
Det-CGD: Compressed Gradient Descent with Matrix Stepsizes for Non-Convex Optimization

Anonymous Author(s)

Affiliation

Address

email

Abstract

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

12 1 Introduction

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

23 In this paper, we consider the general minimization problem:

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

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

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

27 The stochastic gradient descent (SGD) algorithm [MB11, B⁺15, GLQ⁺19] is one of the most
28 common algorithms to solve this problem. In its most general form, it can be written as

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

29 where $g(x^k)$ is a stochastic estimator of $\nabla f(x^k)$ and $\gamma > 0$ is a positive scalar stepsize. A particular
30 case of interest is the compressed gradient descent (CGD) algorithm [KFJ18], where the estimator g

31 is taken as a compressed alternative of the initial gradient:

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

32 and the compressor \mathcal{C} is chosen to be a "sparser" estimator that aims to reduce the communication
 33 overhead in distributed or federated settings. This is crucial, as highlighted in the seminal paper
 34 by [KMY⁺16], which showed that the bottleneck of distributed optimization algorithms is the
 35 communication complexity. In order to deal with the limited resources of current devices, there are
 36 various compression objectives that are practical to achieve. These include also compressing the
 37 model broadcasted from server to clients for local training, and reducing the computational burden
 38 of local training. These objectives are mostly complementary, but compressing gradients has the
 39 potential for the greatest practical impact due to slower upload speeds of client connections and the
 40 benefits of averaging [KMA⁺21]. In this paper we will focus on this latter problem.

41 An important subclass of compressors are the sketches. Sketches are linear operators defined on
 42 \mathbb{R}^d , i.e., $\mathcal{C}(y) = Sy$ for every $y \in \mathbb{R}^d$, where S is a random matrix. A standard example of such
 43 a compressor is the Rand- k compressor, which randomly chooses k entries of its argument and
 44 scales them with a scalar multiplier to make the estimator unbiased. Instead of communicating all d
 45 coordinates of the gradient, one communicates only a subset of size k , thus reducing the number of
 46 communicated bits by a factor of d/k . Formally, Rand- k is defined as follows: $S = \sum_{j=1}^k \frac{d}{k} e_{i_j} e_{i_j}^\top$,
 47 where i_j are the selected coordinates of the input vector. We refer the reader to [SSR22] for an
 48 overview on compressions.

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

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

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

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

54 The assumption of matrix smoothness, which is a generalization of scalar smoothness, has been
 55 shown to be a more powerful tool for improving supervised model training. In [SHR21], the authors
 56 proposed using smoothness matrices and suggested a novel communication sparsification strategy to
 57 reduce communication complexity in distributed optimization. The technique was adapted to three
 58 distributed optimization algorithms in the convex setting, resulting in significant communication
 59 complexity savings and consistently outperforming the baselines. The results of this study demonstrate
 60 the efficacy of the matrix smoothness assumption in improving distributed optimization algorithms.

61 The case of block-diagonal smoothness matrices is particularly relevant in various applications, such
 62 as neural networks (NN). In this setting, each block corresponds to a layer of the network, and we
 63 characterize the smoothness with respect to nodes in the i -th layer by a corresponding matrix L_i .
 64 Unlike in the scalar setting, we favor the similarity of certain entries of the argument over the others.
 65 This is because the information carried by the layers becomes more complex, while the nodes in the
 66 same layers are similar. This phenomenon has been observed visually in various studies, such as
 67 those by [YCN⁺15] and [ZCAW17].

68 Another motivation for using a layer-dependent stepsize has its roots in physics. In nature, the
 69 propagation speed of light in media of different densities varies due to frequency variations. Similarly,
 70 different layers in neural networks carry different information, metric systems, and scaling. Thus, the
 71 stepsizes need to be picked accordingly to achieve optimal convergence.

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

77 1.1 Related work

78 Many successful convex optimization techniques have been adapted for use in the non-convex
 79 setting. Here is a non-exhaustive list: adaptivity [DOG⁺19, ZKV⁺20], variance reduction [JRSPS16,

LBZR21], and acceleration [GNDG19]. A paper of particular importance for our work is that of [KR20], 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 [JK⁺17, DDG⁺22].

A classical example of a matrix stepsize method is Newton’s method. This method has been popular in the optimization community for a long time [GT74, Mie80, Yam87]. 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 [Bro65, DM77, ABK07, ABSM14]. An example is the Newton-Star algorithm [IQR21], which we discuss in Section 2.

[GR15] 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, [HMR18] proposed a variance-reduced version of this method.

Leveraging the layer-wise structure of neural networks has been widely studied for optimizing the training loss function. For example, [ZTJY19] propose SGD with different scalar stepsizes for each layer, [YHL⁺17, GCH⁺19] propose layer-wise normalization for Stochastic Normalized Gradient Descent, and [DBA⁺20, WSR22] propose layer-wise compression in the distributed setting.

DCGD, proposed by [KFJ18], has since been improved in various ways, such as in [HHH⁺19, LKQR20]. There is also a large body of literature on other federated learning algorithms with unbiased compressors [AGL⁺17, MGTR19, GBLR21, MMSR22, MSR22, HKM⁺23].

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 shows 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 the that the server-to-client communication is less expensive [KMY⁺16, KMA⁺21], 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 I_d we denote the $d \times d$ identity matrix, and by 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 $Q \in \mathbb{S}_{++}^d$ and $x \in \mathbb{R}^d$, we write $\|x\|_Q := \sqrt{\langle Qx, x \rangle}$, where $\langle \cdot, \cdot \rangle$ is the standard Euclidean inner product on \mathbb{R}^d . For a matrix $A \in \mathbb{S}_{++}^d$, we define by $\lambda_{\max}(A)$ (resp. $\lambda_{\min}(A)$) the largest (resp. smallest) eigenvalue of the matrix A . Let $A_i \in \mathbb{R}^{d_i \times d_i}$ and $d = d_1 + \dots + d_\ell$. Then the matrix $A = \text{Diag}(A_1, \dots, A_\ell)$ is defined as a block diagonal $d \times d$ matrix where the i -th block is equal to A_i . We will use $\text{diag}(A) \in \mathbb{R}^{d \times d}$ to denote the diagonal of any matrix $A \in \mathbb{R}^{d \times d}$. 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 $\nabla^2 f(x)$.

2 The algorithms

Below we define our two main algorithms:

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

and

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

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

$$S^k, T^k \in \mathbb{S}_+^d, \quad S^k \stackrel{iid}{\sim} S \quad \text{and} \quad T^k \stackrel{iid}{\sim} T \\ \mathbb{E}[S^k] = \mathbb{E}[T^k] = I_d, \quad \text{for every } k \in \mathbb{N}.$$

A simple instance of [det-CGD1](#) and [det-CGD2](#) is the vanilla GD. Indeed, if $S^k = T^k = I_d$ and $D = \gamma I_d$, 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 [\[IQR21\]](#):

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

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 D depending on the smoothness matrix 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 S and the stepsize D are commutative w.r.t. matrix product, the algorithms become equivalent. In general, a simple calculation shows that if we take

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

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

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

However, in general $D \mathbb{E}[S^k] D^{-1}$ 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:

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

and

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

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 $k \geq 0$*

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

if one of the below conditions is true:

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

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

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

166 This normalization is meaningful because adjusting the weight matrix to $\frac{\mathbf{D}}{\det(\mathbf{D})^{1/d}}$ allows its determi-
 167 nant to be 1, making the norm on the left-hand side comparable to the standard Euclidean norm. It is
 168 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
 169 not depend on the choice of $\mathbf{D} \in \mathbb{S}_{++}^d$. Therefore, the results of (9) are comparable across different
 170 \mathbf{D} in the sense that the right-hand side of (9) measures the volume of the ellipsoid containing the
 171 gradient.

172 3.1 Optimal matrix stepsize

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

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

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

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

178 This explicit condition on \mathbf{D} can assist in determining the optimal stepsize. Since both \mathbf{D} and
 179 $(\mathbf{T}^k \mathbf{L} \mathbf{T}^k)^{-1}$ are positive definite, then the right-hand side of (10) is minimized exactly when

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

180 The situation is different for [det-CGD1](#). According to (10), the optimal \mathbf{D} is defined as the solution
 181 of the following constrained optimization problem:

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

182

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

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

189 One could consider using the CVXPY [DB16] package to solve (13), provided that it is first transformed
 190 into a Disciplined Convex Programming (DCP) form [GBY06]. Nevertheless, (7) is not recognized
 191 as a DCP constraint in the general case. To make CVXPY applicable, additional steps tailored to the
 192 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 D_i here for **det-CGD1** is W_i with the optimal scaling determined using Theorem 2, for **det-CGD2** it is the optimal stepsize matrix defined in (11). 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 $\tilde{L}_{i,k}$ is defined as $\frac{d-k}{d-1} \text{diag}(\mathbf{L}_i) + \frac{k-1}{d-1} \mathbf{L}_i$.

No.	The method	(S_i^k, 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	$(I_d, \gamma \mathbf{L}_i^{-1})$	$d \cdot \det(\mathbf{L})^{1/d}$	$d \cdot \det(\mathbf{L})^{1/d}$
2.	det-CGD1	$(I_d, \gamma \text{diag}^{-1}(\mathbf{L}_i))$	$d \cdot \det(\text{diag}(\mathbf{L}))^{1/d}$	$d \cdot \det(\text{diag}(\mathbf{L}))^{1/d}$
3.	det-CGD1	$(I_d, \gamma I_{d_i})$	$d \cdot \left(\prod_{i=1}^{\ell} \lambda_{\max}^{d_i}(\mathbf{L}_i) \right)^{1/d}$	$d \cdot \lambda_{\max}(\mathbf{L})$
4.	det-CGD1	$(\text{rand-1}, \gamma I_{d_i})$	$\ell \cdot \left(\prod_{i=1}^{\ell} d_i^{d_i} (\max_j (\mathbf{L}_i)_{jj})^{d_i} \right)^{1/d}$	$d \cdot \max_j (\mathbf{L}_{jj})$
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	$(\text{rand-1}, \gamma \text{diag}^{-1}(\mathbf{L}_i))$	$\ell \cdot \left(\frac{\prod_{i=1}^{\ell} d_i^{d_i}}{\prod_{j=1}^d (\mathbf{L}_{jj})^{1/d}} \right)^{1/d}$	$d \cdot \det(\text{diag}(\mathbf{L}))^{1/d}$
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	(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	$(\text{rand-}k, \frac{k_i}{d_i} \tilde{\mathbf{L}}_{i,k_i}^{-1})$	$k \cdot \left(\prod_{i=1}^{\ell} \left(\frac{d_i}{k_i} \right)^{\frac{d_i}{d}} \right) \left(\prod_{i=1}^{\ell} \det(\tilde{\mathbf{L}}_{i,k_i}) \right)^{1/d}$	$d \cdot \det(\tilde{\mathbf{L}}_{1,k})$
12.	det-CGD2	$(\text{Bern-}q_i, q_i \mathbf{L}_i^{-1})$	$\left(\sum_{i=1}^{\ell} q_i d_i \right) \cdot \prod_{i=1}^{\ell} \left(\frac{1}{q_i} \right)^{\frac{d_i}{d}} \det(\mathbf{L})^{1/d}$	$d \cdot \det(\mathbf{L})^{1/d}$
13.	GD	$(I_d, \lambda_{\max}^{-1}(\mathbf{L}) I_d)$	N/A	$d \cdot \lambda_{\max}(\mathbf{L})$

193 4 Leveraging the layer-wise structure

194 In this section we focus on the block-diagonal case of \mathbf{L} for both **det-CGD1** and **det-CGD2**. In
195 particular, we propose hyper-parameters of **det-CGD1** designed specifically for training NNs. Let
196 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
197 classical smoothness condition, as in the latter case $\mathbf{L}_i = L I_{d_i}$ for all $i = 1, \dots, \ell$. Respectively,
198 we choose both the sketches and the stepsize to be block diagonal: $\mathbf{D} = \text{Diag}(\mathbf{D}_1, \dots, \mathbf{D}_{\ell})$ and
199 $\mathbf{S}^k = \text{Diag}(\mathbf{S}_1^k, \dots, \mathbf{S}_{\ell}^k)$, where $\mathbf{D}_i, \mathbf{S}_i^k \in \mathbb{S}_{++}^{d_i}$.

200 Let us notice that the left hand side of the inequality constraint in (13) has quadratic dependence on
201 \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

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

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

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

203 The following theorem is based on this simple fact applied to the corresponding blocks of the matrices
204 $\mathbf{D}, \mathbf{L}, \mathbf{S}^k$ for **det-CGD1**.

205 **Theorem 2.** Let $f : \mathbb{R}^d \rightarrow \mathbb{R}$ satisfy Assumptions 1 and 2, with \mathbf{L} admitting the layer-separable struc-
206 ture $\mathbf{L} = \text{Diag}(\mathbf{L}_1, \dots, \mathbf{L}_{\ell})$, where $\mathbf{L}_1, \dots, \mathbf{L}_{\ell} \in \mathbb{S}_{++}^{d_i}$. Choose random matrices $\mathbf{S}_1^k, \dots, \mathbf{S}_{\ell}^k \in \mathbb{S}_{++}^{d_i}$
207 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
208 matrices $\mathbf{W}_1, \dots, \mathbf{W}_{\ell} \in \mathbb{S}_{++}^{d_i}$ and scalars $\gamma_1, \dots, \gamma_{\ell} > 0$ such that

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

209 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

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

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

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

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

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

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

220 Jensen's inequality implies that

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

221 The equality is attained when $q_i = q$ for all $i \in [\ell]$. The expected bits transferred per iteration of
 222 this algorithm is then equal to $k_{\text{exp}} = qd$ and the communication complexity equals $d \det(\mathbf{L})^{1/d}$.
 223 Comparing with the results for **det-CGD2** with $\text{rand-}k_{\text{exp}}$ on row 11 and using the fact that $\det(\mathbf{L}) \leq$
 224 $\det(\text{diag}(\mathbf{L}))$, we deduce that the Bernoulli scheme is better than the uniform sampling scheme.
 225 Notice also, the communication complexity matches the one for the uncompressed **det-CGD2**
 226 displayed on row 9. This, in particular means that using the Bern- q sketches we can compress the
 227 gradients for free. The latter means that we reduce the number of bits broadcasted at each iteration
 228 without losing in the total communication complexity. In particular, when all the layers have the same
 229 width d_i , the number of broadcasted bits for each iteration is reduced by a factor of q .

230 5 Distributed setting

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

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

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

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

237 This assumption also implies that f is of matrix smoothness with $\bar{\mathbf{L}} \in \mathbb{S}_{++}^d$, where $\bar{\mathbf{L}} = \frac{1}{n} \sum_{i=1}^n \mathbf{L}_i$.
 238 Following the standard FL framework [KMY⁺16, MMR⁺17, KFJ18], we assume that the i -th
 239 component function f_i is stored on the i -th client. At each iteration, the clients in parallel compute
 240 and compress the local gradient ∇f_i and communicate it to the central server. The server, then
 241 aggregates the compressed gradients, computes the next iterate, and in parallel broadcasts it to the
 242 clients. See the algorithms below for the pseudo-codes.

Algorithm 1 Distributed **det-CGD1**

1: **Input:** Starting point x_0 , stepsize matrix D ,
 number of iterations K
 2: **for** $k = 0, 1, 2, \dots, K - 1$ **do**
 3: The devices in parallel:
 4: sample $S_i^k \sim \mathcal{S}$;
 5: compute $S_i^k \nabla f_i(x_k)$;
 6: broadcast $S_i^k \nabla f_i(x_k)$.
 7: The server:
 8: combines $g_k = \frac{D}{n} \sum_i S_i^k \nabla f_i(x_k)$;
 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 D ,
 number of iterations K
 2: **for** $k = 0, 1, 2, \dots, K - 1$ **do**
 3: The devices in parallel:
 4: sample $T_i^k \sim \mathcal{T}$;
 5: compute $T_i^k D \nabla f_i(x_k)$;
 6: broadcast $T_i^k D \nabla f_i(x_k)$.
 7: The server:
 8: combines $g_k = \frac{1}{n} \sum_i T_i^k D \nabla f_i(x_k)$;
 9: computes $x_{k+1} = x_k - g_k$;
 10: broadcasts x_{k+1} .
 11: **end for**
 12: **Return:** x_K

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 L . If the stepsize satisfies

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

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

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

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

$$\lambda_D := \max_i \left\{ \lambda_{\max} \left(\mathbb{E} \left[L_i^{\frac{1}{2}} (S_i^k - I_d) DLD (S_i^k - I_d) L_i^{\frac{1}{2}} \right] \right) \right\}.$$

The same result is true for Algorithm 2 with a different constant λ_D . The proof of Theorem 3 and its analogue for Algorithm 2 are presented in the Appendix. The analysis is largely inspired by [KR20, 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, $1 + \lambda_D/n$, depends on the stepsize D . Furthermore, it has a second-order dependence on D , implying that $\lambda_{\alpha D} = \alpha^2 \lambda_D$, as opposed to $\det(\alpha D)^{1/d}$, which is linear in α . Therefore, we can choose a small enough coefficient α to ensure that λ_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 denominator 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 error level of ε^2 in (20) if the following conditions are satisfied:

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

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

$$\begin{aligned} & \text{minimize} && \log \det(D^{-1}) \\ & \text{subject to} && D \text{ satisfies (21).} \end{aligned}$$

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 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.

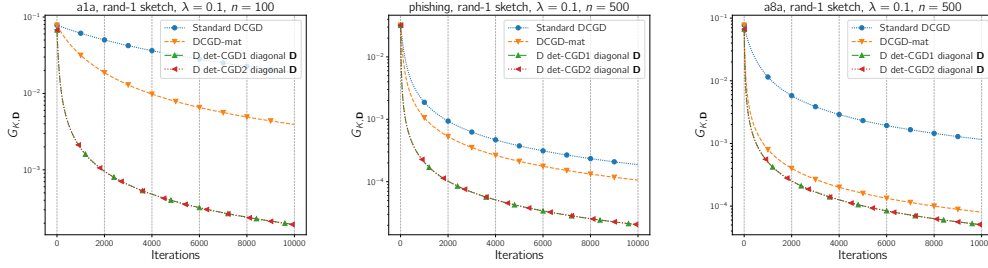


Figure 1: Comparison of standard DCGD, DCGD with matrix smoothness, D-det-CGD1 and D-det-CGD2 with optimal diagonal stepsizes under rand-1 sketch. The stepsize for standard DCGD is determined using [KR20, Proposition 4], the stepsize for DCGD with matrix smoothness along with D_1 , D_2 is determined using Corollary 1, the error level is set to be $\varepsilon^2 = 0.0001$. Here $G_{K,D} := \frac{1}{K} \left(\sum_{k=0}^{K-1} \|\nabla f(x^k)\|_{D/\det(D)^{1/d}}^2 \right)$.

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 λ_D to be upper-bounded by n/K . Thus, 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 [KR20] for the scalar smooth DCGD with scalar stepsizes. This means that the parameters of algorithms are $L_i = L_i I_d$, $L = L I_d$, $D = \gamma I_d$, $\omega = \lambda_{\max} \left(\mathbb{E} \left[(S_i^k)^\top S_i^k \right] \right) - 1$. One may check that (21) reduces to

$$\gamma \leq \min \left\{ \frac{1}{L}, \sqrt{\frac{n}{K L_{\max} L \omega}}, \frac{n \varepsilon^2}{4 \Delta^{\text{inf}} L_{\max} L \omega} \right\} \quad \text{and} \quad K \gamma \geq \frac{12(f(x^0) - f^{\text{inf}})}{\varepsilon^2} \quad (22)$$

As expected, this coincides with the results from [KR20, Corollary 1]. See the Appendix for the details on the analysis of [KR20]. 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 stepsized DCGD. For more details on the experiments we refer the reader to the corresponding section in the Appendix.

6 Conclusion

6.1 Limitations

It is worth noting that every point in \mathbb{R}^d can be enclosed within some volume 1 ellipsoid. To see this, let $0 \neq v \in \mathbb{R}^d$ and define $Q := \frac{\alpha}{\|v\|^2} v v^\top + \beta \sum_{i=1}^d v_i v_i^\top$, where $v_1 = \frac{v}{\|v\|}$, v_2, \dots, v_d form an orthonormal basis. The eigenvalues of Q are β (with multiplicity $d-1$) and α (with multiplicity 1), so we have $\det(Q) = \beta^{d-1} \alpha \leq 1$. Furthermore, we have $\|v\|_Q^2 = v^\top Q v = \alpha \|v\|^2$. By choosing $\alpha = \frac{1}{\|v\|^2}$ and $\beta = \|v\|^{2/(d-1)}$, we can obtain $\det(Q) = 1$ while $\|v\|_Q^2 \leq 1$. Therefore, having the average D -norm of the gradient bounded by a small number does not guarantee that the average Euclidean norm is small. This implies that the theory does not guarantee stationarity in the Euclidean sense.

6.2 Future work

Matrix stepsize gradient methods are still not well studied and require further analysis. Although many important algorithms have been proposed using scalar stepsizes and are known to have good performance, their matrix analogs have yet to be thoroughly examined. The distributed algorithms proposed in Section 5 follow the structure of DCGD by [KFJ18]. However, other federated learning mechanisms such as MARINA, which has variance reduction [GBLR21], or EF21 by [RSF21], which has powerful practical performance, should also be explored.

References

- [ABK07] 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.
- [ABSM14] 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.
- [AGL⁺17] 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.
- [B⁺15] Sébastien Bubeck et al. Convex optimization: Algorithms and complexity. *Foundations and Trends® in Machine Learning*, 8(3-4):231–357, 2015.
- [Bro65] Charles G Broyden. A class of methods for solving nonlinear simultaneous equations. *Mathematics of computation*, 19(92):577–593, 1965.
- [CL11] 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.
- [DB16] 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.
- [DBA⁺20] 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, pages 3817–3824, 2020.
- [DDG⁺22] 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*, pages 79–163. Springer, 2022.
- [DM77] John E Dennis, Jr and Jorge J Moré. Quasi-Newton methods, motivation and theory. *SIAM review*, 19(1):46–89, 1977.
- [DOG⁺19] 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.
- [GBLR21] Eduard Gorbunov, Konstantin P Burlachenko, Zhize Li, and Peter Richtárik. MARINA: Faster non-convex distributed learning with compression. In *International Conference on Machine Learning*, pages 3788–3798. PMLR, 2021.
- [GBY06] Michael Grant, Stephen Boyd, and Yinyu Ye. Disciplined convex programming. *Global optimization: From theory to implementation*, pages 155–210, 2006.
- [GCH⁺19] Boris Ginsburg, Patrice Castonguay, Oleksii Hrinchuk, Oleksii Kuchaiev, Vitaly Lavrukhin, 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.
- [GLQ⁺19] 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*, pages 5200–5209. PMLR, 2019.
- [GNDG19] 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, pages 125–128. Springer, 2019.

348 [GR15] Robert M Gower and Peter Richtárik. Randomized iterative methods for linear systems.
349 *SIAM Journal on Matrix Analysis and Applications*, 36(4):1660–1690, 2015.

350 [GT74] William B Gragg and Richard A Tapia. Optimal error bounds for the Newton–
351 Kantorovich theorem. *SIAM Journal on Numerical Analysis*, 11(1):10–13, 1974.

352 [HHH⁺19] Samuel Horváth, Chen-Yu Ho, Ludovit Horvath, Atal Narayan Sahu, Marco Canini,
353 and Peter Richtárik. Natural compression for distributed deep learning. *CoRR*,
354 abs/1905.10988, 2019.

355 [HKM⁺23] Samuel Horváth, Dmitry Kovalev, Konstantin Mishchenko, Peter Richtárik, and Sebas-
356 tian Stich. Stochastic distributed learning with gradient quantization and double-variance
357 reduction. *Optimization Methods and Software*, 38(1):91–106, 2023.

358 [HMR18] Filip Hanzely, Konstantin Mishchenko, and Peter Richtárik. SEGA: Variance reduction
359 via gradient sketching. *Advances in Neural Information Processing Systems*, 31, 2018.

360 [IQR21] Rustem Islamov, Xun Qian, and Peter Richtárik. Distributed second order methods
361 with fast rates and compressed communication. In *International conference on machine*
362 *learning*, pages 4617–4628. PMLR, 2021.

363 [JK⁺17] Prateek Jain, Purushottam Kar, et al. Non-convex optimization for machine learning.
364 *Foundations and Trends® in Machine Learning*, 10(3-4):142–363, 2017.

365 [JRSPS16] Sashank J Reddi, Suvrit Sra, Barnabas Poczos, and Alexander J Smola. Proximal
366 stochastic methods for nonsmooth nonconvex finite-sum optimization. *Advances in*
367 *neural information processing systems*, 29, 2016.

368 [KFJ18] Sarit Khirirat, Hamid Reza Feyzmahdavian, and Mikael Johansson. Distributed learning
369 with compressed gradients. *arXiv preprint arXiv:1806.06573*, 2018.

370 [KMA⁺21] Peter Kairouz, H Brendan McMahan, Brendan Avent, Aurélien Bellet, Mehdi Bennis,
371 Arjun Nitin Bhagoji, Kallista Bonawitz, Zachary Charles, Graham Cormode, Rachel
372 Cummings, et al. Advances and open problems in federated learning. *Foundations and*
373 *Trends® in Machine Learning*, 14(1–2):1–210, 2021.

374 [KMY⁺16] Jakub Konečný, H Brendan McMahan, Felix X Yu, Peter Richtárik, Ananda Theertha
375 Suresh, and Dave Bacon. Federated learning: Strategies for improving communication
376 efficiency. *arXiv preprint arXiv:1610.05492*, 2016.

377 [KR20] Ahmed Khaled and Peter Richtárik. Better theory for SGD in the nonconvex world.
378 *arXiv preprint arXiv:2002.03329*, 2020.

379 [LBZR21] Zhize Li, Hongyan Bao, Xiangliang Zhang, and Peter Richtárik. PAGE: A simple and
380 optimal probabilistic gradient estimator for nonconvex optimization. In *International*
381 *conference on machine learning*, pages 6286–6295. PMLR, 2021.

382 [LKQR20] Zhize Li, Dmitry Kovalev, Xun Qian, and Peter Richtárik. Acceleration for com-
383 pressed gradient descent in distributed and federated optimization. *arXiv preprint*
384 *arXiv:2002.11364*, 2020.

385 [MB11] Eric Moulines and Francis Bach. Non-asymptotic analysis of stochastic approximation
386 algorithms for machine learning. *Advances in neural information processing systems*,
387 24, 2011.

388 [MGTR19] Konstantin Mishchenko, Eduard Gorbunov, Martin Takáč, and Peter Richtárik.
389 Distributed learning with compressed gradient differences. *arXiv preprint*
390 *arXiv:1901.09269*, 2019.

391 [Mie80] George J Miel. Majorizing sequences and error bounds for iterative methods. *Mathe-*
392 *matics of Computation*, 34(149):185–202, 1980.

393 [MMR⁺17] H. Brendan McMahan, Eider Moore, Daniel Ramage, Seth Hampson, and Blaise Aguera
394 y Arcas. Communication-efficient learning of deep networks from decentralized data. In
395 *Proceedings of the 20th International Conference on Artificial Intelligence and Statistics*
396 (AISTATS), 2017.

397 [MMSR22] Konstantin Mishchenko, Grigory Malinovsky, Sebastian Stich, and Peter Richtárik.
398 ProxSkip: Yes! Local gradient steps provably lead to communication acceleration!
399 Finally! In Kamalika Chaudhuri, Stefanie Jegelka, Le Song, Csaba Szepesvari, Gang
400 Niu, and Sivan Sabato, editors, *Proceedings of the 39th International Conference on*
401 *Machine Learning*, volume 162 of *Proceedings of Machine Learning Research*, pages
402 15750–15769. PMLR, 17–23 Jul 2022.

403 [MSR22] Artavazd Maranjyan, Mher Safaryan, and Peter Richtárik. GradSkip: Communication-
404 Accelerated Local Gradient Methods with Better Computational Complexity. *arXiv*
405 *preprint arXiv:2210.16402*, 2022.

406 [RSF21] Peter Richtárik, Igor Sokolov, and Ilyas Fatkhullin. EF21: A new, simpler, theoretically
407 better, and practically faster error feedback. *Advances in Neural Information Processing*
408 *Systems*, 34:4384–4396, 2021.

409 [SHR21] Mher Safaryan, Filip Hanzely, and Peter Richtárik. Smoothness matrices beat smooth-
410 ness constants: Better communication compression techniques for distributed optimiza-
411 tion. *Advances in Neural Information Processing Systems*, 34:25688–25702, 2021.

412 [SSR22] Mher Safaryan, Egor Shulgin, and Peter Richtárik. Uncertainty principle for communi-
413 cation compression in distributed and federated learning and the search for an optimal
414 compressor. *Information and Inference: A Journal of the IMA*, 11(2):557–580, 2022.

415 [Sti19] Sebastian U Stich. Unified optimal analysis of the (stochastic) gradient method. *arXiv*
416 *preprint arXiv:1907.04232*, 2019.

417 [WSR22] Bokun Wang, Mher Safaryan, and Peter Richtárik. Theoretically better and numeri-
418 cally faster distributed optimization with smoothness-aware quantization techniques.
419 *Advances in Neural Information Processing Systems*, 35:9841–9852, 2022.

420 [Yam87] Tetsuro Yamamoto. A convergence theorem for newton-like methods in banach spaces.
421 *Numerische Mathematik*, 51:545–557, 1987.

422 [YCN⁺15] Jason Yosinski, Jeff Clune, Anh Nguyen, Thomas Fuchs, and Hod Lipson. Under-
423 standing neural networks through deep visualization. *arXiv preprint arXiv:1506.06579*,
424 2015.

425 [YHL⁺17] Adams Wei Yu, Lei Huang, Qihang Lin, Ruslan Salakhutdinov, and Jaime Carbonell.
426 Block-normalized gradient method: An empirical study for training deep neural network.
427 *arXiv preprint arXiv:1707.04822*, 2017.

428 [ZCAW17] Luisa M Zintgraf, Taco S Cohen, Tameem Adel, and Max Welling. Visualizing deep neu-
429 ral network decisions: Prediction difference analysis. *arXiv preprint arXiv:1702.04595*,
430 2017.

431 [ZKV⁺20] Jingzhao Zhang, Sai Praneeth Karimireddy, Andreas Veit, Seungyeon Kim, Sashank
432 Reddi, Sanjiv Kumar, and Suvrit Sra. Why are adaptive methods good for attention
433 models? *Advances in Neural Information Processing Systems*, 33:15383–15393, 2020.

434 [ZTJY19] Qinghe Zheng, Xinyu Tian, Nan Jiang, and Mingqiang Yang. Layer-wise learning based
435 stochastic gradient descent method for the optimization of deep convolutional neural
436 network. *Journal of Intelligent & Fuzzy Systems*, 37(4):5641–5654, 2019.

437 Contents

438	A Single node case	14
439	A.1 Proof of Theorem 1	14
440	A.2 Proof of Proposition 1	15
441	B Layer-wise case	16
442	B.1 Proof of Theorem 2	16
443	B.2 Bernoulli sketch for det-CGD2	16
444	B.3 General cases for det-CGD1	18
445	B.4 General cases for det-CGD2	18
446	B.5 Interpretations of Table 1	18
447	B.5.1 Comparison of row 5 and 7	18
448	B.5.2 Comparison of row 6 and 7	19
449	C Distributed case	19
450	C.1 Proof of Theorem 3	19
451	C.2 Convexity of the constraints	22
452	C.2.1 Proof of Corollary 1	24
453	C.3 Distributed det-CGD2	24
454	C.3.1 Analysis of distributed det-CGD2	24
455	C.3.2 Optimal stepsize	26
456	C.4 DCGD with constant stepsize	26
457	D Proofs of technical lemmas	27
458	D.1 Proof of Lemma 1	27
459	D.2 Proof of Lemma 2	28
460	D.3 Proof of Lemma 3	28
461	D.4 Proof of Lemma 4	28
462	D.5 Proof of Lemma 5	28
463	D.6 Proof of Lemma 6	29
464	D.7 Proof of Lemma 7	29
465	E Experiments	29
466	E.1 Single node case	29
467	E.1.1 Comparison to CGD with scalar stepsize, scalar smoothness constant . . .	30
468	E.1.2 Comparison of the two algorithms under the same stepsize	31
469	E.2 Distributed case	31
470	E.2.1 Comparison to standard DCGD in the distributed case	31

471 A Single node case

472 A.1 Proof of Theorem 1

473 i) Using Assumption 2 with $x = x^{k+1} = x^k - \mathbf{D}\mathbf{S}^k \nabla f(x^k)$ and $y = x^k$, we get

$$\begin{aligned} & \mathbb{E} [f(x^{k+1}) \mid x^k] \\ & \leq \mathbb{E} \left[f(x^k) + \langle \nabla f(x^k), -\mathbf{D}\mathbf{S}^k \nabla f(x^k) \rangle + \frac{1}{2} \langle \mathbf{L}(-\mathbf{D}\mathbf{S}^k \nabla f(x^k)), -\mathbf{D}\mathbf{S}^k \nabla f(x^k) \rangle \mid x^k \right] \\ & = f(x^k) - \langle \nabla f(x^k), \mathbf{D}\mathbb{E} [\mathbf{S}^k] \nabla f(x^k) \rangle + \frac{1}{2} \langle \mathbb{E} [\mathbf{S}^k \mathbf{D} \mathbf{L} \mathbf{D} \mathbf{S}^k] \nabla f(x^k), \nabla f(x^k) \rangle. \end{aligned}$$

474 From the unbiasedness of the sketch \mathbf{S}^k

$$\begin{aligned} & \mathbb{E} [f(x^{k+1}) \mid x^k] \\ & \leq f(x^k) - \langle \nabla f(x^k), \mathbf{D}\nabla f(x^k) \rangle + \frac{1}{2} \langle \mathbb{E} [\mathbf{S}^k \mathbf{D} \mathbf{L} \mathbf{D} \mathbf{S}^k] \nabla f(x^k), \nabla f(x^k) \rangle \\ & \stackrel{(7)}{\leq} f(x^k) - \langle \nabla f(x^k), \mathbf{D}\nabla f(x^k) \rangle + \frac{1}{2} \langle \mathbf{D}\nabla f(x^k), \nabla f(x^k) \rangle \\ & = f(x^k) - \frac{1}{2} \langle \nabla f(x^k), \mathbf{D}\nabla f(x^k) \rangle \\ & = f(x^k) - \frac{1}{2} \|\nabla f(x^k)\|_{\mathbf{D}}^2. \end{aligned} \tag{23}$$

475 Next, by subtracting f^{inf} from both sides of (23), taking expectation and applying the tower property,
476 we get

$$\begin{aligned} \mathbb{E} [f(x^{k+1})] - f^{\text{inf}} &= \mathbb{E} [\mathbb{E} [f(x^{k+1}) \mid x^k]] - f^{\text{inf}} \\ & \stackrel{(23)}{\leq} \mathbb{E} \left[f(x^k) - \frac{1}{2} \|\nabla f(x^k)\|_{\mathbf{D}}^2 \right] - f^{\text{inf}} \\ &= \mathbb{E} [f(x^k)] - f^{\text{inf}} - \frac{1}{2} \mathbb{E} [\|\nabla f(x^k)\|_{\mathbf{D}}^2]. \end{aligned}$$

477 Letting $\Delta^k := \mathbb{E} [f(x^k)] - f^{\text{inf}}$, the last inequality can be written as $\Delta^{k+1} \leq \Delta^k -$
478 $\frac{1}{2} \mathbb{E} [\|\nabla f(x^k)\|_{\mathbf{D}}^2]$. Summing these inequalities for $k = 0, 1, \dots, K-1$, we get a telescoping
479 effect leading to

$$\Delta^K \leq \Delta^0 - \frac{1}{2} \sum_{k=0}^{K-1} \mathbb{E} [\|\nabla f(x^k)\|_{\mathbf{D}}^2].$$

480 It remains to rearrange the terms of this inequality, divide both sides by $K \det(\mathbf{D})^{1/d}$, and use the
481 inequality $\Delta^K \geq 0$.

482 ii) Similar to the previous case, using matrix smoothness for $x = x^{k+1} = x^k - \mathbf{T}^k \mathbf{D} \nabla f(x^k)$ and
483 $y = x^k$, we get

$$\begin{aligned} & \mathbb{E} [f(x^{k+1}) \mid x^k] \\ & \leq \mathbb{E} \left[f(x^k) + \langle \nabla f(x^k), -\mathbf{T}^k \mathbf{D} \nabla f(x^k) \rangle + \frac{1}{2} \langle \mathbf{L}(-\mathbf{T}^k \mathbf{D} \nabla f(x^k)), -\mathbf{T}^k \mathbf{D} \nabla f(x^k) \rangle \mid x^k \right] \\ & = f(x^k) - \langle \nabla f(x^k), \mathbb{E} [\mathbf{T}^k] \mathbf{D} \nabla f(x^k) \rangle + \frac{1}{2} \langle \mathbb{E} [\mathbf{D} (\mathbf{T}^k)^\top \mathbf{L} \mathbf{T}^k \mathbf{D}] \nabla f(x^k), \nabla f(x^k) \rangle. \end{aligned}$$

484 From Assumption 3 and condition (8) we deduce

$$\begin{aligned} \mathbb{E} [f(x^{k+1}) \mid x^k] &\leq f(x^k) - \langle \nabla f(x^k), \mathbf{D}\nabla f(x^k) \rangle + \frac{1}{2} \langle \mathbf{D}\nabla f(x^k), \nabla f(x^k) \rangle \\ &= f(x^k) - \frac{1}{2} \langle \nabla f(x^k), \mathbf{D}\nabla f(x^k) \rangle \\ &= f(x^k) - \frac{1}{2} \|\nabla f(x^k)\|_{\mathbf{D}}^2. \end{aligned} \tag{24}$$

Next, by subtracting f^{inf} from both sides of (24), taking expectation and applying the tower property, we get

$$\begin{aligned}\mathbb{E}[f(x^{k+1})] - f^{\text{inf}} &= \mathbb{E}[\mathbb{E}[f(x^{k+1}) | x^k]] - f^{\text{inf}} \\ &\stackrel{(24)}{\leq} \mathbb{E}\left[f(x^k) - \frac{1}{2}\|\nabla f(x^k)\|_{\mathbf{D}}^2\right] - f^{\text{inf}} \\ &= \mathbb{E}[f(x^k)] - f^{\text{inf}} - \frac{1}{2}\mathbb{E}[\|\nabla f(x^k)\|_{\mathbf{D}}^2].\end{aligned}$$

Following the steps from the first part, we conclude the proof.

A.2 Proof of Proposition 1

We first present the following lemma about the convexity of a specific function.

Lemma 1. For every matrix $\mathbf{R} \in \mathbb{S}_{++}^d$, we define

$$f(\mathbf{D}) = \text{tr}(\mathbf{L}^{\frac{1}{2}} \mathbf{D} \mathbf{R} \mathbf{D} \mathbf{L}^{\frac{1}{2}}), \quad (25)$$

where $\mathbf{L}, \mathbf{D} \in \mathbb{S}_{++}^d$. Then function $f : \mathbb{S}_{++}^d \rightarrow \mathbb{R}$ is a convex function.

Proof. Let us rewrite (7) using quadratic forms. That is for every non-zero $v \in \mathbb{R}^d$, the following inequality must be true:

$$v^\top \mathbb{E}[\mathbf{S}^k \mathbf{D} \mathbf{L} \mathbf{D} \mathbf{S}^k] v \leq v^\top \mathbf{D} v, \quad \forall v \neq 0$$

Notice that both sides of this inequality are real numbers, thus can be written equivalently as

$$\text{tr}(v^\top \mathbb{E}[\mathbf{S}^k \mathbf{D} \mathbf{L} \mathbf{D} \mathbf{S}^k] v) \leq \text{tr}(v^\top \mathbf{D} v), \quad \forall v \neq 0$$

The LHS can be modified in the following way

$$\begin{aligned}\text{tr}(v^\top \mathbb{E}[\mathbf{S}^k \mathbf{D} \mathbf{L} \mathbf{D} \mathbf{S}^k] v) &\stackrel{\text{I}}{=} \text{tr}(\mathbb{E}[v^\top \mathbf{S}^k \mathbf{D} \mathbf{L} \mathbf{D} \mathbf{S}^k v]) \\ &\stackrel{\text{II}}{=} \mathbb{E}[\text{tr}(v^\top \mathbf{S}^k \mathbf{D} \mathbf{L} \mathbf{D} \mathbf{S}^k v)] \\ &\stackrel{\text{III}}{=} \mathbb{E}\left[\text{tr}(\mathbf{L}^{\frac{1}{2}} \mathbf{D} \mathbf{S}^k v v^\top \mathbf{S}^k \mathbf{D} \mathbf{L}^{\frac{1}{2}})\right] \\ &\stackrel{\text{IV}}{=} \text{tr}\left(\mathbb{E}\left[\mathbf{L}^{\frac{1}{2}} \mathbf{D} \mathbf{S}^k v v^\top \mathbf{S}^k \mathbf{D} \mathbf{L}^{\frac{1}{2}}\right]\right) \\ &\stackrel{\text{V}}{=} \text{tr}\left(\mathbf{L}^{\frac{1}{2}} \mathbf{D} \mathbb{E}[\mathbf{S}^k v v^\top \mathbf{S}^k] \mathbf{D} \mathbf{L}^{\frac{1}{2}}\right),\end{aligned}$$

where I, V are due to the linearity of expectation, II, IV are due to the linearity of trace operator, III is obtained using cyclic property of trace. Therefore, we can write the condition (7) equivalently as

$$\text{tr}\left(\mathbf{L}^{\frac{1}{2}} \mathbf{D} \mathbb{E}[\mathbf{S}^k v v^\top \mathbf{S}^k] \mathbf{D} \mathbf{L}^{\frac{1}{2}}\right) \leq \text{tr}(v v^\top \mathbf{D}), \quad \forall v \neq 0.$$

We then define function $g_v : \mathbb{S}_{++}^d \rightarrow \mathbb{R}$ for some fixed $v \neq 0$ as

$$g_v(\mathbf{D}) := \text{tr}\left(\mathbf{L}^{\frac{1}{2}} \mathbf{D} \mathbb{E}[\mathbf{S}^k v v^\top \mathbf{S}^k] \mathbf{D} \mathbf{L}^{\frac{1}{2}}\right) - \text{tr}(v v^\top \mathbf{D}). \quad (26)$$

We want to show that for every fixed $v \neq 0$, g is a convex function w.r.t \mathbf{D} , so that in this case, the sub-level set $\{\mathbf{D} \in \mathbb{S}_{++}^d \mid g_v(\mathbf{D}) \leq 0\}$ is convex.

- Notice that $v v^\top$ is a rank-1 matrix whose eigenvalues are all zero except one of them is $\|v\|^2 > 0$. We also have $(v v^\top)^\top = (v^\top)^\top v^\top = v v^\top$, so it is also a symmetric matrix. Thus we conclude that $v v^\top \in \mathbb{S}_+^d$ for every choice of v , we use $\mathbf{V} = v v^\top$ to denote it.
- If $\mathbf{S}^k = \mathbf{O}_d$, then the first term is equal to \mathbf{O}_d and the function $g_v(\mathbf{D})$ is linear, thus, also convex. Now, let us assume \mathbf{S}^k is nonzero. Similarly $\mathbf{S}^k v v^\top \mathbf{S}^k = \mathbf{S}^k v (\mathbf{S}^k v)^\top$ is also a symmetric positive semi-definite matrix whose eigenvalues are all 0 except one of them is $\|\mathbf{S}^k v\|^2$, this tells us that its expectation over \mathbf{S}^k is still a symmetric positive semi-definite matrix, we use $\mathbf{R} = \mathbb{E}[\mathbf{S}^k v v^\top \mathbf{S}^k]$ to denote it.

509 Now we can write function g_v as

$$g_v(\mathbf{D}) = \text{tr}(\mathbf{L}^{\frac{1}{2}} \mathbf{D} \mathbf{R} \mathbf{D} \mathbf{L}^{\frac{1}{2}}) - \text{tr}(\mathbf{V} \mathbf{D}).$$

510 According to Lemma 1, the first term of $g_v(\mathbf{D})$ is a convex function, and we know that the second
 511 term is linear in \mathbf{D} . As a result, $g_v(\mathbf{D})$ is a convex function w.r.t \mathbf{D} for every $v \neq 0$, thus the
 512 sub-level set $\{\mathbf{D} \in \mathbb{S}_{++}^d \mid g_v(\mathbf{D}) \leq 0\}$ is a convex set for every $v \neq 0$. The intersection of all those
 513 convex sets corresponding to every $v \neq 0$ is still a convex set, which tells us the original condition
 514 (7) is convex. \square

515 B Layer-wise case

516 Here in this section, we provide interpretations about some of the results and conclusions we had in
 517 Section 4.

518 B.1 Proof of Theorem 2

519 Note that $\mathbb{E}[\mathbf{S}^k \mathbf{D} \mathbf{L} \mathbf{D} \mathbf{S}^k] = \text{Diag}(\gamma_1^2 \mathbb{E}[\mathbf{S}_1^k \mathbf{W}_1 \mathbf{L}_1 \mathbf{W}_1 \mathbf{S}_1^k], \dots, \gamma_\ell^2 \mathbb{E}[\mathbf{S}_\ell^k \mathbf{W}_\ell \mathbf{L}_\ell \mathbf{D}_\ell \mathbf{S}_\ell^k])$, i.e.,

$$\mathbb{E}[\mathbf{S}^k \mathbf{D} \mathbf{L} \mathbf{D} \mathbf{S}^k] = \begin{pmatrix} \gamma_1^2 \mathbb{E}[\mathbf{S}_1^k \mathbf{W}_1 \mathbf{L}_1 \mathbf{W}_1 \mathbf{S}_1^k] & 0 & \cdots & 0 \\ 0 & \gamma_2^2 \mathbb{E}[\mathbf{S}_2^k \mathbf{W}_2 \mathbf{L}_2 \mathbf{W}_2 \mathbf{S}_2^k] & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & \gamma_\ell^2 \mathbb{E}[\mathbf{S}_\ell^k \mathbf{W}_\ell \mathbf{L}_\ell \mathbf{D}_\ell \mathbf{S}_\ell^k] \end{pmatrix},$$

520 which means that (7) holds if and only if $\gamma_i^2 \mathbb{E}[\mathbf{S}_i^k \mathbf{W}_i \mathbf{L}_i \mathbf{W}_i \mathbf{S}_i^k] \preceq \gamma_i \mathbf{W}_i$ for all $i \in [\ell]$, which holds
 521 if and only if (15) holds. So, Theorem 1 applies, and we conclude that

$$\frac{1}{K} \sum_{k=0}^{K-1} \mathbb{E}[\|\nabla f(x^k)\|_{\Gamma \mathbf{W}}^2] \leq \frac{2(f(x^0) - f^{\inf})}{K}, \quad (27)$$

522 To obtain (16), it remains to multiply both sides of (27) by $\frac{1}{\det(\Gamma \mathbf{W})^{1/d}}$.

523 B.2 Bernoulli sketch for det-CGD2

524 We formulate the following corollary regarding the communication complexity of det-CGD2 with
 525 Bernoulli sketches in the block diagonal setting.

526 **Corollary 2.** *If we are using Bern- q_i sketch \mathbf{T}_i^k for the i -th layer in det-CGD2 which is defined as*

$$\mathbf{T}_i^k = \frac{\eta_i}{q_i} \mathbf{I}_{d_i}, \quad \text{where } \eta_i \sim \text{Bernoulli}(q_i). \quad (28)$$

527 *Then in this case, the communication complexity of the algorithm is given by, if we leave out the*
 528 *constant factor $2(f(x^0) - f^{\inf})/\epsilon^2$,*

$$\left(\sum_{i=1}^{\ell} q_i d_i \right) \prod_{i=1}^{\ell} \left(\frac{1}{q_i} \right)^{\frac{d_i}{d}} \det(\mathbf{L})^{\frac{1}{d}}. \quad (29)$$

529 *In addition, the communication complexity is minimized if the probabilities $\{q_i\}_{i=1}^{\ell}$ satisfy*

$$q_i = q, \quad \forall i \in [\ell],$$

530 *The minimized communication complexity is*

$$d \cdot \det(\mathbf{L})^{\frac{1}{d}}. \quad (30)$$

531 *Proof.* For det-CGD2, its convergence requires (11). We are using Bernoulli sketch here, so we
 532 deduce that

$$\begin{aligned} \mathbb{E}[\mathbf{T}^k \mathbf{L} \mathbf{T}^k] &= \mathbb{E}[\text{Diag}(\mathbf{T}_1^k \mathbf{L}_1 \mathbf{T}_1^k, \dots, \mathbf{T}_\ell^k \mathbf{L}_\ell \mathbf{T}_\ell^k)] \\ &= \text{Diag}(\mathbb{E}[\mathbf{T}_1^k \mathbf{L}_1 \mathbf{T}_1^k], \dots, \mathbb{E}[\mathbf{T}_\ell^k \mathbf{L}_\ell \mathbf{T}_\ell^k]). \end{aligned}$$

533 Using the fact that for each block, we have

$$\mathbb{E} [\mathbf{T}_i^k \mathbf{L}_i \mathbf{T}_i^k] = (1 - q_i) \mathbf{O}_{d_i} \mathbf{L}_i \mathbf{O}_{d_i} + q_i \cdot \frac{1}{q_i^2} \mathbf{I}_{d_i} \mathbf{L}_i \mathbf{I}_{d_i} = \frac{\mathbf{L}_i}{q_i},$$

534 we obtain

$$\mathbb{E} [\mathbf{T}^k \mathbf{L} \mathbf{T}^k] = \text{Diag} \left(\frac{\mathbf{L}_1}{q_1}, \dots, \frac{\mathbf{L}_\ell}{q_\ell} \right).$$

535 Recalling (11), the best stepsize possible is therefore given by

$$\begin{aligned} \mathbf{D} &= (\mathbb{E} [\mathbf{T}^k \mathbf{L} \mathbf{T}^k])^{-1} \\ &= \text{Diag}^{-1} \left(\frac{\mathbf{L}_1}{q_1}, \dots, \frac{\mathbf{L}_\ell}{q_\ell} \right) \\ &= \text{Diag} (q_1 \mathbf{L}_1^{-1}, \dots, q_\ell \mathbf{L}_\ell^{-1}). \end{aligned}$$

536 From (10), we know that in order for **det-CGD2** to converge to ϵ^2 error level, we need

$$\frac{2(f(x^0) - f^{\inf})}{\det(\mathbf{D})^{\frac{1}{d}} K} \leq \epsilon^2,$$

537 which means that we need

$$K \geq \frac{2(f(x^0) - f^{\inf})}{\det(\mathbf{D})^{\frac{1}{d}} \epsilon^2} = \frac{1}{\det(\mathbf{D})^{\frac{1}{d}}} \cdot \frac{2(f(x^0) - f^{\inf})}{\epsilon^2},$$

538 iterations. For each iteration, the number of bits sent in expectation is equal to $\sum_{i=1}^{\ell} q_i d_i$. As a result,
539 the communication complexity is given by, if we leave out the constant factor $2(f(x^0) - f^{\inf})/\epsilon^2$,

$$\begin{aligned} \left(\sum_{i=1}^{\ell} q_i d_i \right) \cdot \frac{1}{\det(\mathbf{D})^{\frac{1}{d}}} &= \left(\sum_{i=1}^{\ell} q_i d_i \right) \cdot \det(\mathbf{D}^{-1})^{\frac{1}{d}} \\ &= \left(\sum_{i=1}^{\ell} q_i d_i \right) \cdot \left(\prod_{i=1}^{\ell} \det\left(\frac{\mathbf{L}_i}{q_i}\right) \right)^{\frac{1}{d}} \\ &= \left(\sum_{i=1}^{\ell} q_i d_i \right) \cdot \prod_{i=1}^{\ell} \left(\frac{1}{q_i} \right)^{\frac{d_i}{d}} \left(\prod_{i=1}^{\ell} \det(\mathbf{L}_i) \right)^{\frac{1}{d}} \\ &= \left(\sum_{i=1}^{\ell} q_i d_i \right) \cdot \prod_{i=1}^{\ell} \left(\frac{1}{q_i} \right)^{\frac{d_i}{d}} \cdot \det(\mathbf{L})^{\frac{1}{d}}. \end{aligned}$$

540 To obtain the optimal probability q_i , we can do the following transformation

$$\left(\sum_{i=1}^{\ell} q_i d_i \right) \cdot \frac{1}{\det(\mathbf{D})^{\frac{1}{d}}} = \left(\sum_{i=1}^{\ell} q_i \frac{d_i}{d} \right) \cdot \prod_{i=1}^{\ell} \left(\frac{1}{q_i} \right)^{\frac{d_i}{d}} \cdot d \det(\mathbf{L})^{\frac{1}{d}},$$

541 therefore, it is equivalent to minimize the coefficient

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

542 If we denote $\alpha_i = \frac{d_i}{d}$, then we know that $\alpha_i \in (0, 1]$ and $\sum_{i=1}^{\ell} \alpha_i = 1$, the above coefficient turns
543 into

$$\left(\sum_{i=1}^{\ell} \alpha_i q_i \right) \prod_{i=1}^{\ell} \left(\frac{1}{q_i} \right)^{\alpha_i}.$$

544 From the weighted AM-GM inequality (or the strict concavity of $\log(\cdot)$), we have

$$\left(\sum_{i=1}^{\ell} \alpha_i q_i \right) \geq \prod_{i=1}^{\ell} q_i^{\alpha_i},$$

545 where identity is obtained if and only if $q_i = q_j$, for all $i \neq j$. Thus we get

$$\left(\sum_{i=1}^{\ell} \alpha_i q_i \right) \prod_{i=1}^{\ell} \left(\frac{1}{q_i} \right)^{\alpha_i} \geq 1,$$

546 which in its turn implies that the minimum of expected communication complexity is equal to
 547 $d \cdot \det(\mathbf{L})^{\frac{1}{d}}$. The equality is achieved when the probabilities are equal. This concludes the proof. \square

548 B.3 General cases for **det-CGD1**

549 The first part (row 1 to row 8) of Table 1 records the communication complexities of **det-CGD1** in the
 550 block diagonal setting and in the general setting. Depending on the types of sketches \mathbf{S}_i^k and matrices
 551 \mathbf{W}_i we are using, we can calculate the optimal scaling factor γ_i using Theorem 2. According to (10),
 552 in order to reach an error level of ϵ^2 , we need

$$K \geq \frac{1}{\det(\mathbf{D})^{\frac{1}{d}}} \cdot \frac{2(f(x^0) - f^{\inf})}{\epsilon^2}, \quad (31)$$

553 where K is the number of iterations in total. We can then obtain the communication complexity
 554 taking into account the number of bits transferred in each iteration in the block diagonal case, the
 555 same apply to the general case which can be viewed as a special case of block diagonal setting where
 556 there is only 1 block.

557 B.4 General cases for **det-CGD2**

558 The second part of Table 1 (row 9 to row 12) records the communication complexities of **det-CGD2**.
 559 Different from **det-CGD1**, we can always obtain the best stepsize matrix \mathbf{D} here if the sketch \mathbf{S}^k is
 560 given. The communication complexity can then be obtained in the same way as previous case using
 561 (31) in combination with the number of bits sent per iteration.

562 B.5 Interpretations of Table 1

563 Compared to GD (row 13), **det-CGD1** and **det-CGD2** using matrix stepsize without compression (row
 564 1, row 9) is better in terms of both iteration complexity and communication complexity. By utilizing
 565 the block diagonal structure, we are able to design special sketches that allow us to compress for free.
 566 This can be seen from row 12, where the communication complexity of using Bernoulli compressor
 567 with equal probabilities for **det-CGD2** in expectation is the same with GD, but the number of bits
 568 sent per iteration is reduced. There are some results in the table needs careful analysis, especially for
 569 **det-CGD1**.

570 B.5.1 Comparison of row 5 and 7

571 Here we show that the communication complexity given in row 5 is always worse than that of row 7.
 572 This can be seen from the following corollary.

573 **Corollary 3.** *For any matrix $\mathbf{L} \in \mathbb{S}_{++}^d$, the following inequality holds*

$$\lambda_{\max} \left(\mathbf{L}^{\frac{1}{2}} \text{diag}(\mathbf{L}^{-1}) \mathbf{L}^{\frac{1}{2}} \right) \cdot \det(\mathbf{L})^{\frac{1}{d}} \geq \det(\text{diag}(\mathbf{L}))^{\frac{1}{d}}.$$

574 *Proof.* The inequality given in Corollary 3 can be reformulated equivalently to

$$\lambda_{\max}(\mathbf{L} \text{diag}(\mathbf{L}^{-1})) \geq \det(\mathbf{L}^{-1} \text{diag}(\mathbf{L}))^{\frac{1}{d}}.$$

575 We use the notation

$$\mathbf{M}_1 = \mathbf{L} \text{diag}(\mathbf{L}^{-1}), \quad \mathbf{M}_2 = \mathbf{L}^{-1} \text{diag}(\mathbf{L}),$$

576 and notice that for any $i \in [d]$, we have

$$(\mathbf{M}_1)_{ii} = (\mathbf{L})_{ii} \cdot (\mathbf{L}^{-1})_{ii} = (\mathbf{M}_2)_{ii}.$$

577 As a result

$$\lambda_{\max}(\mathbf{M}_1) \geq \left(\prod_{i=1}^d (\mathbf{M}_1)_{ii} \right)^{\frac{1}{d}} = \left(\prod_{i=1}^d (\mathbf{M}_2)_{ii} \right)^{\frac{1}{d}} \geq \det(\mathbf{M}_2)^{\frac{1}{d}},$$

578 where the first inequality is due to the fact that each diagonal element is upper-bounded by the
 579 maximum eigenvalue value, and the last is obtained using the fact that the product of the diagonal
 580 elements is an upper bound of the determinant. \square

581 From Corollary 3, it immediately follows that the result in row 7 is better than row 5 in terms of both
 582 communication and iteration complexity.

583 B.5.2 Comparison of row 6 and 7

584 In this section we want to give a simple example that tells us results in row 6 and 7 are not comparable
 585 in general. Consider a simple matrix $\mathbf{L} \in \mathbb{S}_{++}^2$, if we pick

$$\mathbf{L} = \begin{pmatrix} 16 & 0 \\ 0 & 1 \end{pmatrix},$$

586 then

$$\begin{aligned} \det(\text{diag}(\mathbf{L}))^{\frac{1}{d}} &= 4; \\ \lambda_{\max}^{\frac{1}{2}}(\mathbf{L}) \det(\mathbf{L})^{\frac{1}{2d}} &= 8. \end{aligned}$$

587 However, if we pick

$$\mathbf{L} = \begin{pmatrix} 16 & 3.9 \\ 3.9 & 1 \end{pmatrix},$$

588 then

$$\begin{aligned} \det(\text{diag}(\mathbf{L}))^{\frac{1}{d}} &= 4; \\ \lambda_{\max}^{\frac{1}{2}}(\mathbf{L}) \det(\mathbf{L})^{\frac{1}{2d}} &\simeq 3.88. \end{aligned}$$

589 From this example, we can see that the relation between the results in row 6 and 7 may vary depending
 590 on the value of \mathbf{L} .

591 C Distributed case

592 C.1 Proof of Theorem 3

593 We first present following technical lemmas.

594 **Lemma 2.** For any sketch \mathbf{S}_i^k of client i drawn randomly from some distribution \mathcal{S} over \mathbb{S}_+^d which
 595 satisfies

$$\mathbb{E}[\mathbf{S}_i^k] = \mathbf{I}_d$$

596 the following bound holds for any $x \in \mathbb{R}^d$ and each client i ,

$$\mathbb{E} \left[\|\mathbf{S}_i^k x - x\|_{\mathbf{DLD}}^2 \right] \leq \lambda_{\max} \left(\mathbf{L}_i^{\frac{1}{2}} \mathbb{E} [(\mathbf{S}_i^k - \mathbf{I}_d) \mathbf{DLD} (\mathbf{S}_i^k - \mathbf{I}_d)] \mathbf{L}_i^{\frac{1}{2}} \right) \cdot \|x\|_{\mathbf{L}_i^{-1}}^2. \quad (32)$$

597 Here $\mathbf{D} \in \mathbb{S}_{++}^d$ is the stepsize matrix, $\mathbf{L}, \mathbf{L}_i \in \mathbb{S}_{++}^d$ are the smoothness matrices for f and f_i ,
 598 respectively.

599 **Lemma 3.** (Variance Decomposition) For any random vector $x \in \mathbb{R}^d$, and any matrix $\mathbf{M} \in \mathbb{S}_+^d$, the
 600 following identity holds

$$\mathbb{E} \left[\|x - \mathbb{E}[x]\|_{\mathbf{M}}^2 \right] = \mathbb{E} \left[\|x\|_{\mathbf{M}}^2 \right] - \|\mathbb{E}[x]\|_{\mathbf{M}}^2. \quad (33)$$

601 **Lemma 4.** Assume $\{a_i\}_{i=1}^n$ is a set of independent random vectors in \mathbb{R}^d , which satisfy

$$\mathbb{E}[a_i] = 0, \quad \forall i \in [n].$$

602 Then, for any $M \in \mathbb{S}_{++}^d$, we have

$$\mathbb{E} \left[\left\| \frac{1}{n} \sum_{i=1}^n a_i \right\|_M^2 \right] = \frac{1}{n^2} \sum_{i=1}^n \mathbb{E} [\|a_i\|_M^2]. \quad (34)$$

603 **Lemma 5.** (Property of Sketch) For any vector $x \in \mathbb{R}^d$, and sketch matrix $S \in \mathbb{S}_+^d$ taken from some
604 distribution \mathcal{S} over \mathbb{S}_+^d , which satisfies

$$\mathbb{E}[S] = I_d.$$

605 Then for any matrix $M \in \mathbb{S}_{++}^d$, we have the following identity holds,

$$\mathbb{E} [\|Sx - x\|_M^2] = \|x\|_{\mathbb{E}[SM S] - M}^2. \quad (35)$$

606 **Lemma 6.** If we have a differentiable function $f : \mathbb{R}^d \rightarrow \mathbb{R}$, that is L matrix smooth and lower
607 bounded by f^{\inf} , if we assume $L \in \mathbb{S}_{++}^d$, then the following inequality holds

$$\langle \nabla f(x), L^{-1} \nabla f(x) \rangle \leq 2(f(x) - f^{\inf}). \quad (36)$$

608 We start by defining

$$g(x) := \frac{1}{n} \sum_{i=1}^n S_i^k \nabla f_i(x), \quad (37)$$

609 as a result, **det-CGD1** in the distributed case can then be written as

$$x^{k+1} = x^k - Dg(x^k).$$

610 Notice that we have

$$\mathbb{E} [g(x^k) \mid x^k] = \frac{1}{n} \sum_{i=1}^n \mathbb{E} [S_i^k] \nabla f_i(x^k) = \nabla f(x^k). \quad (38)$$

611 We start with L matrix smoothness of f .

$$\begin{aligned} f(x^{k+1}) &\leq f(x^k) + \langle \nabla f(x^k), x^{k+1} - x^k \rangle + \frac{1}{2} \langle L(x^{k+1} - x^k), x^{k+1} - x^k \rangle \\ &= f(x^k) + \langle \nabla f(x^k), -Dg(x^k) \rangle + \frac{1}{2} \langle L(-Dg(x^k)), -Dg(x^k) \rangle \\ &= f(x^k) - \langle \nabla f(x^k), Dg(x^k) \rangle + \frac{1}{2} \langle LDg(x^k), Dg(x^k) \rangle. \end{aligned}$$

612 Taking expectation conditioned on x^k , we get

$$\begin{aligned} &\mathbb{E} [f(x^{k+1}) \mid x^k] \\ &\leq f(x^k) - \langle \nabla f(x^k), D\mathbb{E} [g(x^k) \mid x^k] \rangle + \frac{1}{2} \mathbb{E} [\langle LDg(x^k), Dg(x^k) \rangle \mid x^k] \\ &\stackrel{(38)}{=} f(x^k) - \langle \nabla f(x^k), D\nabla f(x^k) \rangle + \frac{1}{2} \mathbb{E} [\langle LDg(x^k), Dg(x^k) \rangle \mid x^k] \\ &= f(x^k) - \|\nabla f(x^k)\|_D^2 + \frac{1}{2} \underbrace{\mathbb{E} [\langle LDg(x^k), Dg(x^k) \rangle \mid x^k]}_{:=T}. \end{aligned} \quad (39)$$

613 The last term T of above can be bounded by

$$\begin{aligned} T &= \mathbb{E} [\|g(x^k)\|_{DL D}^2 \mid x^k] \\ &\stackrel{(33)}{=} \mathbb{E} [\|g(x^k) - \mathbb{E} [g(x^k) \mid x^k]\|_{DL D}^2 \mid x^k] + \|\mathbb{E} [g(x^k) \mid x^k]\|_{DL D}^2. \end{aligned}$$

614 We have already shown that $\mathbb{E} [g(x^k) \mid x^k] = \nabla f(x^k)$,

$$\begin{aligned} T &= \mathbb{E} \left[\left\| g(x^k) - \nabla f(x^k) \right\|_{DL D}^2 \mid x^k \right] + \left\| \nabla f(x^k) \right\|_{DL D}^2 \\ &= \mathbb{E} \left[\left\| \frac{1}{n} \sum_{i=1}^n (\mathbf{S}_i^k \nabla f_i(x^k) - \nabla f_i(x^k)) \right\|_{DL D}^2 \mid x^k \right] + \left\| \nabla f(x^k) \right\|_{DL D}^2. \end{aligned}$$

615 Using Lemma 4, we have

$$\begin{aligned} T &= \frac{1}{n^2} \sum_{i=1}^n \mathbb{E} \left[\left\| \mathbf{S}_i^k \nabla f_i(x^k) - \nabla f_i(x^k) \right\|_{DL D}^2 \mid x^k \right] + \left\| \nabla f(x^k) \right\|_{DL D}^2 \\ &\leq \frac{1}{n^2} \sum_{i=1}^n \mathbb{E} \left[\left\| \mathbf{S}_i^k \nabla f_i(x^k) - \nabla f_i(x^k) \right\|_{DL D}^2 \mid x^k \right] + \left\| \nabla f(x^k) \right\|_D^2, \end{aligned} \quad (40)$$

616 where the last inequality holds because of the inequality $DL D \preceq D$. Plug (40) into (39), we get

$$\begin{aligned} &\mathbb{E} [f(x^{k+1}) \mid x^k] \\ &\leq f(x^k) - \frac{1}{2} \left\| \nabla f(x^k) \right\|_D^2 + \frac{1}{2n^2} \sum_{i=1}^n \mathbb{E} \left[\left\| \mathbf{S}_i^k \nabla f_i(x^k) - \nabla f_i(x^k) \right\|_{DL D}^2 \mid x^k \right]. \\ &\stackrel{(32)}{\leq} f(x^k) - \frac{1}{2} \left\| \nabla f(x^k) \right\|_D^2 \\ &\quad + \frac{1}{2n^2} \sum_{i=1}^n \lambda_{\max} \left(\mathbb{E} \left[\mathbf{L}_i^{\frac{1}{2}} (\mathbf{S}_i^k - \mathbf{I}_d) DL D (\mathbf{S}_i^k - \mathbf{I}_d) \mathbf{L}_i^{\frac{1}{2}} \right] \right) \left\| \nabla f_i(x^k) \right\|_{\mathbf{L}_i^{-1}}^2 \\ &\stackrel{(36)}{\leq} f(x^k) - \frac{1}{2} \left\| \nabla f(x^k) \right\|_D^2 \\ &\quad + \frac{1}{n^2} \sum_{i=1}^n \lambda_{\max} \left(\mathbb{E} \left[\mathbf{L}_i^{\frac{1}{2}} (\mathbf{S}_i^k - \mathbf{I}_d) DL D (\mathbf{S}_i^k - \mathbf{I}_d) \mathbf{L}_i^{\frac{1}{2}} \right] \right) (f_i(x^k) - f_i^{\inf}). \end{aligned}$$

617 Let

$$\lambda_D = \max_i \left\{ \lambda_{\max} \left(\mathbb{E} \left[\mathbf{L}_i^{\frac{1}{2}} (\mathbf{S}_i^k - \mathbf{I}_d) DL D (\mathbf{S}_i^k - \mathbf{I}_d) \mathbf{L}_i^{\frac{1}{2}} \right] \right) \right\}, \quad (41)$$

618 the above bound on $f(x^{k+1})$ turns into

$$\begin{aligned} &\mathbb{E} [f(x^{k+1}) \mid x^k] \\ &\leq f(x^k) - \frac{1}{2} \left\| \nabla f(x^k) \right\|_D^2 + \frac{1}{n^2} \sum_{i=1}^n \lambda_D (f_i(x^k) - f_i^{\inf}) \\ &= f(x^k) - \frac{1}{2} \left\| \nabla f(x^k) \right\|_D^2 + \frac{\lambda_D}{n} \left(\frac{1}{n} \sum_{i=1}^n f_i(x^k) - \frac{1}{n} \sum_{i=1}^n f_i^{\inf} \right) \\ &= f(x^k) - \frac{1}{2} \left\| \nabla f(x^k) \right\|_D^2 + \frac{\lambda_D}{n} (f(x^k) - f^{\inf}) + \frac{\lambda_D}{n} \left(f^{\inf} - \frac{1}{n} \sum_{i=1}^n f_i^{\inf} \right). \end{aligned}$$

619 Subtracting f^{\inf} from both sides, we get

$$\begin{aligned} &\mathbb{E} [f(x^{k+1}) - f^{\inf} \mid x^k] \\ &\leq f(x^k) - f^{\inf} - \frac{1}{2} \left\| \nabla f(x^k) \right\|_D^2 + \frac{\lambda_D}{n} (f(x^k) - f^{\inf}) + \frac{\lambda_D}{n} \left(f^{\inf} - \frac{1}{n} \sum_{i=1}^n f_i^{\inf} \right). \end{aligned}$$

620 Taking expectation, applying tower property and rearranging terms, we get

$$\begin{aligned} &\mathbb{E} [f(x^{k+1}) - f^{\inf}] \\ &\leq \left(1 + \frac{\lambda_D}{n} \right) \mathbb{E} [f(x^k) - f^{\inf}] - \frac{1}{2} \mathbb{E} [\left\| \nabla f(x^k) \right\|_D^2] + \frac{\lambda_D}{n} \left(f^{\inf} - \frac{1}{n} \sum_{i=1}^n f_i^{\inf} \right). \end{aligned} \quad (42)$$

621 If we denote

$$\delta^k = \mathbb{E} [f(x^k) - f^{\text{inf}}], \quad r^k = \mathbb{E} [\|\nabla f(x^k)\|_D^2], \quad \Delta^{\text{inf}} = f^{\text{inf}} - \frac{1}{n} \sum_{i=1}^n f_i^{\text{inf}}.$$

622 Then (42) becomes

$$\frac{1}{2} r^k \leq \left(1 + \frac{\lambda_D}{n}\right) \delta^k - \delta^{k+1} + \frac{\lambda_D \Delta^{\text{inf}}}{n} \quad (43)$$

623 In order to approach the final result, we now follow [Sti19], [KR20] and define an exponentially
624 decaying weighting sequence $\{w_k\}_{k=-1}^K$, where K is the total number of iterations. We fix $w_{-1} > 0$
625 and define

$$w_k = \frac{w_{k-1}}{1 + \lambda_D/n}, \quad \text{for all } k \geq 0.$$

626 By multiplying both sides of the recursion (43) by w_k , we get

$$\frac{1}{2} w_k r^k \leq w_{k-1} \delta^k - w_k \delta^{k+1} + \frac{\lambda_D \Delta^{\text{inf}}}{n} w_k$$

627 Summing up the inequalities from $k = 0, \dots, K-1$, we get

$$\frac{1}{2} \sum_{k=0}^{K-1} w_k r^k \leq w_{-1} \delta^0 - w_{K-1} \delta^K + \frac{\lambda_D \Delta^{\text{inf}}}{n} \sum_{k=0}^{K-1} w_k.$$

628 Define $W_K = \sum_{k=0}^{K-1} w_k$, and divide both sides by W_K , we get

$$\frac{1}{2} \min_{0 \leq k \leq K-1} r^k \leq \frac{1}{2} \frac{\sum_{k=0}^{K-1} w_k r^k}{W_K} \leq \frac{w_{-1}}{W_K} \delta^0 + \frac{\lambda_D \Delta^{\text{inf}}}{n},$$

629 Notice that from the definition of w_k , we know that the following inequality holds,

$$\frac{w_{-1}}{W_K} \leq \frac{w_{-1}}{K w_{K-1}} = \frac{(1 + \frac{\lambda_D}{n})^K}{K}.$$

630 As a result, we have

$$\min_{0 \leq k \leq K-1} r^k \leq \frac{2(1 + \frac{\lambda_D}{n})^K}{K} \delta^0 + \frac{2\lambda_D \Delta^{\text{inf}}}{n}.$$

631 Recalling the definition for r^k and δ^k , we get the following result,

$$\min_{0 \leq k \leq K-1} \mathbb{E} [\|\nabla f(x^k)\|_D^2] \leq \frac{2(1 + \frac{\lambda_D}{n})^K (f(x^0) - f^{\text{inf}})}{K} + \frac{2\lambda_D \Delta^{\text{inf}}}{n}.$$

632 Then we do determinant normalization and get,

$$\min_{0 \leq k \leq K-1} \mathbb{E} [\|\nabla f(x^k)\|_{D/\det(D)^{1/d}}^2] \leq \frac{2(1 + \frac{\lambda_D}{n})^K (f(x^0) - f^{\text{inf}})}{\det(D)^{1/d} K} + \frac{2\lambda_D \Delta^{\text{inf}}}{\det(D)^{1/d} n}. \quad (44)$$

633 This concludes the proof.

634 C.2 Convexity of the constraints

635 **Proposition 2.** *The set of matrices D that satisfy (21) is convex.*

636 *Proof.* The first inequality in (21) can be reformulated into

$$D \preceq L^{-1},$$

637 which is linear in D thus convex. For the second constraint in (21),

$$\max_i \left\{ \lambda_{\max} \left(\mathbb{E} \left[L_i^{\frac{1}{2}} (S_i^k - I_d) D L D (S_i^k - I_d) L_i^{\frac{1}{2}} \right] \right) \right\} \leq \frac{n}{K}, \quad (45)$$

638 we can reformulate it into n constraints, one for each client i ,

$$\begin{aligned}
& \lambda_{\max} \left(\mathbb{E} \left[\mathbf{L}_i^{\frac{1}{2}} (\mathbf{S}_i^k - \mathbf{I}_d) \mathbf{D} \mathbf{L} \mathbf{D} (\mathbf{S}_i^k - \mathbf{I}_d) \mathbf{L}_i^{\frac{1}{2}} \right] \right) \leq \frac{n}{K}, \quad \forall i \\
\Leftrightarrow & \mathbb{E} \left[\mathbf{L}_i^{\frac{1}{2}} (\mathbf{S}_i^k - \mathbf{I}_d) \mathbf{D} \mathbf{L} \mathbf{D} (\mathbf{S}_i^k - \mathbf{I}_d) \mathbf{L}_i^{\frac{1}{2}} \right] \preceq \frac{n}{K} \mathbf{I}_d, \quad \forall i \\
\Leftrightarrow & \mathbf{L}_i^{\frac{1}{2}} \mathbb{E} \left[(\mathbf{S}_i^k - \mathbf{I}_d) \mathbf{D} \mathbf{L} \mathbf{D} (\mathbf{S}_i^k - \mathbf{I}_d) \right] \mathbf{L}_i^{\frac{1}{2}} \preceq \frac{n}{K} \mathbf{I}_d, \quad \forall i \\
\Leftrightarrow & \mathbb{E} \left[(\mathbf{S}_i^k - \mathbf{I}_d) \mathbf{D} \mathbf{L} \mathbf{D} (\mathbf{S}_i^k - \mathbf{I}_d) \right] \preceq \frac{n}{K} \mathbf{L}_i^{-1}, \quad \forall i.
\end{aligned}$$

639 We then look at the individual condition for one client i ,

$$\mathbb{E} \left[(\mathbf{S}_i^k - \mathbf{I}_d) \mathbf{D} \mathbf{L} \mathbf{D} (\mathbf{S}_i^k - \mathbf{I}_d) \right] \preceq \frac{n}{K} \mathbf{L}_i^{-1}, \quad (46)$$

640 that is for any vector $0 \neq u \in \mathbb{R}^d$, we require

$$\begin{aligned}
& u^\top \mathbb{E} \left[(\mathbf{S}_i^k - \mathbf{I}_d) \mathbf{D} \mathbf{L} \mathbf{D} (\mathbf{S}_i^k - \mathbf{I}_d) \right] u \leq \frac{n}{K} u^\top \mathbf{L}_i^{-1} u, \quad \forall u \neq 0 \\
\Leftrightarrow & \text{tr} \left(u^\top \mathbb{E} \left[(\mathbf{S}_i^k - \mathbf{I}_d) \mathbf{D} \mathbf{L} \mathbf{D} (\mathbf{S}_i^k - \mathbf{I}_d) \right] u \right) \leq \frac{n}{K} \text{tr}(u^\top \mathbf{L}_i^{-1} u), \quad \forall u \neq 0 \\
\Leftrightarrow & \mathbb{E} \left[\text{tr}(u^\top (\mathbf{S}_i^k - \mathbf{I}_d) \mathbf{D} \mathbf{L} \mathbf{D} (\mathbf{S}_i^k - \mathbf{I}_d) u) \right] \leq \text{tr}(u^\top \mathbf{L}_i^{-1} u), \quad \forall u \neq 0 \\
\Leftrightarrow & \text{tr}(\mathbf{L}_i^{\frac{1}{2}} \mathbf{D} \mathbb{E} \left[(\mathbf{S}_i^k - \mathbf{I}_d) u u^\top (\mathbf{S}_i^k - \mathbf{I}_d) \right] \mathbf{D} \mathbf{L}_i^{\frac{1}{2}}) \leq \text{tr}(u^\top \mathbf{L}_i^{-1} u), \quad \forall u \neq 0.
\end{aligned}$$

641 We now define function $g_u : \mathbb{S}_{++}^d \rightarrow \mathbb{R}$ for every fixed $u \neq 0$,

$$g_u(\mathbf{D}) = \text{tr}(\mathbf{L}_i^{\frac{1}{2}} \mathbf{D} \mathbb{E} \left[(\mathbf{S}_i^k - \mathbf{I}_d) u u^\top (\mathbf{S}_i^k - \mathbf{I}_d) \right] \mathbf{D} \mathbf{L}_i^{\frac{1}{2}}), \quad (47)$$

642 notice that $u u^\top$ is a rank-1 matrix that is positive semi-definite, so for every $y \in \mathbb{R}^d$,

$$((\mathbf{S}_i^k - \mathbf{I}_d) y)^\top u u^\top ((\mathbf{S}_i^k - \mathbf{I}_d) y) \geq 0,$$

643 which means that $(\mathbf{S}_i^k - \mathbf{I}_d) u u^\top (\mathbf{S}_i^k - \mathbf{I}_d)$ is positive semi-definite, and thus is $\mathbf{R} :=$
644 $\mathbb{E} \left[(\mathbf{S}_i^k - \mathbf{I}_d) u u^\top (\mathbf{S}_i^k - \mathbf{I}_d) \right]$. Using Lemma 1, we know that $g_u(\mathbf{D})$ is a convex function for
645 every $0 \neq u \in \mathbb{R}^d$, thus its sub-level set $\{\mathbf{D} \in \mathbb{S}_{++}^d \mid g_u(\mathbf{D}) \leq \text{tr}(u^\top \mathbf{L}_i^{-1} u)\}$ is a convex set. The
646 intersection of those convex sets corresponding to the individual constraint (46) of client i is convex.
647 Again the intersection of those convex sets for each client i , which corresponds to (45), is still convex.

648 For the third constraint in (21), we can transform it using similar steps as we obtain (45) into

$$\mathbb{E} \left[(\mathbf{S}_i^k - \mathbf{I}_d) \mathbf{D} \mathbf{L} \mathbf{D} (\mathbf{S}_i^k - \mathbf{I}_d) \right] \preceq \frac{n\varepsilon^2}{4\Delta_{\inf}} \det(\mathbf{D})^{1/d} \mathbf{L}_i^{-1}, \quad \forall i. \quad (48)$$

649 If we look at each individual constraint, we can write in quadratic forms for any $0 \neq u \in \mathbb{R}^d$,

$$u^\top \mathbb{E} \left[(\mathbf{S}_i^k - \mathbf{I}_d) \mathbf{D} \mathbf{L} \mathbf{D} (\mathbf{S}_i^k - \mathbf{I}_d) \right] u \leq \frac{n\varepsilon^2}{4\Delta_{\inf}} \det(\mathbf{D})^{1/d} \cdot u^\top \mathbf{L}_i^{-1} u, \quad \forall u \neq 0.$$

650 Using the linearity of expectation and the trace operator with the trace trick, we can transform the
651 above condition into,

$$\text{tr}(\mathbf{L}_i^{\frac{1}{2}} \mathbf{D} \mathbb{E} \left[(\mathbf{S}_i^k - \mathbf{I}_d) u u^\top (\mathbf{S}_i^k - \mathbf{I}_d) \right] \mathbf{D} \mathbf{L}_i^{\frac{1}{2}}) \leq \frac{n\varepsilon^2}{4\Delta_{\inf}} \det(\mathbf{D})^{\frac{1}{d}} \text{tr}(u^\top \mathbf{L}_i^{-1} u) \quad \forall u \neq 0.$$

652 notice that we have already shown that $\mathbf{R} = \mathbb{E} \left[(\mathbf{S}_i^k - \mathbf{I}_d) u u^\top (\mathbf{S}_i^k - \mathbf{I}_d) \right] \in \mathbb{S}_+^d$, thus if we apply
653 Lemma 1, we know that the LHS is a convex function, while we know that $\det(\mathbf{D})^{\frac{1}{d}}$ is a concave
654 function for hermitian positive definite matrices \mathbf{D} . So the set of \mathbf{D} satisfying the constraint here
655 for every $u \in \mathbb{R}^d$ is convex, thus the intersection of them is convex, which means that the set of
656 \mathbf{D} satisfying the constraint for each client i is convex. Thus the intersection of those convex sets
657 corresponding to different clients, which corresponds to (48), is still convex. Now we know that the
658 set of \mathbf{D} satisfying each of the three constraints in (21) is convex, thus the intersection of them is
659 convex. This concludes the proof. \square

660 C.2.1 Proof of Corollary 1

661 For the first term in the RHS of convergence bound (20) under condition (21), we know that

$$2(1 + \frac{\lambda_D}{n})^K \leq 2 \cdot \exp(\lambda_D \cdot \frac{K}{n}) \leq 2 \cdot \exp(1) \leq 6,$$

662 thus

$$\begin{aligned} \frac{2(1 + \frac{\lambda_D}{n})^K (f(x^0) - f^{\inf})}{\det(\mathbf{D})^{1/d} K} &\leq \frac{6 (f(x^0) - f^{\inf})}{\det(\mathbf{D})^{1/d} K} \\ &\leq \frac{6 (f(x^0) - f^{\inf})}{\det(\mathbf{D})^{1/d}} \cdot \frac{\varepsilon^2 \det(\mathbf{D})^{\frac{1}{d}}}{12 (f(x^0) - f^{\inf})} \\ &= \frac{\varepsilon^2}{2}. \end{aligned}$$

663 While for the second term of RHS in (20), we have

$$\begin{aligned} \frac{2\lambda_D \Delta^{\inf}}{\det(\mathbf{D})^{1/d} n} &\leq \frac{2\Delta^{\inf}}{\det(\mathbf{D})^{1/d} n} \cdot \frac{\varepsilon^2 (\det(\mathbf{D}))^{1/d} n}{4\Delta^{\inf}} \\ &\leq \frac{\varepsilon^2}{2}. \end{aligned}$$

664 Thus we know that the left hand side of (20) is upper bounded by

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

665 This concludes the proof.

666 C.3 Distributed det-CGD2

667 We also extend det-CGD2 to the distributed case. Consider the method

$$x^{k+1} = x^k - \frac{1}{n} \sum_{i=1}^n \mathbf{T}_i^k \mathbf{D} \nabla f_i(x^k), \quad (49)$$

668 where $\mathbf{D} \in \mathbb{S}_{++}^d$ is the stepsize matrix, and each \mathbf{T}_i^k is a sequence of sketch matrices drawn randomly
669 from some distribution \mathcal{T} over \mathbb{S}_+^d independent of each other, satisfying

$$\mathbb{E} [\mathbf{T}_i^k] = \mathbf{I}_d. \quad (50)$$

670 C.3.1 Analysis of distributed det-CGD2

671 In this section, we present the theory for Algorithm 2, which is an analogue of Theorem 3. We first
672 present the following lemma which is necessary for our analysis.

673 **Lemma 7.** For any sketch \mathbf{T}_i^k of client i drawn randomly from some distribution \mathcal{T} over \mathbb{S}_+^d which
674 satisfies

$$\mathbb{E} [\mathbf{T}_i^k] = \mathbf{I}_d,$$

675 the following inequality holds for any $x \in \mathbb{R}^d$ for each client i ,

$$\mathbb{E} \left[\left\| \mathbf{T}_i^k \mathbf{D} x - \mathbf{D} x \right\|_{\mathbf{L}}^2 \right] \leq \lambda_{\max} \left(\mathbf{L}_i^{\frac{1}{2}} \mathbf{D} \mathbb{E} [(\mathbf{T}_i^k - \mathbf{I}_d) \mathbf{L} (\mathbf{T}_i^k - \mathbf{I}_d)] \mathbf{D} \mathbf{L}_i^{\frac{1}{2}} \right) \cdot \|x\|_{\mathbf{L}_i^{-1}}^2. \quad (51)$$

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

$$\mathbf{D} \mathbf{L} \mathbf{D} \preceq \mathbf{D}, \quad (52)$$

678 then the following convergence bound is true for the iteration of Algorithm 2

$$\min_{0 \leq k \leq K-1} \mathbb{E} \left[\left\| \nabla f(x^k) \right\|_{\frac{\mathbf{D}}{\det(\mathbf{D})^{1/d}}}^2 \right] \leq \frac{2(1 + \frac{\lambda'_D}{n})^K (f(x^0) - f^{\inf})}{\det(\mathbf{D})^{1/d} K} + \frac{2\lambda'_D \Delta^{\inf}}{\det(\mathbf{D})^{1/d} n}, \quad (53)$$

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

$$\lambda'_D := \max_i \left\{ \lambda_{\max} \left(\mathbb{E} \left[\mathbf{L}_i^{\frac{1}{2}} \mathbf{D} (\mathbf{T}_i^k - \mathbf{I}_d) \mathbf{L} (\mathbf{T}_i^k - \mathbf{I}_d) \mathbf{D} \mathbf{L}_i^{\frac{1}{2}} \right] \right) \right\}.$$

680 *Proof.* We first define function $g(x)$ as follows,

$$g(x) = \frac{1}{n} \sum_{i=1}^n \mathbf{T}_i^k \mathbf{D} \nabla f_i(x^k).$$

681 As a result, Algorithm 2 can be written as

$$x^{k+1} = x^k - g(x^k).$$

682 Notice that

$$\mathbb{E}[g(x)] = \frac{1}{n} \sum_{i=1}^n \mathbb{E}[\mathbf{T}_i^k] \mathbf{D} \nabla f_i(x) = \mathbf{D} \nabla f(x). \quad (54)$$

683 We then start with the \mathbf{L} matrix smoothness of function f ,

$$\begin{aligned} f(x^{k+1}) &\leq f(x^k) + \langle \nabla f(x^k), x^{k+1} - x^k \rangle + \frac{1}{2} \langle \mathbf{L}(x^{k+1} - x^k), x^{k+1} - x^k \rangle \\ &= f(x^k) + \langle \nabla f(x^k), -g(x^k) \rangle + \frac{1}{2} \langle \mathbf{L}(-g(x^k)), -g(x^k) \rangle \\ &= f(x^k) - \langle \nabla f(x^k), g(x^k) \rangle + \frac{1}{2} \langle \mathbf{L}g(x^k), g(x^k) \rangle. \end{aligned}$$

684 We then take expectation conditioned on x^k ,

$$\begin{aligned} &\mathbb{E}[f(x^{k+1}) | x^k] \\ &\leq f(x^k) - \langle \nabla f(x^k), \mathbb{E}[g(x^k) | x^k] \rangle + \frac{1}{2} \mathbb{E}[\langle \mathbf{L}g(x^k), g(x^k) \rangle | x^k] \\ &= f(x^k) - \langle \nabla f(x^k), \mathbf{D} \nabla f(x^k) \rangle + \frac{1}{2} \underbrace{\mathbb{E}[\langle \mathbf{L}g(x^k), g(x^k) \rangle | x^k]}_{:=T}. \end{aligned} \quad (55)$$

685 We then upper bound the last term T in the following way

$$\begin{aligned} T &= \mathbb{E}[\|g(x^k)\|_{\mathbf{L}}^2 | x^k] \\ &\stackrel{(33)}{=} \mathbb{E}[\|g(x^k) - \mathbb{E}[g(x^k) | x^k]\|_{\mathbf{L}}^2 | x^k] + \|\mathbb{E}[g(x^k) | x^k]\|_{\mathbf{L}}^2. \end{aligned}$$

686 From (54) we deduce

$$\begin{aligned} T &= \mathbb{E}[\|g(x^k) - \mathbf{D} \nabla f(x^k)\|_{\mathbf{L}}^2 | x^k] + \|\mathbf{D} \nabla f(x^k)\|_{\mathbf{L}}^2 \\ &= \mathbb{E} \left[\left\| \frac{1}{n} \sum_{i=1}^n \mathbf{T}_i^k \mathbf{D} \nabla f_i(x^k) - \frac{\mathbf{D}}{n} \sum_{i=1}^n \nabla f_i(x^k) \right\|_{\mathbf{L}}^2 | x^k \right] + \|\nabla f(x^k)\|_{\mathbf{D} \mathbf{L} \mathbf{D}}^2 \\ &= \mathbb{E} \left[\left\| \frac{1}{n} \sum_{i=1}^n (\mathbf{T}_i^k \mathbf{D} - \mathbf{D}) \nabla f_i(x^k) \right\|_{\mathbf{L}}^2 | x^k \right] + \|\nabla f(x^k)\|_{\mathbf{D} \mathbf{L} \mathbf{D}}^2. \end{aligned}$$

687 Recalling (34) we obtain

$$\begin{aligned} T &= \frac{1}{n^2} \sum_{i=1}^n \mathbb{E}[\|\mathbf{T}_i^k \mathbf{D} \nabla f_i(x^k) - \mathbf{D} \nabla f_i(x^k)\|_{\mathbf{L}}^2 | x^k] + \|\nabla f(x^k)\|_{\mathbf{D} \mathbf{L} \mathbf{D}}^2 \\ &\stackrel{(52)}{\leq} \frac{1}{n^2} \sum_{i=1}^n \mathbb{E}[\|\mathbf{T}_i^k \mathbf{D} \nabla f_i(x^k) - \mathbf{D} \nabla f_i(x^k)\|_{\mathbf{L}}^2 | x^k] + \|\nabla f(x^k)\|_{\mathbf{D}}^2. \end{aligned}$$

688 By applying Lemma 7, we get

$$\begin{aligned} T &\leq \frac{1}{n^2} \sum_{i=1}^n \lambda_{\max} \left(\mathbf{L}_i^{\frac{1}{2}} \mathbf{D} \mathbb{E}[(\mathbf{T}_i^k - \mathbf{I}_d) \mathbf{L}(\mathbf{T}_i^k - \mathbf{I}_d)] \mathbf{D} \mathbf{L}_i^{\frac{1}{2}} \right) \|\nabla f_i(x^k)\|_{\mathbf{L}_i^{-1}}^2 + \|\nabla f(x^k)\|_{\mathbf{D}}^2 \\ &\stackrel{(36)}{\leq} \lambda'_D \cdot \frac{2}{n} \left(f(x^k) - \frac{1}{n} \sum_{i=1}^n f_i^{\inf} \right) + \|\nabla f(x^k)\|_{\mathbf{D}}^2. \end{aligned}$$

Then we plug the upper bound of T back into (55), we get

$$\begin{aligned} & \mathbb{E} [f(x^{k+1}) | x^k] \\ & \leq f(x^k) - \frac{1}{2} \|\nabla f(x^k)\|_D^2 + \frac{\lambda'_D}{n} (f(x^k) - f^{\inf}) + \frac{\lambda'_D}{n} (f^{\inf} - \frac{1}{n} \sum_{i=1}^n f_i^{\inf}). \end{aligned}$$

Taking expectation, subtracting f^{\inf} from both sides, and using tower property, we get

$$\begin{aligned} & \mathbb{E} [f(x^{k+1}) - f^{\inf}] \\ & \leq \mathbb{E} [f(x^k) - f^{\inf}] - \frac{1}{2} \mathbb{E} [\|\nabla f(x^k)\|_D^2] + \frac{\lambda'_D}{n} \mathbb{E} [f(x^k) - f^{\inf}] + \frac{\lambda'_D}{n} \Delta^{\inf}. \end{aligned}$$

Then following similar steps as in the proof of Theorem 3, we are able to get

$$\min_{0 \leq k \leq K-1} \mathbb{E} \left[\|\nabla f(x^k)\|_{\frac{D}{\det(D)^{1/d}}}^2 \right] \leq \frac{2(1 + \frac{\lambda'_D}{n})^K (f(x^0) - f^{\inf})}{\det(D)^{1/d} K} + \frac{2\lambda'_D \Delta^{\inf}}{\det(D)^{1/d} n}.$$

This concludes the proof. \square

Due to the same reason for Algorithm 1, we don't need to worry about exponential blow-up in convergence bound (53) as we can always upper bound the exponential term by some constant by carefully controlling the stepsize matrix. Similarly, one corollary can be formulated to sum up the convergence conditions for Algorithm 2.

Corollary 4. We reach an error level of ε^2 in (53) if the following conditions are satisfied:

$$DL D \preceq D, \quad \lambda'_D \leq \min \left\{ \frac{n}{K}, \frac{n\varepsilon^2}{4\Delta^{\inf} \det(D)^{1/d}} \right\}, \quad K \geq \frac{12(f(x^0) - f^{\inf})}{\det(D)^{1/d} \varepsilon^2}. \quad (56)$$

The proof of this corollary is exactly the same as Corollary 1.

C.3.2 Optimal stepsize

In order to minimize the iteration complexity for Algorithm 2, the following optimization problem needs to be solved

$$\begin{aligned} & \min \quad \log \det(D^{-1}) \\ & \text{subject to} \quad D \text{ satisfies (56)} \end{aligned}$$

Following similar techniques in the proof of Proposition 2, we are able to prove that the above optimization problem is still a convex optimization problem. One simple way to find stepsize matrices is that we first fix $W \in \mathbb{S}_{++}^d$ and we find the optimal $0 < \gamma \in \mathbb{R}$, such that $D = \gamma W$ satisfies (56).

C.4 DCGD with constant stepsize

In this section we describe the convergence result for DCGD from [KR20]. We assume that the component functions f_i satisfy Assumption 4 with $L_i = L_i I_d$ and f satisfies Assumption 1 and 2 with $L = L I_d$. [KR20] proposed a unified analysis for non-convex optimization algorithms based on a generic upper bound on the second moment of the gradient estimator $g(x^k)$:

$$\mathbb{E} [\|g(x^k)\|^2] \leq 2A (f(x^k) - f^{\inf}) + B \|\nabla f(x^k)\|^2 + C, \quad (57)$$

In our case the gradient estimator is defined as follows

$$g_{\text{DCGD}}(x^k) = \frac{1}{n} \sum_{i=1}^n S_i^k \nabla f_i(x^k). \quad (58)$$

Here each S_i^k is the sketch matrix on the i -th client at the k -th iteration. One may check that g_{DCGD} satisfies (57) with the following constants:

$$A = \frac{\omega L_{\max}}{n}, \quad B = 1, \quad C = \frac{2\omega L_{\max}}{n} \Delta^{\inf}. \quad (59)$$

713 The constant L_{\max} is defined as the maximum of all L_i and $\omega = \lambda_{\max} \left(\mathbb{E} \left[(S_i^k)^\top S_i^k \right] \right) - 1$.
 714 Applying Corollary 1 from [KR20], we deduce the following. If

$$\gamma \leq \min \left\{ \frac{1}{L}, \frac{\sqrt{n}}{\sqrt{\omega L L_{\max} K}}, \frac{n\epsilon^2}{4LL_{\max}\omega\Delta_{\inf}} \right\} \text{ and } \gamma K \geq \frac{12(f(x^0) - f^{\inf})}{\epsilon^2}, \quad (60)$$

715 then

$$\min_{k=0, \dots, K-1} \mathbb{E} \left[\|\nabla f(x^k)\|^2 \right] \leq \epsilon^2. \quad (61)$$

716 D Proofs of technical lemmas

717 D.1 Proof of Lemma 1

718 Let us pick any two matrices $D_1, D_2 \in \mathbb{S}_{++}^d$, scalar α satisfying $0 \leq \alpha \leq 1$ and show that the
 719 following inequality holds regardless of the choice of R ,

$$f(\alpha D_1 + (1 - \alpha) D_2) \leq \alpha f(D_1) + (1 - \alpha) f(D_2). \quad (62)$$

720 For the LHS, we have

$$\begin{aligned} & f(\alpha D_1 + (1 - \alpha) D_2) \\ &= \text{tr} \left(L^{\frac{1}{2}} (\alpha D_1 + (1 - \alpha) D_2) R (\alpha D_1 + (1 - \alpha) D_2) L^{\frac{1}{2}} \right) \\ &= \alpha^2 \text{tr} (L^{\frac{1}{2}} D_1 R D_1 L^{\frac{1}{2}}) + (1 - \alpha)^2 \text{tr} (L^{\frac{1}{2}} D_2 R D_2 L^{\frac{1}{2}}) \\ &\quad + \alpha(1 - \alpha) \text{tr} (L^{\frac{1}{2}} D_1 R D_2 L^{\frac{1}{2}}) + \alpha(1 - \alpha) \text{tr} (L^{\frac{1}{2}} D_2 R D_1 L^{\frac{1}{2}}). \end{aligned}$$

721 and for the RHS, we have

$$\alpha f(D_1) + (1 - \alpha) f(D_2) = \alpha \text{tr} (L^{\frac{1}{2}} D_1 R D_1 L^{\frac{1}{2}}) + (1 - \alpha) \text{tr} (L^{\frac{1}{2}} D_2 R D_2 L^{\frac{1}{2}}).$$

722 Thus (62) can be simplified to the following inequality after rearranging terms

$$\begin{aligned} & \alpha(1 - \alpha) \text{tr} (L^{\frac{1}{2}} D_1 R D_2 L^{\frac{1}{2}}) + \alpha(1 - \alpha) \text{tr} (L^{\frac{1}{2}} D_2 R D_1 L^{\frac{1}{2}}) \\ & \leq \alpha(1 - \alpha) \text{tr} (L^{\frac{1}{2}} D_1 R D_1 L^{\frac{1}{2}}) + \alpha(1 - \alpha) \text{tr} (L^{\frac{1}{2}} D_2 R D_2 L^{\frac{1}{2}}). \end{aligned}$$

723 This is equivalent to

$$\text{tr} (L^{\frac{1}{2}} D_1 R D_1 L^{\frac{1}{2}}) + \text{tr} (L^{\frac{1}{2}} D_2 R D_2 L^{\frac{1}{2}}) - \text{tr} (L^{\frac{1}{2}} D_1 R D_2 L^{\frac{1}{2}}) - \text{tr} (L^{\frac{1}{2}} D_2 R D_1 L^{\frac{1}{2}}) \geq 0.$$

724 To show that the above inequality holds, we do the following transformation for the LHS

$$\begin{aligned} & \text{tr} (L^{\frac{1}{2}} D_1 R D_1 L^{\frac{1}{2}}) + \text{tr} (L^{\frac{1}{2}} D_2 R D_2 L^{\frac{1}{2}}) - \text{tr} (L^{\frac{1}{2}} D_1 R D_2 L^{\frac{1}{2}}) - \text{tr} (L^{\frac{1}{2}} D_2 R D_1 L^{\frac{1}{2}}) \\ &= \text{tr} (L^{\frac{1}{2}} D_1 R (D_1 - D_2) L^{\frac{1}{2}}) + \text{tr} (L^{\frac{1}{2}} D_2 R (D_2 - D_1) L^{\frac{1}{2}}) \\ &= \text{tr} (L^{\frac{1}{2}} (D_1 - D_2) R (D_1 - D_2) L^{\frac{1}{2}}). \end{aligned}$$

725 Since $R \in \mathbb{S}_+^d$ and $D_1 - D_2, L$ are symmetric, for any vector $u \in \mathbb{R}^d$

$$u^\top L^{\frac{1}{2}} (D_1 - D_2) R (D_1 - D_2) L^{\frac{1}{2}} u = \left((D_1 - D_2) L^{\frac{1}{2}} u \right)^\top R \left((D_1 - D_2) L^{\frac{1}{2}} u \right) \geq 0. \quad (63)$$

726 Thus, $L^{\frac{1}{2}} (D_1 - D_2) R (D_1 - D_2) L^{\frac{1}{2}} \in \mathbb{S}_+^d$, which yields the positivity of its trace. Therefore, (62)
 727 holds, thus $f(D)$ is a convex function. This concludes the proof.

728 **D.2 Proof of Lemma 2**

Proof.

$$\begin{aligned}
\mathbb{E} \left[\|S_i^k x - x\|_{DL D}^2 \right] &= \mathbb{E} [\langle (S_i^k - I_d)x, DL D(S_i^k - I_d)x \rangle] \\
&= \mathbb{E} [x^\top (S_i^k - I_d) DL D(S_i^k - I_d)x] \\
&= x^\top \mathbb{E} [(S_i^k - I_d) DL D(S_i^k - I_d)] x \\
&= x^\top L_i^{-\frac{1}{2}} \left(L_i^{\frac{1}{2}} \mathbb{E} [(S_i^k - I_d) DL D(S_i^k - I_d)] L_i^{\frac{1}{2}} \right) L_i^{-\frac{1}{2}} x \\
&\leq \lambda_{\max} \left(L_i^{\frac{1}{2}} \mathbb{E} [(S_i^k - I_d) DL D(S_i^k - I_d)] L_i^{\frac{1}{2}} \right) \|L_i^{-\frac{1}{2}} x\|^2 \\
&= \lambda_{\max} \left(L_i^{\frac{1}{2}} \mathbb{E} [(S_i^k - I_d) DL D(S_i^k - I_d)] L_i^{\frac{1}{2}} \right) \|x\|_{L_i^{-1}}^2.
\end{aligned}$$

729 This completes the proof. \square

730 **D.3 Proof of Lemma 3**

731 *Proof.* We have

$$\begin{aligned}
\mathbb{E} \left[\|x - \mathbb{E}[x]\|_M^2 \right] &= \mathbb{E} [\langle x - \mathbb{E}[x], M(x - \mathbb{E}[x]) \rangle] \\
&= \mathbb{E} [(x - \mathbb{E}[x])^\top M(x - \mathbb{E}[x])] \\
&= \mathbb{E} [x^\top Mx - \mathbb{E}[x]^\top Mx - x^\top M\mathbb{E}[x] + \mathbb{E}[x]^\top M\mathbb{E}[x]] \\
&= \mathbb{E} [x^\top Mx] - 2\mathbb{E}[x]^\top M\mathbb{E}[x] + \mathbb{E}[x]^\top M\mathbb{E}[x] \\
&= \mathbb{E} [x^\top Mx] - \mathbb{E}[x]^\top M\mathbb{E}[x] \\
&= \mathbb{E} [\|x\|_M^2] - \|\mathbb{E}[x]\|_M^2,
\end{aligned}$$

732 which concludes the proof. \square

733 **D.4 Proof of Lemma 4**

734 *Proof.* We have

$$\begin{aligned}
\mathbb{E} \left[\left\| \frac{1}{n} \sum_{i=1}^n a_i \right\|_M^2 \right] &= \frac{1}{n^2} \sum_{i=1}^n \mathbb{E} [\langle a_i, M a_i \rangle] + \frac{1}{n^2} \sum_{i \neq j} \mathbb{E} [\langle a_i, M a_j \rangle] \\
&= \frac{1}{n^2} \sum_{i=1}^n \mathbb{E} [\|a_i\|_M^2] + \frac{1}{n} \sum_{i \neq j} \langle \mathbb{E}[a_i], M \mathbb{E}[a_j] \rangle \\
&= \frac{1}{n^2} \sum_{i=1}^n \mathbb{E} [\|a_i\|_M^2].
\end{aligned}$$

735 This concludes the proof. \square

736 **D.5 Proof of Lemma 5**

737 *Proof.* Notice that

$$\mathbb{E}[Sx] = \mathbb{E}[S]x = x.$$

738 We start with variance decomposition in the matrix norm,

$$\begin{aligned}
\mathbb{E} \left[\|Sx - x\|_M^2 \right] &\stackrel{(33)}{=} \mathbb{E} \left[\|Sx\|_M^2 \right] - \|x\|_M^2 \\
&= \mathbb{E} [\langle Sx, MSx \rangle] - \langle x, Mx \rangle \\
&= \langle x, \mathbb{E} [SMS] x \rangle - \langle x, Mx \rangle \\
&= \langle x, (\mathbb{E} [SMS] - M) x \rangle \\
&= \|x\|_{\mathbb{E} [SMS] - M}^2.
\end{aligned}$$

739 This concludes the proof. \square

740 D.6 Proof of Lemma 6

741 *Proof.* We follow the definition of \mathbf{L} matrix smoothness of function f , that for any $x^+, x \in \mathbb{R}^d$, we
742 have

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

743 We plug in $x^+ = x - \mathbf{L}^{-1} \nabla f(x)$, and get

$$f^{\text{inf}} \leq f(x^+) \leq f(x) - \langle \nabla f(x), \mathbf{L}^{-1} \nabla f(x) \rangle + \frac{1}{2} \langle \nabla f(x), \mathbf{L}^{-1} \nabla f(x) \rangle.$$

744 Rearranging terms we get

$$\|\nabla f(x)\|_{\mathbf{L}^{-1}}^2 \leq 2(f(x) - f^{\text{inf}}), \quad (64)$$

745 which completes the proof. \square

746 D.7 Proof of Lemma 7

Proof.

$$\begin{aligned}
\mathbb{E} \left[\|T_i^k D x - D x\|_{\mathbf{L}}^2 \right] &= \mathbb{E} [\langle (T_i^k - I_d) D x, \mathbf{L}(T_i^k - I_d) D x \rangle] \\
&= \mathbb{E} [x^\top D (T_i^k - I_d) \mathbf{L} (T_i^k - I_d) D x] \\
&= x^\top D \mathbb{E} [(T_i^k - I_d) \mathbf{L} (T_i^k - I_d)] D x \\
&= x^\top \mathbf{L}_i^{-\frac{1}{2}} \left(\mathbf{L}_i^{\frac{1}{2}} D \mathbb{E} [(T_i^k - I_d) \mathbf{L} (T_i^k - I_d)] D \mathbf{L}_i^{\frac{1}{2}} \right) \mathbf{L}_i^{-\frac{1}{2}} x \\
&\leq \lambda_{\max} \left(\mathbf{L}_i^{\frac{1}{2}} D \mathbb{E} [(T_i^k - I_d) \mathbf{L} (T_i^k - I_d)] D \mathbf{L}_i^{\frac{1}{2}} \right) \left\| \mathbf{L}_i^{-\frac{1}{2}} x \right\|^2 \\
&= \lambda_{\max} \left(\mathbf{L}_i^{\frac{1}{2}} D \mathbb{E} [(T_i^k - I_d) \mathbf{L} (T_i^k - I_d)] D \mathbf{L}_i^{\frac{1}{2}} \right) \|x\|_{\mathbf{L}_i^{-1}}^2.
\end{aligned}$$

747 This completes the proof. \square

748 E Experiments

749 In this section, we describe the settings and results of numerical experiments to demonstrate the
750 effectiveness of our method. We perform several experiments under single node case and distributed
751 case. The code is available at https://anonymous.4open.science/r/detCGD_Code-A87D/.

752 E.1 Single node case

753 For single node case, we study the logistic regression problem with non-convex regularizer. The
754 objective is given as

$$f(x) = \frac{1}{n} \sum_{i=1}^n \log \left(1 + e^{-b_i \cdot \langle a_i, x \rangle} \right) + \lambda \cdot \sum_{j=1}^d \frac{x_j^2}{1 + x_j^2},$$

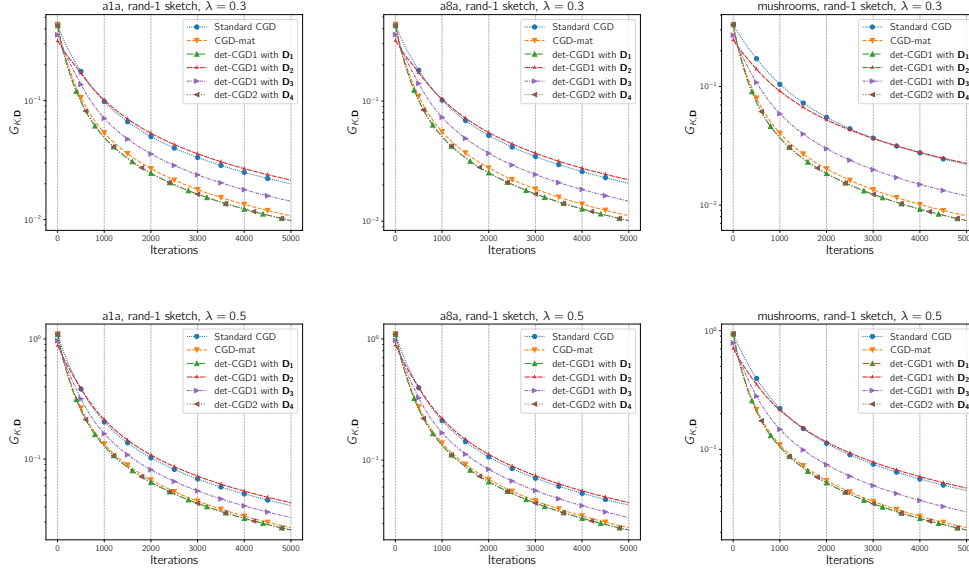


Figure 2: Comparison of standard CGD, CGD-mat, **det-CGD1** with $D_1 = \gamma_1 \cdot \text{diag}^{-1}(L)$, **det-CGD1** with $D_2 = \gamma_2 \cdot L^{-1}$, **det-CGD1** with $D_3 = \gamma_3 \cdot L^{-1/2}$ and **det-CGD2** with $D_4 = \gamma_4 \cdot \text{diag}^{-1}(L)$, where $\gamma_1, \gamma_2, \gamma_3$ are the optimal scaling factors for **det-CGD1** in that case, D_4 is the optimal matrix stepsize for **det-CGD2**. Rand-1 sketch is used in all the methods through out the experiments. The notation $G_{K,D}$ in the y -axis is defined in (65).

755 where $x \in \mathbb{R}^d$ is the model, $(a_i, b_i) \in \mathbb{R}^d \times \{-1, +1\}$ is one data point in the dataset whose size
 756 is n . $\lambda > 0$ is a constant associated with the regularizer. We conduct numerical experiments using
 757 several datasets from the LibSVM repository [CL11]. We estimate the smoothness matrix of function
 758 f here as

$$L = \frac{1}{n} \sum_{i=1}^n \frac{a_i a_i^\top}{4} + 2\lambda \cdot I_d.$$

759 E.1.1 Comparison to CGD with scalar stepsize, scalar smoothness constant

760 The purpose of the first experiment is to show that by using matrix stepsize, **det-CGD1** and **det-CGD2**
 761 will have better iteration and communication complexities compared to standard CGD (which uses
 762 scalar stepsize γ and scalar smoothness constant $L = \lambda_{\max}(L)$) and CGD with scalar stepsize $\gamma \cdot I_d$,
 763 smoothness matrix L . We use standard CGD to refer to CGD with scalar stepsize, scalar smoothness
 764 constant, and CGD-mat to refer to CGD with scalar stepsize, smoothness matrix in Figure 2, 3. The
 765 notation $G_{K,D}$ appears in the label of y axis is defined as

$$G_{K,D} := \frac{1}{K} \left(\sum_{k=0}^{K-1} \left\| \nabla f(x^k) \right\|_{\frac{D}{\det(D)^{1/d}}}^2 \right), \quad (65)$$

766 it is the average matrix norm of the gradient of f over the first $K - 1$ iterations in log scale. The
 767 weight matrix here has determinant 1, and thus it is comparable to the standard Euclidean norm. The
 768 result is meaningful in this sense.

769 The result presented in Figure 2 suggests that compared to standard CGD [KFJ18], CGD that uses
 770 smoothness matrices performs better in terms of both iteration complexity and communication
 771 complexity, while **det-CGD1** and **det-CGD2** with best diagonal matrix stepsizes outperform both of
 772 CGD and CGD with matrix smoothness which confirms our theory. The scaling factors $\gamma_1, \gamma_2, \gamma_3$
 773 here for **det-CGD1** are determined using Theorem 2 with $\ell = 1$. The matrix stepsize for **det-CGD2**
 774 is determined through (11). **det-CGD1** and **det-CGD2** with diagonal matrix stepsizes perform very
 775 similarly in the experiment, this is expected since we are using rand-1 sketch, which means that the

stepsize matrix and the sketch matrix are commutable since they are both diagonal. We also notice that **det-CGD1** with $D_2 = \gamma_2 \cdot L^{-1}$ is always worse than $D_4 = \gamma_4 \cdot \text{diag}^{-1}(L)$, this is also expected since we mentioned in Appendix B.5.1 that the result row 5 (corresponding to D_2) in Table 1 is always worse than row 7 (corresponding to D_4).

780 E.1.2 Comparison of the two algorithms under the same stepsize

781 The purpose of the second experiment is to compare the performance of **det-CGD1** and **det-CGD2** in
782 terms of iteration complexity and communication complexity. We know the conditions for **det-CGD1**
783 and **det-CGD2** to converge are given by (7) and (8) respectively, as a result, we are able to obtain the
784 optimal matrix stepsize for **det-CGD2** if we are using rand- τ sparsification. It is given by

$$D_2^* = \frac{\tau}{d} \left(\frac{d-\tau}{d-1} \text{diag}(L) + \frac{\tau-1}{d-1} L \right)^{-1},$$

785 according to (11). The definition of $G_{K,D}$ is given in (65), τ here for random sparsification is set to
786 be integers around $\{\frac{d}{4}, \frac{d}{2}, \frac{3d}{4}\}$, where d is the dimension of the model.

787 It can be observed from the result presented in Figure 3, that in almost all cases in this experiment, 2
788 with $D = D_2^*$ outperforms the other methods. Compared to standard CGD and CGD with matrix
789 stepsize, **det-CGD1** and **det-CGD2** are always better. This provides numerical evidence in support of
790 our theory. In this case, the stepsize matrix is not diagonal for **det-CGD1** and **det-CGD2**, so we do
791 not expect them to perform similarly. Notice that in dataset phishing, the four algorithms behave
792 very similarly, this is because the smoothness matrix L here has a very centralized spectrum.

793 E.2 Distributed case

794 For distributed case, we still use the logistic regression problem with non-convex regularizer as our
795 experiment setting. The objective is given similarly as

$$f(x) = \frac{1}{n} \sum_{i=1}^n f_i(x); \quad f_i(x) = \frac{1}{m_i} \sum_{j=1}^{m_i} \log \left(1 + e^{-b_{i,j} \cdot \langle a_{i,j}, x \rangle} \right) + \lambda \cdot \sum_{t=1}^d \frac{x_t^2}{1 + x_t^2},$$

796 where $x \in \mathbb{R}^d$ is the model, $(a_{i,j}, b_{i,j}) \in \mathbb{R}^d \times \{-1, +1\}$ is one data point in the dataset of client i
797 whose size is m_i . $\lambda > 0$ is a constant associated with the regularizer. For each dataset used in the
798 distributed setting, we randomly reshuffled the dataset before splitting it equally to each client. We
799 estimate the smoothness matrices of function f and each individual function f_i here as

$$\begin{aligned} L_i &= \frac{1}{m_i} \sum_{i=1}^{m_i} \frac{a_i a_i^\top}{4} + 2\lambda \cdot I_d; \\ L &= \frac{1}{n} \sum_{i=1}^n L_i. \end{aligned}$$

800 The value of Δ^{inf} here is determined in the following way, we first perform gradient descent on f
801 and record the minimum value in the entire run, f^{inf} , as the estimate of its global minimum, then we
802 do the same procedure for each f_i to obtain the estimate of its global minimum f_i^{inf} . After that we
803 estimate Δ^{inf} using its definition.

804 E.2.1 Comparison to standard DCGD in the distributed case

805 This experiment is designed to show that D-**det-CGD1** and D-**det-CGD2** will have better iteration
806 complexity and communication complexity compared to standard DCGD [KFJ18] and DCGD with
807 scalar stepsize, smoothness matrix. We will use standard DCGD here to refer to DCGD with scalar
808 stepsize, scalar smoothness constant, and DCGD-mat here to refer to DCGD with scalar stepsize,
809 smoothness matrix. Rand-1 sparsifier is used in all the algorithms throughout the experiment. The
810 error level is fixed as $\varepsilon^2 = 0.0001$, the conditions for standard DCGD to converge can be deduced
811 using Proposition 4 in [KR20], we use the largest possible scalar stepsize here for standard DCGD.
812 The optimal scalar stepsize for DCGD-mat, optimal diagonal matrix stepsize D_1 for D-**det-CGD1**
813 and D_2 for D-**det-CGD2** can be determined using Corollary 1.

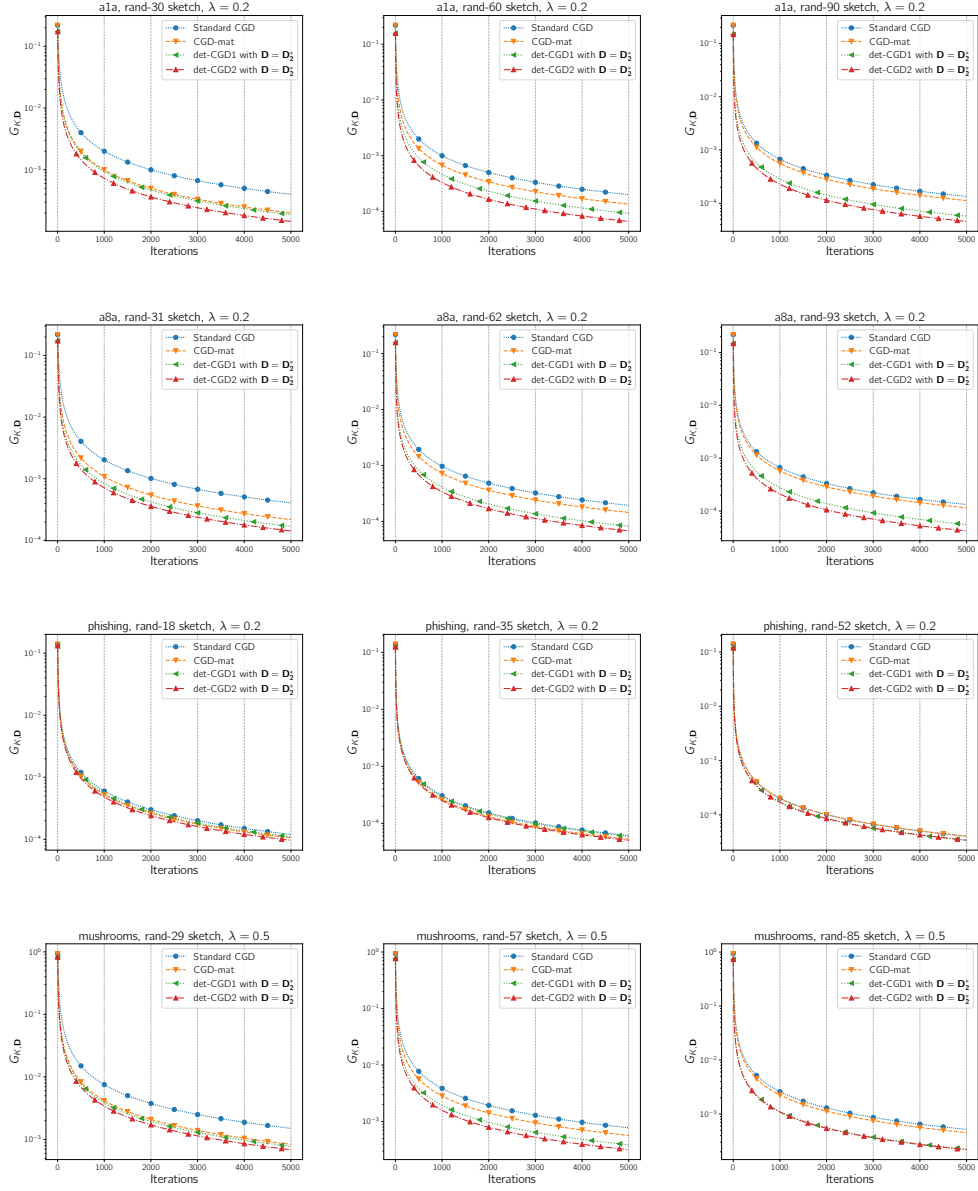


Figure 3: Comparison of standard CGD, CGD-mat **det-CGD1** with stepsize $D = D_2^*$ and **det-CGD2** with stepsize $D = D_2^*$, where D_2^* is the optimal stepsize matrix for **det-CGD1** and the optimal diagonal stepsize matrix for **det-CGD2**. Rand- τ sketch is used in all the algorithms throughout the experiments. The notation $G_{K,D}$ in the y -axis is defined in (65).

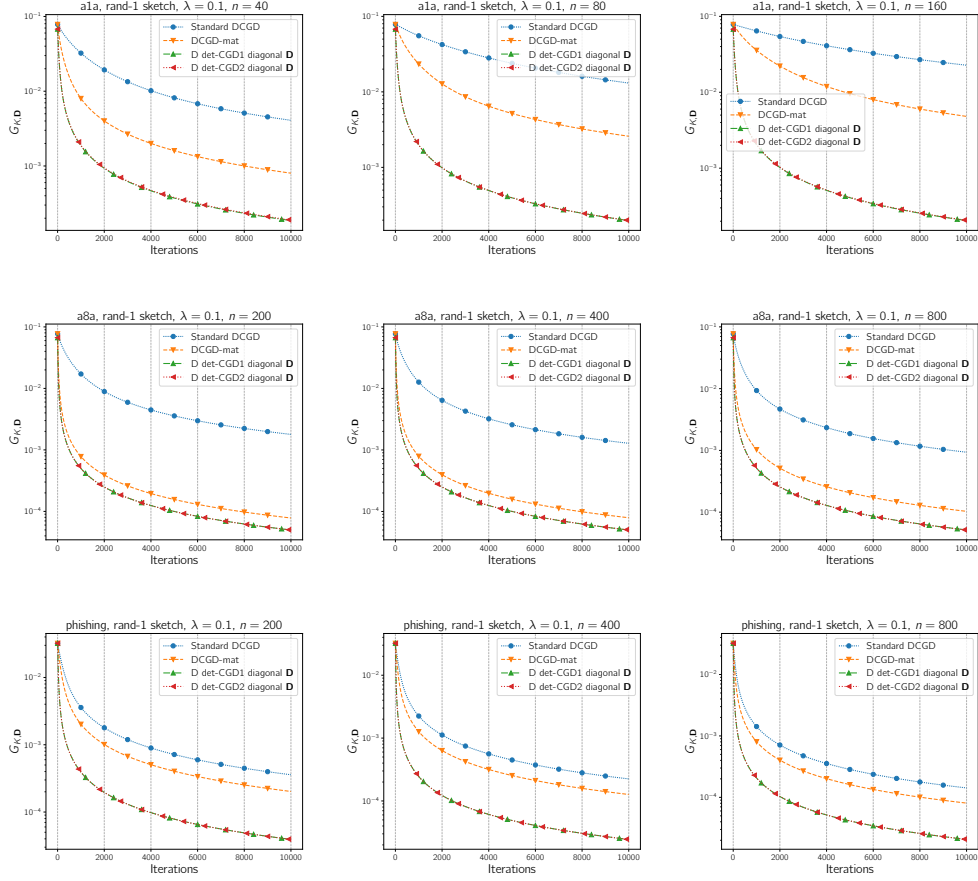


Figure 4: Comparison of standard DCGD, DCGD-mat, D-det-CGD1 with matrix stepsize \mathbf{D}_1 and D-det-CGD2 with matrix stepsize \mathbf{D}_2 , where $\mathbf{D}_1, \mathbf{D}_2$ are the optimal diagonal matrix stepsizes for D-det-CGD1 and D-det-CGD2 respectively. Rand-1 sketch is used in all the algorithms throughout the experiment. The notation $G_{K,D}$ in the y -axis is defined in (65).

From the result of Figure 4, we are able to see that both D-det-CGD1 and D-det-CGD2 outperform standard DCGD and DCGD-mat in terms of iteration complexity and communication complexity, which confirms our theory. Notice that D-det-CGD1, D-det-CGD2 are expected to perform very similarly because the stepsize matrix and sketches are diagonal which means that they are commutable. We also plot the corresponding standard Euclidean norm of iterates of D-det-CGD1 and D-det-CGD2 in Figure 5, the E_K here appears in the y -axis is defined as,

$$E_K := \frac{1}{K} \sum_{k=0}^{K-1} \|\nabla f(x^k)\|^2. \quad (66)$$

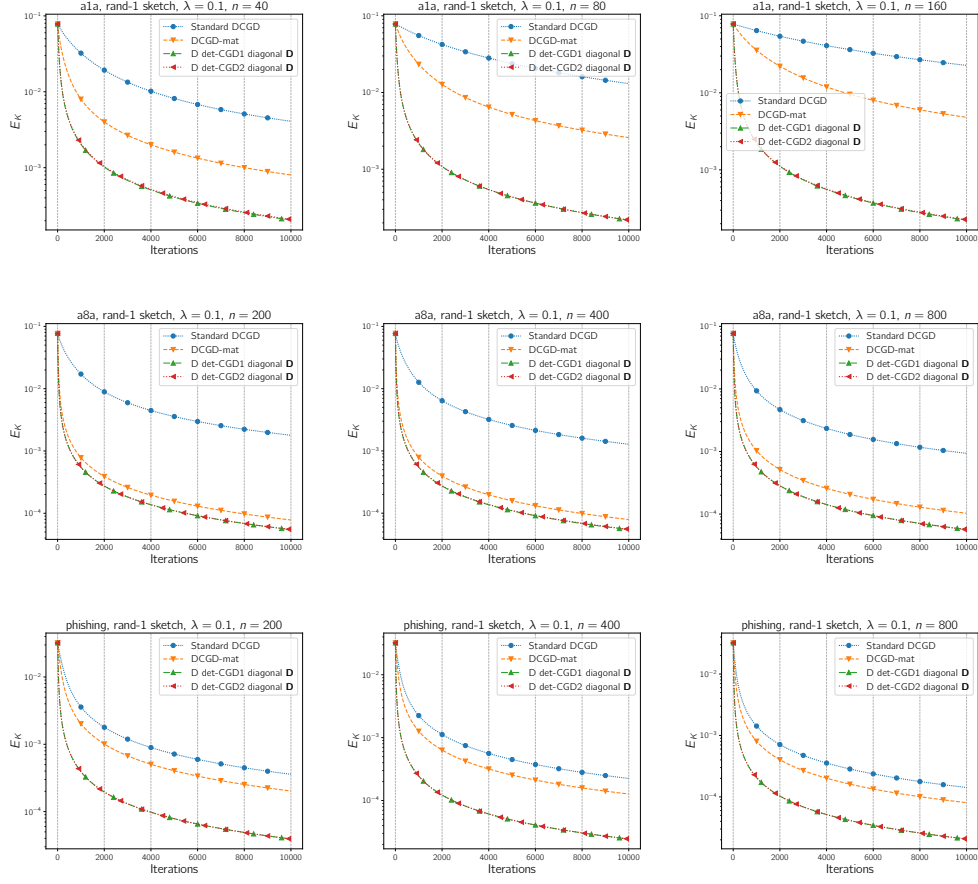


Figure 5: Comparison of standard DCGD, DCGD-mat, D-det-CGD1 with matrix stepsize D_1 and D-det-CGD2 with matrix stepsize D_2 , where D_1, D_2 are the optimal diagonal matrix stepsizes for D-det-CGD1 and D-det-CGD2 respectively. Rand-1 sketch is used in all the algorithms throughout the experiment. The y -axis is now standard Euclidean norm defined in (66).