- Bokun Wang, Mher Safaryan, and Peter Richtárik. Theoretically better and numerically faster distributed optimization with smoothness-aware quantization techniques. *Advances in Neural Information Processing Systems*, 35:9841–9852, 2022.
- Tetsuro Yamamoto. A convergence theorem for Newton-like methods in Banach spaces. *Numerische Mathematik*, 51:545–557, 1987.
- Jason Yosinski, Jeff Clune, Anh Nguyen, Thomas Fuchs, and Hod Lipson. Understanding neural networks through deep visualization. *arXiv preprint arXiv:1506.06579*, 2015.
- Adams Wei Yu, Lei Huang, Qihang Lin, Ruslan Salakhutdinov, and Jaime Carbonell. Block-normalized gradient method: An empirical study for training deep neural network. *arXiv preprint arXiv:1707.04822*, 2017.
- Jingzhao Zhang, Sai Praneeth Karimireddy, Andreas Veit, Seungyeon Kim, Sashank Reddi, Sanjiv Kumar, and Suvrit Sra. Why are adaptive methods good for attention models? *Advances in Neural Information Processing Systems*, 33:15383–15393, 2020.
- Qinghe Zheng, Xinyu Tian, Nan Jiang, and Mingqiang Yang. Layer-wise learning based stochastic gradient descent method for the optimization of deep convolutional neural network. *Journal of Intelligent & Fuzzy Systems*, 37(4):5641–5654, 2019.
- Luisa M Zintgraf, Taco S Cohen, Tameem Adel, and Max Welling. Visualizing deep neural network decisions: Prediction difference analysis. *arXiv preprint arXiv:1702.04595*, 2017.