

VARIANCE REDUCED DISTRIBUTED NONCONVEX OPTIMIZATION USING MATRIX STEPSIZES

Anonymous authors

Paper under double-blind review

ABSTRACT

Matrix-stepsized gradient descent algorithms have been shown to have superior performance in non-convex optimization problems compared to their scalar counterparts. The **det-CGD** algorithm, as introduced by Li et al. (2024b), leverages matrix stepsizes to perform compressed gradient descent for non-convex objectives and matrix-smooth problems in a federated manner. The authors establish the algorithm’s convergence to a neighborhood of a weighted stationarity point under a convex condition for the symmetric and positive-definite matrix stepsize. In this paper, we propose two variance-reduced versions of the **det-CGD** algorithm, incorporating **MARINA** and **DASHA** methods. Notably, we establish theoretically and empirically, that **det-MARINA** and **det-DASHA** outperform **MARINA**, **DASHA** and the distributed **det-CGD** algorithms in terms of iteration and communication complexities.

1 INTRODUCTION

We focus on optimizing the finite sum non-convex objective

$$\min_{x \in \mathbb{R}^d} \left\{ f(x) := \frac{1}{n} \sum_{i=1}^n f_i(x) \right\}. \quad (1)$$

In this context, each function $f_i : \mathbb{R}^d \rightarrow \mathbb{R}$ is differentiable and bounded from below. This type of objective function finds extensive application in various practical machine learning algorithms, which increase not only in terms of the data size but also in the model size and overall complexity as well. As a result, most neural network architectures result in highly non-convex empirical losses, which need to be minimized. In addition, it becomes computationally infeasible to train these models on one device, often excessively large, and one needs to redistribute them amongst different devices/clients. This redistribution results in a high communication overhead, which often becomes the bottleneck in this framework.

In other words, we have the following setting. The data is partitioned into n clients, where the i -th client has access to the component function f_i and its derivatives. The clients are connected to each other through a central device, called the server. In this work, we are going to study iterative gradient descent-based algorithms that operate as follows. The clients compute the local gradients in parallel. Then they compress these gradients to reduce the communication cost and send them to the server in parallel. The server then aggregates these vectors and broadcasts the iterate update back to the clients. This meta-algorithm is called federated learning. We refer the readers to Konečný et al. (2016); McMahan et al. (2017); Kairouz et al. (2021) for a more thorough introduction to federated learning.

1.1 CONTRIBUTIONS

In this paper, we introduce two novel federated learning algorithms named **det-MARINA** and **det-DASHA**. These algorithms extend a recent method called **det-CGD** (Li et al., 2024b), which aims to solve problem (1) using matrix stepsized gradient descent. Under the matrix smoothness assumption proposed by Safaryan et al. (2021), the authors demonstrate that the matrix stepsized

version of the Distributed Compressed Gradient Descent (Khirirat et al., 2018) algorithm enhances communication complexity compared to its scalar counterpart. However, in their analysis, Li et al. (2024b) show stationarity only within a certain neighborhood due to stochastic compressors. Our algorithm addresses this issue by incorporating previously known variance reduction schemes, namely, **MARINA** (Gorbunov et al., 2021) and **DASHA** (Tyurin & Richtárik, 2024). We establish theoretically and empirically, that both algorithms outperform their scalar alternatives, as well as the distributed **det-CGD** algorithms. In addition, we describe specific matrix stepsize choices, for which our algorithms beat **MARINA**, **DASHA** and distributed **det-CGD** both in theory and in practice.

2 BACKGROUND AND MOTIVATION

For a given $\varepsilon > 0$, finding an approximately global optimum, that is x_ε such that $f(x_\varepsilon) - \min_x f(x) < \varepsilon$, is known to be NP-hard (Jain et al., 2017; Danilova et al., 2022). However, gradient descent based methods are still useful in this case. When these methods are applied to non-convex objectives, they treat the function f as locally convex and aim to converge to a local minimum. Despite this simplification, such methods have gained popularity in practice due to their superior performance compared to other approaches for non-convex optimization, such as convex relaxation-based methods (Tibshirani, 1996; Cai et al., 2010).

2.1 STOCHASTIC GRADIENT DESCENT

Arguably, one of the most prominent meta-methods for tackling non-convex optimization problems is stochastic gradient descent (**SGD**). The formulation of **SGD** is presented as the following iterative algorithm: $x^{k+1} = x^k - \gamma g^k$. Here, $g^k \in \mathbb{R}^d$ serves as a stochastic estimator of the gradient $\nabla f(x^k)$. **SGD** essentially mimics the classical gradient descent algorithm, and recovers it when $g^k = \nabla f(x^k)$. In this scenario, the method approximates the objective function f using a linear function and takes a step of size γ in the direction that maximally reduces this approximation. When the stepsize is sufficiently small, and the function f is suitably smooth, it can be demonstrated that the function value decreases, as discussed in (Bubeck et al., 2015; Gower et al., 2019).

However, computing the full gradient can often be computationally expensive. In such cases, stochastic approximations of the gradient come into play. Stochastic estimators of the gradient can be employed for various purposes, leading to the development of different methods. These include stochastic batch gradient descent (Nemirovski et al., 2009; Johnson & Zhang, 2013; Defazio et al., 2014), randomized coordinate descent (Nesterov, 2012; Wright, 2015), and compressed gradient descent (Alistarh et al., 2017; Khirirat et al., 2018; Mishchenko et al., 2019). The latter, compressed gradient descent, holds particular relevance to this paper, and we will delve into a more detailed discussion of it in subsequent sections.

2.2 SECOND ORDER METHODS

The stochastic gradient descent is considered as a first-order method as it uses only the first order derivative information. Although being immensely popular, the first order methods are not always optimal. Not surprisingly, using higher order derivatives in deciding update direction can yield to faster algorithms. A simple instance of such algorithms is the Newton Star algorithm (Islamov et al., 2021):

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

where x^* is the minimum point of the objective function. The authors establish that under specific conditions, the algorithm’s convergence to the unique solution x^* in the convex scenario occurs at a local quadratic rate. Nonetheless, its practicality is limited since we do not have prior knowledge of the Hessian matrix at the optimal point. Despite being proposed recently, the Newton-Star algorithm gives a deeper insight on the generic Newton method (Gragg & Tapia, 1974; Miel, 1980; Yamamoto, 1987):

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

Here, the unknown Hessian of the Newton-Star algorithm, is estimated progressively along the iterations. The latter causes elevated computational costs, as the inverting a large square matrix is

expensive. As an alternative, quasi-Newton methods replace the inverse of the Hessian at the iterate with a computationally cheaper estimate (Broyden, 1965; Dennis & Moré, 1977; Al-Baali & Khalfan, 2007; Al-Baali et al., 2014).

2.3 FIXED MATRIX STEPSIZES

The **det-CGD** algorithm falls into this framework of the second order methods as well. Proposed by Li et al. (2024b)¹, the algorithm suggests using a uniform “upper bound” on the inverse Hessian matrix. Assuming matrix smoothness of the objective (Safaryan et al., 2021), they replace the scalar stepsize with a positive definite matrix \mathbf{D} . The algorithm is given as follows:

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

Matrix \mathbf{D} . Here, \mathbf{D} plays the role of the stepsize. Essentially, it uniformly lower bounds the inverse Hessian. The standard **SGD** is a particular case of this method, as the scalar stepsize γ can be seen as a matrix $\gamma\mathbf{I}_d$, where \mathbf{I}_d is the d -dimensional identity matrix. An advantage of using a matrix stepsize is more evident if we take the perspective of the second order methods. Indeed, the scalar stepsize $\gamma\mathbf{I}_d$ uniformly estimates the largest eigenvalue of the Hessian matrix, while \mathbf{D} can capture the Hessian more accurately. The authors show both theoretical and empirical improvement that comes with matrix stepsizes.

Matrix \mathbf{S}^k . We assume that \mathbf{S}^k is a positive semi-definite, stochastic sketch matrix. Furthermore, it is unbiased: $\mathbb{E}[\mathbf{S}^k] = \mathbf{I}_d$. We notice that **det-CGD** can be seen as a matrix stepsize instance of **SGD**, with $g^k = \mathbf{S}^k \nabla f(x^k)$. The sketch matrix can be seen as a linear compressing operator, hence the name of the algorithm: Compressed Gradient Descent (**CGD**) (Alistarh et al., 2017; Khirirat et al., 2018). A commonly used example of such a compressor is the Rand- τ compressor. This compressor randomly selects τ entries from its input and scales them using a scalar multiplier to ensure an unbiased estimation. By adopting this approach, instead of using all d coordinates of the gradient, only a subset of size τ is communicated. Formally, Rand- τ is defined as follows:

$$\mathbf{S} = \frac{d}{\tau} \sum_{j=1}^{\tau} e_{i_j} e_{i_j}^{\top}. \quad (2)$$

Here, e_{i_j} denotes the i_j -th standard basis vector in \mathbb{R}^d . For a more comprehensive understanding of compression techniques, we refer to Safaryan et al. (2022b).

2.4 THE NEIGHBORHOOD OF THE DISTRIBUTED DET-CGD1

The distributed version of **det-CGD** follows the standard federated learning paradigm (McMahan et al., 2017). The pseudocode of the method, as well as the convergence result of Li et al. (2024b), can be found in Appendix F. Informally, their convergence result can be written as

$$\min_{k=1, \dots, K} \mathbb{E} \left[\|\nabla f(x^k)\|_{\mathbf{D}}^2 \right] \leq \mathcal{O} \left(\frac{(1 + \alpha)^K}{K} \right) + \mathcal{O}(\alpha),$$

where $\alpha > 0$ is a constant that can be controlled. The crucial insight from this result is that the error bound does not diminish as the number of iterations increases. In fact, by controlling α and considering a large K , it is impossible to make the second term smaller than ε . This implies that the algorithm converges to a certain neighborhood surrounding the (local) optimum. This phenomenon is a common occurrence in **SGD** and is primarily attributable to the variance associated with the stochastic gradient estimator. In the case of **det-CGD** the stochasticity comes from the sketch \mathbf{S}^k .

2.5 VARIANCE REDUCTION

To eliminate this neighborhood, various techniques for reducing variance are employed. One of the simplest techniques applicable to **CGD** is gradient shifting. By replacing $\mathbf{S}^k \nabla f(x^k)$ with

¹In the original paper, the algorithm is referred to as **det-CGD**, as there is a variant of the same algorithm named **det-CGD2**. Since we are going to use only the first one and our framework is applicable to both, we will remove the number in the end for the sake of brevity.

162 $S^k(\nabla f(x^k) - \nabla f(x^*)) + \nabla f(x^*)$, the neighborhood effect is removed from the general **CGD**. This
 163 algorithm is an instance of a more commonly known method called **SGD*** (Gower et al., 2020).
 164 However, since the exact optimum x^* is typically unknown, this technique encounters similar chal-
 165 lenges as the Newton-Star algorithm mentioned earlier. Fortunately, akin to quasi-Newton methods,
 166 one can employ methods that iteratively learn the optimal shift (Shulgin & Richtárik, 2022). A line
 167 of research focuses on variance reduction for **CGD** based algorithms on this insight.

168 To eliminate the neighborhood in the distributed version of **CGD**, denoted as **det-CGD1**, we apply
 169 a technique called **MARINA** (Gorbunov et al., 2021). **MARINA** cleverly combines the general
 170 shifting (Shulgin & Richtárik, 2022) technique with loopless variance reduction techniques (Qian
 171 et al., 2021). This approach introduces an alternative gradient estimator specifically designed for
 172 the federated learning setting. Thanks to its structure, it allows to establish an upper bound on the
 173 stationarity error that diminishes significantly with a large number of iterations. In this paper, we
 174 construct the analog of the this algorithm called **det-MARINA**, using matrix stepsizes and sketch
 175 gradient compressors. For this new method, we prove a convergence guarantee similar to the results
 176 of Li et al. (2024b) without a neighborhood term.

177 Furthermore, we also propose **det-DASHA**, which is the extension of **DASHA** in the matrix step-
 178 size setting. The latter was proposed by Tyurin & Richtárik (2024) and it combines **MARINA**
 179 with momentum variance reduction techniques (Cutkosky & Orabona, 2019). **DASHA** offers better
 180 practicality compared to **MARINA**, as it always sends compressed gradients and does not need to
 181 synchronize among all the nodes.

182 2.6 ORGANIZATION OF THE PAPER

183 The rest of the paper is organized as follows. Section 3 discusses the general mathematical frame-
 184 work. Section 4 and Section 5 present the **det-MARINA** and **det-DASHA** algorithms, respectively.
 185 We show the superior theoretical performance of our algorithms compared to the relevant existing
 186 algorithms, that is **MARINA**, **DASHA** and **det-CGD** in Section 6. The experimental results validat-
 187 ing our theoretical findings are presented in Section 7, with additional details and setups available in
 188 the Appendix. We conclude the paper by outlining several directions of future work in Section 8.

191 3 MATHEMATICAL FRAMEWORK

192 In this section we present the assumptions that we further require in the analysis.

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

194 This is a standard assumption in optimization, as otherwise the problem of minimizing the objective
 195 would not be correct mathematically. We then need a matrix version of Lipschitz continuity for the
 196 gradient.

197 **Definition 1.** (**L**-Lipschitz Gradient) *Assume that $f : \mathbb{R}^d \rightarrow \mathbb{R}$ is a continuously differentiable
 198 function and matrix $\mathbf{L} \in \mathbb{S}_{++}^d$. We say the gradient of f is **L**-Lipschitz if for all $x, y \in \mathbb{R}^d$*

$$199 \|\nabla f(x) - \nabla f(y)\|_{\mathbf{L}^{-1}} \leq \|x - y\|_{\mathbf{L}}. \quad (3)$$

200 In the following, we will assume that (3) is satisfied for component functions f_i .

201 **Assumption 2.** *Each function f_i is \mathbf{L}_i -gradient Lipschitz, while f is **L**-gradient Lipschitz.*

202 In fact, the second half of the assumption is a consequence of the first one. Below, we formalize this
 203 claim.

204 **Proposition 1.** *If f_i is \mathbf{L}_i -gradient Lipschitz for every $i = 1, \dots, n$, then function f has **L**-Lipschitz
 205 gradient with $\mathbf{L} \in \mathbb{S}_{++}^d$ satisfying*

$$206 \frac{1}{n} \sum_{i=1}^n \lambda_{\max}(\mathbf{L}^{-1}) \cdot \lambda_{\max}(\mathbf{L}_i) \cdot \lambda_{\max}(\mathbf{L}_i \mathbf{L}^{-1}) = 1.$$

207 **Remark 1.** *In the scalar case, where $\mathbf{L} = L\mathbf{I}_d$, $\mathbf{L}_i = L_i\mathbf{I}_d$, the relation becomes $L^2 = \frac{1}{n} \sum_{i=1}^n L_i^2$.
 208 This corresponds to the statement in Assumption 1.2 in (Gorbunov et al., 2021).*

Nevertheless, the matrix \mathbf{L} found according to Proposition 1 is only an estimate. In principle, there might exist a better $\mathbf{L}_f \preceq \mathbf{L}$ such that f has \mathbf{L}_f -Lipschitz gradient.

More generally, this condition can be interpreted as follows. The gradient of f naturally belongs to the dual space of \mathbb{R}^d , as it is defined as a linear functional on \mathbb{R}^d . In the scalar case, ℓ_2 -norm is self-dual, thus (3) reduces to the standard Lipschitz continuity of the gradient. However, with the matrix smoothness assumption, we are using the \mathbf{L} -norm for the iterates, which naturally induces the \mathbf{L}^{-1} -matrix norm for the gradients in the dual space. This insight, which is originally presented by Nemirovski & Yudin (1983), plays a key role in our analysis.

See Appendix C for a more thorough discussion on the properties of Assumption 2, as well as its connection to matrix smoothness (Safaryan et al., 2021).

4 MARINA-BASED VARIANCE REDUCTION

In this section, we present our algorithm **det-MARINA** with its convergence result. We construct a sequence of vectors g^k which are stochastic estimators of $\nabla f(x^k)$. At each iteration, the server samples a Bernoulli random variable (coin flip) c_k and broadcasts it in parallel to the clients, along with the current gradient estimate g^k . Each client, then, does a **det-CGD**-type update with the stepsize \mathbf{D} and a gradient estimate g^k . The next gradient estimate g^{k+1} is then computed. With a low probability, that is when $c_k = 1$, we take the g^{k+1} to be the full gradient $\nabla f(x^{k+1})$. Otherwise, we update it using the compressed gradient differences at each client. See Algorithm 1 for the pseudocode of **det-MARINA**.

Algorithm 1 **det-MARINA**

```

1: Input: starting point  $x^0$ , stepsize matrix  $\mathbf{D}$ , probability  $p \in (0, 1]$ , number of iterations  $K$ 
2: Initialize  $g^0 = \nabla f(x^0)$ 
3: for  $k = 0, 1, \dots, K - 1$  do
4:   Sample  $c_k \sim \text{Be}(p)$ 
5:   Broadcast  $g^k$  to all workers
6:   for  $i = 1, 2, \dots$  in parallel do
7:      $x^{k+1} = x^k - \mathbf{D} \cdot g^k$ 
8:     if  $c_k = 1$  then
9:        $g_i^{k+1} = \nabla f_i(x^{k+1})$ 
10:    else
11:       $g_i^{k+1} = g^k + \mathbf{S}_i^k (\nabla f_i(x^{k+1}) - \nabla f_i(x^k))$ 
12:    end if
13:  end for
14:   $g^{k+1} = \frac{1}{n} \sum_{i=1}^n g_i^{k+1}$ 
15: end for
16: Return:  $\tilde{x}^K$  chosen uniformly at random from  $\{x^k\}_{k=0}^{K-1}$ 

```

4.1 CONVERGENCE GUARANTEES

In the following theorem, we formulate one of the main results of this paper, which guarantees the convergence of Algorithm 1 under the above-mentioned assumptions.

Theorem 1. *Assume that Assumptions 1 and 2 hold, and the following condition on stepsize matrix $\mathbf{D} \in \mathbb{S}_{++}^d$ holds,*

$$\mathbf{D}^{-1} \succeq \left(\frac{(1-p) \cdot R(\mathbf{D}, \mathcal{S})}{np} + 1 \right) \mathbf{L}, \quad (4)$$

where $R(\mathbf{D}, \mathcal{S}) := \frac{1}{n} \sum_{i=1}^n \lambda_{\max}(\mathbf{L}_i) \lambda_{\max} \left(\mathbf{L}^{-\frac{1}{2}} \mathbf{L}_i \mathbf{L}^{-\frac{1}{2}} \right) \times \lambda_{\max} \left(\mathbb{E} [\mathbf{S}_i^k \mathbf{D} \mathbf{S}_i^k] - \mathbf{D} \right)$. Then, after K iterations of **det-MARINA**, we have

$$\mathbb{E} \left[\|\nabla f(\tilde{x}^K)\|_{\frac{\mathbf{D}}{\det(\mathbf{D})^{1/d}}}^2 \right] \leq \frac{2(f(x^0) - f^*)}{\det(\mathbf{D})^{1/d} \cdot K}. \quad (5)$$

Here, \tilde{x}^K is chosen uniformly randomly from the first K iterates of the algorithm.

The criterion $\|\cdot\|_{\mathbf{D}/\det(\mathbf{D})^{1/d}}^2$ is the same as that used in Li et al. (2024b), known as determinant normalization. The weight matrix of the matrix norm has determinant 1 after normalization, which makes it comparable to the standard Euclidean norm.

Remark 2. We notice that the right-hand side of the algorithm vanishes with the number of iterations, thus solving the neighborhood issue of the distributed *det-CGD*. Therefore, *det-MARINA* is indeed the variance reduced version of *det-CGD* in the distributed setting and has better convergence guarantees.

Remark 3. Theorem 1 implies the following iteration complexity for the algorithm. In order to get an ε^2 stationarity error², the algorithm requires K iterations, with

$$K \geq \frac{2(f(x^0) - f^*)}{\det(\mathbf{D})^{1/d} \cdot \varepsilon^2}.$$

Remark 4. In the case where no compression is applied, that is we have $\mathbf{S}_i^k = \mathbf{I}_d$, condition (4) reduces to $\mathbf{D} \preceq \mathbf{L}^{-1}$. The latter is due to $\mathbb{E}[\mathbf{S}_i^k \mathbf{D} \mathbf{S}_i^k] = \mathbf{D}$, which results in $R(\mathbf{D}, \mathcal{S}) = 0$. This is expected, since in the deterministic case *det-MARINA* reduces to *GD* with matrix stepsize.

The convergence condition and rate of matrix stepsize *GD* can be found in (Li et al., 2024b). Below we do a sanity check to verify that the convergence condition for scalar *MARINA* can be obtained.

Remark 5. Let us consider the scalar case. That is $\mathbf{L}_i = L_i \mathbf{I}_d$, $\mathbf{L} = L \mathbf{I}_d$, $\mathbf{D} = \gamma \mathbf{I}_d$ and $\omega = \lambda_{\max} \left(\mathbb{E} \left[(\mathbf{S}_i^k)^\top \mathbf{S}_i^k \right] \right) - 1$. Then, the condition (4) reduces to

$$\gamma \leq \left[L \left(1 + \sqrt{\frac{(1-p)\omega}{pn}} \right) \right]^{-1}.$$

The latter coincides with the stepsize condition of the convergence result of scalar *MARINA*.

4.2 OPTIMIZING THE MATRIX STEPSIZE

Now let us look at the right-hand side of (5). We notice that it decreases in terms of the determinant of the stepsize matrix. Therefore, one needs to solve the following optimization problem to find the optimal stepsize:

$$\begin{aligned} & \text{minimize} && \log \det(\mathbf{D}^{-1}) \\ & \text{subject to} && \mathbf{D} \text{ satisfying (4)}. \end{aligned}$$

The solution of this constrained minimization problem on \mathbb{S}_{++}^d is not explicit. In theory, one may show that the constraint (4) is convex and attempt to solve the problem numerically. However, as stressed by Li et al. (2024b), the similar stepsize condition for *det-CGD* is not easily computed using solvers like CVXPY (Diamond & Boyd, 2016). Instead, we may relax the problem to certain linear subspaces of \mathbb{S}_{++}^d . In particular, we fix a matrix $\mathbf{W} \in \mathbb{S}_{++}^d$, and define $\mathbf{D} := \gamma \mathbf{W}$. Then, the condition on the matrix \mathbf{D} becomes a condition for the scalar γ , which is given in the following corollary.

Corollary 1. Let $\mathbf{W} \in \mathbb{S}_{++}^d$, defining $\mathbf{D} := \gamma \cdot \mathbf{W}$, where $\gamma \in \mathbb{R}_+$. then the condition in (4) reduces to the following condition on γ

$$\gamma \leq \frac{2\lambda_{\mathbf{W}}}{1 + \sqrt{1 + 4\alpha\beta \cdot \Lambda_{\mathbf{W},\mathcal{S}}\lambda_{\mathbf{W}}}}, \quad (6)$$

where $\Lambda_{\mathbf{W},\mathcal{S}} = \lambda_{\max} \left(\mathbb{E}[\mathbf{S}_i^k \mathbf{W} \mathbf{S}_i^k] - \mathbf{W} \right)$, $\lambda_{\mathbf{W}} = \lambda_{\max}^{-1}(\mathbf{W}^{\frac{1}{2}} \mathbf{L} \mathbf{W}^{\frac{1}{2}})$, $\alpha = \frac{1-p}{np}$ and $\beta = \frac{1}{n} \sum_{i=1}^n \lambda_{\max}(\mathbf{L}_i) \cdot \lambda_{\max}(\mathbf{L}^{-1} \mathbf{L}_i)$.

²We say a (possibly random) vector $x \in \mathbb{R}^d$ is an ε -stationary point of a possibly non-convex function $f: \mathbb{R}^d \mapsto \mathbb{R}$, if $\mathbb{E}[\|\nabla f(x)\|^2] \leq \varepsilon^2$. The expectation is over the randomness of the algorithm

This means that for every fixed \mathbf{W} , we can find the optimal scaling coefficient γ . In section Section 6, we will use this corollary to prove that a suboptimal matrix step size, determined in this efficient way, is already better than the optimal scalar step size.

Extension to det-CGD2. A variant of **det-CGD**, called **det-CGD2**, was also proposed by Li et al. (2024b). This algorithm, has the same structure as **det-CGD** with the sketch and stepsize interchanged. It was shown, that this algorithm has explicit stepsize condition in the single node setting. In Appendix G, we propose the variance reduced extension of the distributed **det-CGD2** following the **MARINA** scheme.

5 DASHA-BASED VARIANCE REDUCTION

In this section, we present our second algorithm based on **DASHA**. The latter utilizes a different type of variance reduction based on momentum (MVR). Compared to **MARINA**, **dasha** makes simpler optimization steps and does not require periodic synchronization with all the nodes. Notice that one may further simplify the notations here used in the algorithm. However, we keep it this way as it is consistent with (Tyurin & Richtárik, 2024).

Algorithm 2 det-DASHA

1: **Input:** starting point $x^0 \in \mathbb{R}^d$, stepsize matrix $\mathbf{D} \in \mathbb{S}_{++}^d$, momentum $a \in (0, 1]$, number of iterations K
2: Initialize $g_i^0, h_i^0 \in \mathbb{R}^d$ on the nodes and $g^0 = \frac{1}{n} \sum_{i=1}^n g_i^0$ on the server
3: **for** $k = 0, 1, \dots, K - 1$ **do**
4: $x^{k+1} = x^k - \mathbf{D} \cdot g^k$
5: Broadcast x^{k+1} to all nodes
6: **for** $i = 1, 2, \dots, n$ in parallel **do**
7: $h_i^{k+1} = \nabla f_i(x^{k+1})$
8: $m_i^{k+1} = \mathbf{S}_i^k (h_i^{k+1} - h_i^k - a (g_i^k - h_i^k))$
9: $g_i^{k+1} = g_i^k + m_i^{k+1}$
10: Send m_i^{k+1} to the server.
11: **end for**
12: $g^{k+1} = g^k + \frac{1}{n} \sum_{i=1}^n m_i^{k+1}$
13: **end for**
14: **Return:** \tilde{x}^K chosen uniformly at random from $\{x^k\}_{k=0}^{K-1}$

5.1 THEORETICAL GUARANTEES

Theorem 2. Suppose that Assumptions 1 and 2 hold. Let us initialize $g_i^0 = h_i^0 = \nabla f_i(x^0)$ for all $i \in [n]$ in Algorithm 2, and define

$$\Lambda_{\mathbf{D}, \mathbf{S}} = \lambda_{\max} (\mathbb{E} [\mathbf{S}_i^k \mathbf{D} \mathbf{S}_i^k] - \mathbf{D}), \quad \omega_{\mathbf{D}} = \lambda_{\max} (\mathbf{D}^{-1}) \cdot \Lambda_{\mathbf{D}, \mathbf{S}}.$$

If $a = \frac{1}{2\omega_{\mathbf{D}} + 1}$, and the following condition on stepsize $\mathbf{D} \in \mathbb{S}_{++}^d$ is satisfied

$$\mathbf{D}^{-1} \succeq \mathbf{L} - \frac{4\lambda_{\max}(\mathbf{D})\omega_{\mathbf{D}}(4\omega_{\mathbf{D}} + 1)}{n^2} \sum_{i=1}^n \lambda_{\max}(\mathbf{L}_i) \mathbf{L}_i,$$

then the following inequality holds for the iterates of Algorithm 2

$$\mathbb{E} \left[\|\nabla f(\tilde{x}^K)\|_{\mathbf{D}/(\det(\mathbf{D}))^{1/d}}^2 \right] \leq \frac{2(f(x^0) - f^*)}{\det(\mathbf{D})^{1/d} \cdot K}.$$

Here \tilde{x}^K is chosen uniformly randomly from the first K iterates of the algorithm.

Remark 6. The term $\Lambda_{\mathbf{D}, \mathbf{S}}$ can be viewed as the matrix version of $\gamma \cdot \omega$, where ω is associated with the sketch, and γ is the scalar stepsize. On the other hand, the $\omega_{\mathbf{D}}$ is the extension of ω in matrix norm. Similar to Remark 5, plugging in scalar arguments in the algorithm, we recover the result from Tyurin & Richtárik (2024).

Following the same scheme as in Section 4, we choose $\mathbf{D} = \gamma_{\mathbf{W}} \cdot \mathbf{W}$, where $\mathbf{W} \in \mathbb{S}_{++}^d$. Thus, for a fixed \mathbf{W} , we relax the problem of finding the optimal stepsize to the problem of finding the optimal scaling factor $\gamma_{\mathbf{W}} > 0$.

Corollary 2. For a fixed $\mathbf{W} \in \mathbb{S}_{++}^d$, the optimal scaling factor $\gamma_{\mathbf{W}} \in \mathbb{R}_+$ is given by

$$\gamma_{\mathbf{W}} = \frac{2\lambda_{\mathbf{W}}}{1 + \sqrt{1 + 16C_{\mathbf{W}}\lambda_{\min}(\mathbf{L}) \cdot \lambda_{\mathbf{W}}}},$$

where $C_{\mathbf{W}} := \lambda_{\max}(\mathbf{W}) \cdot \omega_{\mathbf{W}}(4\omega_{\mathbf{W}} + 1)/n$, and $\lambda_{\mathbf{W}} := \lambda_{\max}^{-1}(\mathbf{L}^{\frac{1}{2}}\mathbf{W}\mathbf{L}^{\frac{1}{2}})$.

We observe that the structure of the optimal scaling factor for obtained above is similar to the one obtained in Corollary 1.

The availability of \mathbf{L} : For both **det-MARINA** and **det-DASHA**, in order to determine the matrix stepsize, the knowledge of \mathbf{L} is needed, if \mathbf{L} is known, better complexities are guaranteed. When \mathbf{L} is unknown, a closed-form solution can be obtained for generalized linear models. In more general cases, \mathbf{L}_i can be treated as hyperparameters and estimated using first-order information via a gradient-based method (Wang et al., 2022). One can think of this as some type of preprocessing step, after which the matrices are learnt.

6 COMPLEXITIES OF THE ALGORITHMS

6.1 DET-MARINA

The following corollary formulates the iteration complexity for **det-MARINA** for $\mathbf{W} = \mathbf{L}^{-1}$.

Corollary 3. If we take $\mathbf{W} = \mathbf{L}^{-1}$, then the condition (6) on γ is given by

$$\gamma \leq 2 \left(1 + \sqrt{1 + 4\alpha\beta \cdot \Lambda_{\mathbf{L}^{-1}, \mathcal{S}}} \right)^{-1}. \quad (7)$$

In order to satisfy ε -stationarity, that is $\mathbb{E} \left[\left\| \nabla f(\tilde{x}^K) \right\|_{\frac{\mathbf{D}}{\det(\mathbf{D})^{1/d}}}^2 \right] \leq \varepsilon^2$, we require

$$K \geq \mathcal{O} \left(\frac{\Delta_0 \cdot \det(\mathbf{L})^{\frac{1}{d}}}{\varepsilon^2} \cdot \left(1 + \sqrt{1 + 4\alpha\beta \cdot \Lambda_{\mathbf{L}^{-1}, \mathcal{S}}} \right) \right),$$

where $\Delta_0 := f(x^0) - f(x^*)$. Moreover, this iteration complexity is always better than the one of **MARINA**.

The proof can be found in the Appendix. In fact, we can show that in cases where we fix $\mathbf{W} = \mathbf{I}_d$ and $\mathbf{W} = \text{diag}^{-1}(\mathbf{L})$, the same conclusion also holds, relevant details can be found in Appendix D.3. This essentially means that **det-MARINA** always has a “larger” stepsize compared to **MARINA**, even if the stepsize is suboptimal for the sake of efficiency, which leads to a better iteration complexity. In addition, because we are using the same compressor for those two algorithms, the communication complexity of **det-MARINA** is also provably better than that of **MARINA**.

In order to compute the communication complexity, we borrow the concept of expected density from Gorbunov et al. (2021).

Definition 2. For a given sketch matrix $\mathbf{S} \in \mathbb{S}_+^d$, the expected density is defined as

$$\zeta_{\mathbf{S}} = \sup_{x \in \mathbb{R}^d} \mathbb{E} [\| \mathbf{S}x \|_0],$$

where $\|x\|_0$ denotes the number of non-zero components of $x \in \mathbb{R}^d$.

In particular, we have $\zeta_{\text{Rand-}\tau} = \tau$. Below, we state the communication complexity of **det-MARINA** with $\mathbf{W} = \mathbf{L}^{-1}$ and the Rand- τ compressor.

Corollary 4. Assume that we are using sketch $\mathbf{S} \sim \mathcal{S}$ with expected density $\zeta_{\mathcal{S}}$. Suppose also we are running *det-MARINA* with probability p and we use the optimal stepsize matrix with respect to $\mathbf{W} = \mathbf{L}^{-1}$. Then the overall communication complexity of the algorithm is given by $\mathcal{O}((Kp + 1)d + (1 - p)K\zeta_{\mathcal{S}})$. Specifically, if we pick $p = \zeta_{\mathcal{S}}/d$, then the communication complexity is given by

$$\mathcal{O}\left(d + \frac{\Delta_0 \det(\mathbf{L})^{\frac{1}{d}}}{\varepsilon^2} \left(\zeta_{\mathcal{S}} + \sqrt{\frac{\beta}{n} \Lambda_{\mathbf{L}^{-1}, \mathcal{S}} \zeta_{\mathcal{S}} (d - \zeta_{\mathcal{S}})}\right)\right).$$

Notice that in case where no compression is applied, the communication complexity reduces to $\mathcal{O}(d\Delta_0 \cdot \det(\mathbf{L})^{\frac{1}{d}}/\varepsilon^2)$. The latter coincides with the rate of matrix stepsize GD (see (Li et al., 2024b)). Therefore, the dependence on ε is not possible to improve further since GD is optimal among first order methods (Carmon et al., 2020).

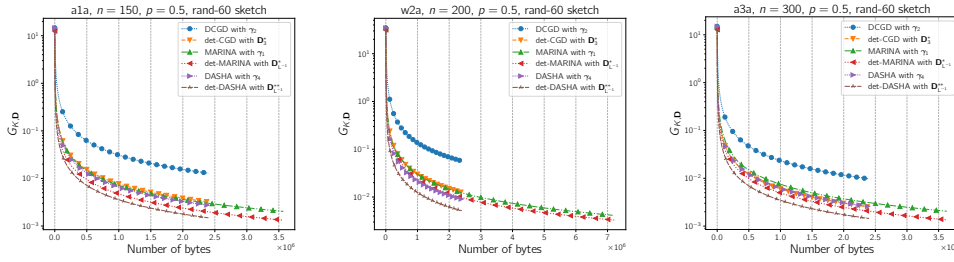


Figure 1: Comparison of DCGD with optimal scalar stepsize, *det-CGD* with matrix stepsize D_3^* , *MARINA* with optimal scalar stepsize, *DASHA* with optimal scalar stepsize, *det-MARINA* with optimal stepsize $D_{L^{-1}}^*$ and *det-DASHA* with optimal stepsize $D_{L^{-1}}^*$. Throughout the experiment, we are using Rand- τ sketch with $\tau = 60$, and each algorithm is run for a fixed number of iterations $K = 10000$. The $G_{K,D}$ in the y-axis is defined in (63), which is the average squared matrix norm of the gradients.

6.2 DET-DASHA

The difference of compression mechanisms, does not allow to have a direct comparison of the complexities of these algorithms. In particular, *det-MARINA* compresses the gradient difference with some probability p , while *det-DASHA* compresses the gradient difference with momentum in each iteration.

Corollary 5. If we pick $\mathbf{D} = \gamma_{L^{-1}} \cdot \mathbf{L}^{-1}$, then in order to reach an ε^2 stationary point, *det-DASHA* needs K iterations with

$$K \geq \frac{f(x^0) - f^*}{\det(\mathbf{L})^{-\frac{1}{d}} \varepsilon^2} \left(1 + \sqrt{1 + 16C_{L^{-1}} \lambda_{\min}(\mathbf{L})}\right).$$

The following corollary compares the complexities of *DASHA* and *det-DASHA*. For the sake of brevity, we defer the complexities and other details to the proof of this corollary.

Corollary 6. Suppose that the conditions in Theorem 2 hold, then compared to *DASHA*, *det-DASHA* with $\mathbf{W} = \mathbf{L}^{-1}$ always has a **better** iteration complexity, therefore, communication complexity as well.

The following corollary suggests that the communication complexity of *det-DASHA* is better than that of *det-MARINA*,

Corollary 7. The iteration complexity of *det-MARINA* with $p = 1/(\omega_{L^{-1}} + 1)$ and *det-DASHA* with momentum $1/(2\omega_{L^{-1}} + 1)$ is the same, therefore the communication complexity of *det-DASHA* is **better than** the communication complexity of *det-MARINA*.

This is expected since the same relation occurs between *MARINA* and *DASHA* as it is described by Tyurin & Richtárik (2024, Table 1). We refer the readers to Appendix E.2.1.

7 EXPERIMENTS

This section contains several plots which confirm our theoretical improvements on the existing methods. Figure 1 shows that the performance in terms of communication complexity of **det-DASHA** and **det-MARINA** is better than their scalar counterpart **DASHA** and **MARINA** respectively. This validates the efficiency of using a matrix stepsize over a scalar stepsize. Further, we notice that **det-DASHA** and **det-MARINA** have better communication complexity in this case, compared to **det-CGD**. This demonstrates the effectiveness of applying variance reduction. Finally, as expected, **det-DASHA** has better communication complexity than **det-MARINA**. We refer the readers to the appendix for more technical details of the experiments.

8 FUTURE WORK

i) In this paper, we have only considered (linear) sketches as the compression operator. However, there exists a variety of compressors which are useful in practice that do not fall into this category. Extending **det-CGD** and **det-MARINA** for general unbiased compressors is a promising future work direction. ii) Additionally, given recent successes with adaptive stepsizes (e.g., (Loizou et al., 2021; Orvieto et al., 2022; Schaipp et al., 2023)), designing an adaptive matrix stepsize tailored to our case could be viable. iii) Finally, recent advances suggest that server step sizes play a key role in accelerating federated learning algorithms (Jhunjunwala et al., 2023; Li et al., 2024a). Designing a matrix version of the server step size could also be interesting.

REFERENCES

- Mehiddin Al-Baali and H Khalfan. An overview of some practical quasi-Newton methods for unconstrained optimization. *Sultan Qaboos University Journal for Science [SQUJS]*, 12(2):199–209, 2007.
- Mehiddin Al-Baali, Emilio Spedicato, and Francesca Maggioni. Broyden’s quasi-Newton methods for a nonlinear system of equations and unconstrained optimization: a review and open problems. *Optimization Methods and Software*, 29(5):937–954, 2014.
- Foivos Alimisis, Peter Davies, and Dan Alistarh. Communication-efficient distributed optimization with quantized preconditioners. In *International Conference on Machine Learning*, pp. 196–206. PMLR, 2021.
- 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.
- Zeyuan Allen-Zhu. Katyusha: The first direct acceleration of stochastic gradient methods. In *Proceedings of the 49th Annual ACM SIGACT Symposium on Theory of Computing*, pp. 1200–1205, 2017.
- Rajendra Bhatia. *Positive definite matrices*. Princeton University Press, 2009.
- Charles G Broyden. A class of methods for solving nonlinear simultaneous equations. *Mathematics of Computation*, 19(92):577–593, 1965.
- Sébastien Bubeck et al. Convex optimization: Algorithms and complexity. *Foundations and Trends® in Machine Learning*, 8(3-4):231–357, 2015.
- Jian-Feng Cai, Emmanuel J Candès, and Zuowei Shen. A singular value thresholding algorithm for matrix completion. *SIAM Journal on Optimization*, 20(4):1956–1982, 2010.
- Yair Carmon, John C Duchi, Oliver Hinder, and Aaron Sidford. Lower bounds for finding stationary points i. *Mathematical Programming*, 184(1-2):71–120, 2020.
- 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.

- 540 Sélím Chraïbi, Ahmed Khaled, Dmitry Kovalev, Peter Richtárik, Adil Salim, and Martin Takáč. Dis-
541 tributed fixed point methods with compressed iterates. *arXiv preprint arXiv:1912.09925*, 2019.
542
- 543 Rixon Crane and Fred Roosta. Dingo: Distributed Newton-type method for gradient-norm optimiza-
544 tion. *Advances in Neural Information Processing Systems*, 32, 2019.
- 545 Ashok Cutkosky and Francesco Orabona. Momentum-based variance reduction in non-convex SGD.
546 *Advances in Neural Information Processing Systems*, 32, 2019.
547
- 548 Marina Danilova, Pavel Dvurechensky, Alexander Gasnikov, Eduard Gorbunov, Sergey Guminov,
549 Dmitry Kamzolov, and Innokentiy Shibaev. Recent theoretical advances in non-convex optimiza-
550 tion. In *High-Dimensional Optimization and Probability: With a View Towards Data Science*, pp.
551 79–163. Springer, 2022.
- 552 Aaron Defazio, Francis Bach, and Simon Lacoste-Julien. SAGA: A fast incremental gradient
553 method with support for non-strongly convex composite objectives. *Advances in Neural Informa-
554 tion Processing Systems*, 27, 2014.
- 555 John E Dennis, Jr and Jorge J Moré. Quasi-Newton methods, motivation and theory. *SIAM Review*,
556 19(1):46–89, 1977.
557
- 558 Steven Diamond and Stephen Boyd. CVXPY: A Python-embedded modeling language for convex
559 optimization. *The Journal of Machine Learning Research*, 17(1):2909–2913, 2016.
560
- 561 Canh T Dinh, Nguyen H Tran, Tuan Dung Nguyen, Wei Bao, Albert Y Zomaya, and Bing B Zhou.
562 Federated learning with proximal stochastic variance reduced gradient algorithms. In *Proceedings
563 of the 49th International Conference on Parallel Processing*, pp. 1–11, 2020.
- 564 Darina Dvinskikh, Aleksandr Ogaltsov, Alexander Gasnikov, Pavel Dvurechensky, Alexander
565 Tyurin, and Vladimir Spokoiny. Adaptive gradient descent for convex and non-convex stochastic
566 optimization. *arXiv preprint arXiv:1911.08380*, 2019.
- 567 Eduard Gorbunov, Konstantin P Burlachenko, Zhize Li, and Peter Richtárik. MARINA: Faster non-
568 convex distributed learning with compression. In *International Conference on Machine Learning*,
569 pp. 3788–3798. PMLR, 2021.
570
- 571 Robert M Gower, Mark Schmidt, Francis Bach, and Peter Richtárik. Variance-reduced methods for
572 machine learning. *Proceedings of the IEEE*, 108(11):1968–1983, 2020.
- 573 Robert Mansel Gower, Nicolas Loizou, Xun Qian, Alibek Sailanbayev, Egor Shulgin, and Peter
574 Richtárik. SGD: General analysis and improved rates. In *International Conference on Machine
575 Learning*, pp. 5200–5209. PMLR, 2019.
576
- 577 William B Gragg and Richard A Tapia. Optimal error bounds for the Newton–Kantorovich theorem.
578 *SIAM Journal on Numerical Analysis*, 11(1):10–13, 1974.
- 579 SV Guminov, Yu E Nesterov, PE Dvurechensky, and AV Gasnikov. Accelerated primal-dual gra-
580 dient descent with linesearch for convex, nonconvex, and nonsmooth optimization problems. In
581 *Doklady Mathematics*, volume 99, pp. 125–128. Springer, 2019.
582
- 583 Filip Hanzely and Peter Richtárik. Federated learning of a mixture of global and local models. *arXiv
584 preprint arXiv:2002.05516*, 2020.
- 585 Samuel Horváth, Chen-Yu Ho, Ludovit Horvath, Atal Narayan Sahu, Marco Canini, and Peter
586 Richtárik. Natural compression for distributed deep learning. In *Mathematical and Scientific
587 Machine Learning*, pp. 129–141. PMLR, 2022.
- 588 Samuel Horváth, Dmitry Kovalev, Konstantin Mishchenko, Peter Richtárik, and Sebastian Stich.
589 Stochastic distributed learning with gradient quantization and double-variance reduction. *Opti-
590 mization Methods and Software*, 38(1):91–106, 2023.
591
- 592 Rustem Islamov, Xun Qian, and Peter Richtárik. Distributed second order methods with fast rates
593 and compressed communication. In *International Conference on Machine Learning*, pp. 4617–
4628. PMLR, 2021.

- 594 Sashank J Reddi, Suvrit Sra, Barnabas Poczos, and Alexander J Smola. Proximal stochastic methods
595 for nonsmooth nonconvex finite-sum optimization. *Advances in Neural Information Processing*
596 *Systems*, 29, 2016.
- 597
- 598 Prateek Jain, Purushottam Kar, et al. Non-convex optimization for machine learning. *Foundations*
599 *and Trends® in Machine Learning*, 10(3-4):142–363, 2017.
- 600 Divyansh Jhunjhunwala, Shiqiang Wang, and Gauri Joshi. FedExP: Speeding up federated averaging
601 via extrapolation. In *International Conference on Learning Representations*, 2023.
- 602
- 603 Rie Johnson and Tong Zhang. Accelerating stochastic gradient descent using predictive variance
604 reduction. *Advances in neural information processing systems*, 26, 2013.
- 605
- 606 Peter Kairouz, H Brendan McMahan, Brendan Avent, Aurélien Bellet, Mehdi Bennis, Arjun Nitin
607 Bhagoji, Kallista Bonawitz, Zachary Charles, Graham Cormode, Rachel Cummings, et al. Ad-
608 vances and open problems in federated learning. *Foundations and Trends® in Machine Learning*,
609 14(1–2):1–210, 2021.
- 610 Ahmed Khaled and Peter Richtárik. Better theory for SGD in the nonconvex world. *Transactions*
611 *on Machine Learning Research*, 2023.
- 612
- 613 Prashant Khanduri, Pranay Sharma, Swatantra Kafle, Saikiran Bulusu, Ketan Rajawat, and
614 Pramod K Varshney. Distributed stochastic non-convex optimization: Momentum-based variance
615 reduction. *arXiv preprint arXiv:2005.00224*, 2020.
- 616 Sarit Khirirat, Hamid Reza Feyzmahdavian, and Mikael Johansson. Distributed learning with com-
617 pressed gradients. *arXiv preprint arXiv:1806.06573*, 2018.
- 618
- 619 Jakub Konečný, H Brendan McMahan, Felix X Yu, Peter Richtárik, Ananda Theertha Suresh, and
620 Dave Bacon. Federated learning: Strategies for improving communication efficiency. *arXiv*
621 *preprint arXiv:1610.05492*, 8, 2016.
- 622 Dmitry Kovalev, Samuel Horváth, and Peter Richtárik. Don’t jump through hoops and remove those
623 loops: SVRG and Katyusha are better without the outer loop. In *Algorithmic Learning Theory*,
624 pp. 451–467. PMLR, 2020.
- 625
- 626 Hanmin Li, Kirill Acharya, and Peter Richtárik. The power of extrapolation in federated learning.
627 *arXiv preprint arXiv:2405.13766*, 2024a.
- 628
- 629 Hanmin Li, Avetik Karagulyan, and Peter Richtárik. Det-CGD: Compressed gradient descent with
630 matrix stepsizes for non-convex optimization. In *International Conference on Learning Repre-*
631 *sentations*, 2024b.
- 632 Junyi Li, Feihu Huang, and Heng Huang. Local stochastic bilevel optimization with momentum-
633 based variance reduction. *arXiv preprint arXiv:2205.01608*, 2022.
- 634
- 635 Zhize Li, Dmitry Kovalev, Xun Qian, and Peter Richtárik. Acceleration for compressed gradient
636 descent in distributed and federated optimization. *arXiv preprint arXiv:2002.11364*, 2020.
- 637 Zhize Li, Hongyan Bao, Xiangliang Zhang, and Peter Richtárik. PAGE: A simple and optimal prob-
638 abilistic gradient estimator for nonconvex optimization. In *International Conference on Machine*
639 *Learning*, pp. 6286–6295. PMLR, 2021.
- 640
- 641 Deyi Liu, Lam M Nguyen, and Quoc Tran-Dinh. An optimal hybrid variance-reduced algorithm for
642 stochastic composite nonconvex optimization. *arXiv preprint arXiv:2008.09055*, 2020.
- 643
- 644 Nicolas Loizou, Sharan Vaswani, Issam Hadj Laradji, and Simon Lacoste-Julien. Stochastic Polyak
645 step-size for SGD: An adaptive learning rate for fast convergence. In *International Conference*
646 *on Artificial Intelligence and Statistics*, pp. 1306–1314. PMLR, 2021.
- 647 Julien Mairal. Incremental majorization-minimization optimization with application to large-scale
machine learning. *SIAM Journal on Optimization*, 25(2):829–855, 2015.

- 648 Artavazd Maranjyan, Mher Safaryan, and Peter Richtárik. GradSkip: Communication-accelerated
649 local gradient methods with better computational complexity. *arXiv preprint arXiv:2210.16402*,
650 2022.
- 651 Brendan McMahan, Eider Moore, Daniel Ramage, Seth Hampson, and Blaise Aguera y Arcas.
652 Communication-efficient learning of deep networks from decentralized data. In *Artificial Intelli-*
653 *gence and Statistics*, pp. 1273–1282. PMLR, 2017.
- 654 George J Miel. Majorizing sequences and error bounds for iterative methods. *Mathematics of*
655 *Computation*, 34(149):185–202, 1980.
- 656 Konstantin Mishchenko, Eduard Gorbunov, Martin Takáč, and Peter Richtárik. Distributed learning
657 with compressed gradient differences. *arXiv preprint arXiv:1901.09269*, 2019.
- 658 Konstantin Mishchenko, Grigory Malinovsky, Sebastian Stich, and Peter Richtarik. ProxSkip: Yes!
659 Local gradient steps provably lead to communication acceleration! Finally! In Kamalika Chaud-
660 huri, Stefanie Jegelka, Le Song, Csaba Szepesvari, Gang Niu, and Sivan Sabato (eds.), *Proceed-*
661 *ings of the 39th International Conference on Machine Learning*, volume 162 of *Proceedings of*
662 *Machine Learning Research*, pp. 15750–15769. PMLR, 17–23 Jul 2022.
- 663 Arkadi Nemirovski, Anatoli Juditsky, Guanghui Lan, and Alexander Shapiro. Robust stochastic
664 approximation approach to stochastic programming. *SIAM Journal on Optimization*, 19(4):1574–
665 1609, 2009.
- 666 Arkadi Semenovič Nemirovski and David Borisovich Yudin. Problem complexity and method effi-
667 ciency in optimization. *Wiley-Interscience, ISSN 0277-2698*, 1983.
- 668 Yu Nesterov. Efficiency of coordinate descent methods on huge-scale optimization problems. *SIAM*
669 *Journal on Optimization*, 22(2):341–362, 2012.
- 670 Antonio Orvieto, Simon Lacoste-Julien, and Nicolas Loizou. Dynamics of SGD with stochastic
671 Polyak stepsizes: Truly adaptive variants and convergence to exact solution. *Advances in Neural*
672 *Information Processing Systems*, 35:26943–26954, 2022.
- 673 Jie Peng, Zhaoxian Wu, Qing Ling, and Tianyi Chen. Byzantine-robust variance-reduced federated
674 learning over distributed non-iid data. *Information Sciences*, 616:367–391, 2022.
- 675 Xun Qian, Zheng Qu, and Peter Richtárik. L-SVRG and L-Katyusha with arbitrary sampling. *The*
676 *Journal of Machine Learning Research*, 22(1):4991–5039, 2021.
- 677 Peter Richtárik and Martin Takáč. Iteration complexity of randomized block-coordinate descent
678 methods for minimizing a composite function. *Mathematical Programming*, 144(1-2):1–38, 2014.
- 679 Mher Safaryan, Filip Hanzely, and Peter Richtárik. Smoothness matrices beat smoothness constants:
680 Better communication compression techniques for distributed optimization. *Advances in Neural*
681 *Information Processing Systems*, 34:25688–25702, 2021.
- 682 Mher Safaryan, Rustem Islamov, Xun Qian, and Peter Richtarik. FedNL: Making Newton-type
683 methods applicable to federated learning. In *International Conference on Machine Learning*, pp.
684 18959–19010. PMLR, 2022a.
- 685 Mher Safaryan, Egor Shulgin, and Peter Richtárik. Uncertainty principle for communication com-
686 pression in distributed and federated learning and the search for an optimal compressor. *Informa-*
687 *tion and Inference: A Journal of the IMA*, 11(2):557–580, 2022b.
- 688 Fabian Schaipp, Robert M Gower, and Michael Ulbrich. A stochastic proximal Polyak step size.
689 *Transactions on Machine Learning Research*, 2023.
- 690 Mark Schmidt, Nicolas Le Roux, and Francis Bach. Minimizing finite sums with the stochastic
691 average gradient. *Mathematical Programming*, 162:83–112, 2017.
- 692 Egor Shulgin and Peter Richtárik. Shifted compression framework: Generalizations and improve-
693 ments. In *Uncertainty in Artificial Intelligence*, pp. 1813–1823. PMLR, 2022.

- 702 Robert Tibshirani. Regression shrinkage and selection via the LASSO. *Journal of the Royal Statis-*
703 *tical Society Series B: Statistical Methodology*, 58(1):267–288, 1996.
- 704
- 705 Quoc Tran-Dinh, Nhan H Pham, Dzung T Phan, and Lam M Nguyen. A hybrid stochastic opti-
706 mization framework for composite nonconvex optimization. *Mathematical Programming*, 191
707 (2):1005–1071, 2022.
- 708 Alexander Tyurin and Peter Richtárik. DASHA: Distributed nonconvex optimization with commu-
709 nication compression and optimal oracle complexity. In *International Conference on Learning*
710 *Representations*, 2024.
- 711 Bokun Wang, Mher Safaryan, and Peter Richtárik. Theoretically better and numerically faster dis-
712 tributed optimization with smoothness-aware quantization techniques. *Advances in Neural Infor-*
713 *mation Processing Systems*, 35:9841–9852, 2022.
- 714
- 715 Shusen Wang, Fred Roosta, Peng Xu, and Michael W Mahoney. Giant: Globally improved approx-
716 imate Newton method for distributed optimization. *Advances in Neural Information Processing*
717 *Systems*, 31, 2018.
- 718
- 719 Stephen J Wright. Coordinate descent algorithms. *Mathematical programming*, 151(1):3–34, 2015.
- 720 Tetsuro Yamamoto. A convergence theorem for Newton-like methods in banach spaces. *Numerische*
721 *Mathematik*, 51:545–557, 1987.
- 722
- 723 Jiaqi Zhang, Keyou You, and Tamer Başar. Achieving globally superlinear convergence for dis-
724 tributed optimization with adaptive Newton method. In *2020 59th IEEE Conference on Decision*
725 *and Control (CDC)*, pp. 2329–2334. IEEE, 2020a.
- 726
- 727 Jingzhao Zhang, Sai Praneeth Karimireddy, Andreas Veit, Seungyeon Kim, Sashank Reddi, Sanjiv
728 Kumar, and Suvrit Sra. Why are adaptive methods good for attention models? *Advances in*
729 *Neural Information Processing Systems*, 33:15383–15393, 2020b.
- 730
- 731
- 732
- 733
- 734
- 735
- 736
- 737
- 738
- 739
- 740
- 741
- 742
- 743
- 744
- 745
- 746
- 747
- 748
- 749
- 750
- 751
- 752
- 753
- 754
- 755

756	CONTENTS	
757		
758	A Additional details	16
759		
760	A.1 Notations	16
761	A.2 Additional prior work	17
762		
763	B Basic facts	18
764		
765	C Properties of matrix smoothness	20
766		
767	C.1 The matrix Lipschitz-continuous gradient	20
768	C.1.1 Quadratics	20
769	C.2 Comparison of the different smoothness conditions	21
770	C.3 Proofs of the propositions regarding smoothness	21
771	C.3.1 Proof of Proposition 1	21
772	C.3.2 Proof of Proposition 2	22
773	C.3.3 Proof of Proposition 3	22
774	C.3.4 Proof of Proposition 4	23
775	C.3.5 Proof of Proposition 5	24
776	C.3.6 Proof of Proposition 6	24
777	C.3.7 Proof of Proposition 7	25
778	C.3.8 Proof of Proposition 8	25
779		
780		
781		
782		
783		
784		
785	D Analysis of <i>det-MARINA</i>	26
786		
787	D.1 Technical lemmas	26
788	D.2 Proof of Theorem 1	26
789	D.3 Comparison of different stepsizes	29
790	D.3.1 The diagonal case	30
791	D.3.2 The identity case	30
792	D.4 Proofs of the corollaries	30
793	D.4.1 Proof of Corollary 1	30
794	D.4.2 Proof of Corollary 3	31
795	D.4.3 Proof of Corollary 4	32
796	D.4.4 Proof of Corollary 8	32
797	D.4.5 Proof of Corollary 9	33
798		
799		
800		
801		
802		
803	E Analysis of <i>det-DASHA</i>	34
804		
805	E.1 Proof of Theorem 2	34
806	E.2 Proofs of the corollaries	36
807	E.2.1 Proof of Corollary 2	36
808	E.2.2 Proof of Corollary 5	37
809		

810	E.2.3	Proof of Corollary 6	37
811	E.2.4	Proof of Corollary 7	37
812			
813			
814	F	Distributed det-CGD	38
815			
816	G	Extension of det-CGD2 in MARINA form	38
817			
818	G.1	Extension of det-CGD2 to its variance reduced counterpart	38
819	G.2	Analysis of Algorithm 4	40
820	G.3	Proof of Theorem 4	40
821			
822			
823	H	Proofs of the technical lemmas	42
824			
825	H.1	Proof of Lemma 1	42
826	H.2	Proof of Lemma 2	43
827	H.3	Proof of Lemma 4	44
828	H.4	Proof of Lemma 5	45
829	H.5	Proof of Lemma 6	46
830	H.6	Proof of Lemma 7	46
831			
832			
833			
834	I	Experiments	47
835			
836	I.1	The setting	47
837	I.2	Comparison of all the methods	47
838	I.3	Improvements over MARINA	47
839	I.4	Improvements on non variance reduced methods	49
840	I.5	Improvements over det-CGD	51
841	I.6	Comparing different stepsize choices	51
842	I.7	Comparing communication complexity	53
843	I.8	Comparison of DASHA and det-DASHA	54
844	I.9	Comparison of DCGD, det-CGD, DASHA and det-DASHA	56
845	I.10	Comparison of det-DASHA and det-CGD with different stepsizes	56
846	I.11	Comparison of different stepsizes of det-DASHA	57
847	I.12	Comparison of det-MARINA and det-DASHA	57
848	I.13	Comparison in terms of function values	58
849			
850			
851			
852			
853			
854			
855	A	ADDITIONAL DETAILS	
856			
857	A.1	NOTATIONS	
858			
859		The standard Euclidean norm on \mathbb{R}^d is defined as $\ \cdot\ $. We use \mathbb{S}_{++}^d (resp. \mathbb{S}_+^d) to denote the positive	
860		definite (resp. semi-definite) cone of dimension d . \mathbb{S}^d is used to denote all symmetric matrices of	
861		dimension d . We use the notation \mathbf{I}_d to denote the identity matrix of size $d \times d$, and \mathbf{O}_d to denote	
862		the zero matrix of size $d \times d$. Given $\mathbf{Q} \in \mathbb{S}_{++}^d$ and $x \in \mathbb{R}^d$,	
863			

$$\|x\|_{\mathbf{Q}} := \sqrt{x^\top \mathbf{Q} x} = \sqrt{\langle x, \mathbf{Q} x \rangle},$$

864 where $\langle \cdot, \cdot \rangle$ is the standard Euclidean inner product on \mathbb{R}^d . For a matrix $\mathbf{A} \in \mathbb{S}^d$, we use $\lambda_{\max}(\mathbf{A})$
 865 (resp. $\lambda_{\min}(\mathbf{A})$) to denote the largest (resp. smallest) eigenvalue of the matrix \mathbf{A} . For a function
 866 $f : \mathbb{R}^d \mapsto \mathbb{R}$, its gradient and its Hessian at a point $x \in \mathbb{R}^d$ are respectively denoted as $\nabla f(x)$ and
 867 $\nabla^2 f(x)$. For the sketch matrices \mathbf{S}_i^k used in the algorithm, we use the superscript k to denote the
 868 iteration and subscript i to denote the client, the matrix \mathbf{S}_i^k is thus sampled for client i in the k -th
 869 iteration from the same distribution \mathcal{S} . For any matrix $\mathbf{A} \in \mathbb{S}^d$, we use the notation $\text{diag}(\mathbf{A}) \in \mathbb{S}^d$
 870 to denote the diagonal of matrix \mathbf{A} .

872 A.2 ADDITIONAL PRIOR WORK

873
 874 Numerous effective convex optimization techniques have been adapted for application in non-convex
 875 scenarios. Here’s a selection of these techniques, although it’s not an exhaustive list: adaptivity
 876 (Dvinskikh et al., 2019; Zhang et al., 2020b), variance reduction (J Reddi et al., 2016; Li et al.,
 877 2021), and acceleration (Guminov et al., 2019). Of particular relevance to our work is the paper by
 878 Khaled & Richtárik (2023), which introduces a unified approach for analyzing stochastic gradient
 879 descent for non-convex objectives. A comprehensive overview of non-convex optimization can be
 880 found in (Jain et al., 2017; Danilova et al., 2022).

881 An illustrative example of a matrix stepsized method is Newton’s method, which has been a long-
 882 standing favorite in the optimization community (Gragg & Tapia, 1974; Miel, 1980; Yamamoto,
 883 1987). However, the computational complexity involved in computing the stepsize as the inverse of
 884 the Hessian of the current iteration is substantial. Instead, quasi-Newton methods employ a readily
 885 computable estimator to replace the inverse Hessian (Broyden, 1965; Dennis & Moré, 1977; Al-
 886 Baali & Khalfan, 2007; Al-Baali et al., 2014). An important direction of research that is relevant to
 887 our work, studies distributed second order methods. Here is a non-exhaustive list of papers in this
 888 area: (Wang et al., 2018; Crane & Roosta, 2019; Zhang et al., 2020a; Islamov et al., 2021; Alimisis
 889 et al., 2021; Safaryan et al., 2022a).

890 The Distributed Compressed Gradient Descent (**DCGD**) algorithm, initially proposed by Khirirat
 891 et al. (2018), has seen improvements in various aspects, as documented in works such as (Li et al.,
 892 2020; Horváth et al., 2022). Its variance reduced version with gradients shifts was studied by Shulgin
 893 & Richtárik (2022) in the (strongly) convex setting. Additionally, there exists a substantial body of
 894 literature on other federated learning algorithms employing unbiased compressors (Alistarh et al.,
 895 2017; Mishchenko et al., 2019; Gorbunov et al., 2021; Mishchenko et al., 2022; Maranjyan et al.,
 896 2022; Horváth et al., 2023).

897 Variance reduction techniques have gained significant attention in the context of stochastic batch
 898 gradient descent that is prevalent in machine learning. Numerous algorithms have been developed
 899 in this regard, including well-known ones like **SVRG** (Johnson & Zhang, 2013), **SAG** (Schmidt
 900 et al., 2017), **SDCA** (Richtárik & Takáč, 2014), **SAGA** (Defazio et al., 2014), **MISO** (Mairal, 2015),
 901 and **Katyusha** (Allen-Zhu, 2017). An overview of more advanced methods can be found in (Gower
 902 et al., 2020). Notably, **SVRG** and **Katyusha** have been extended with loopless variants, namely
 903 **L-SVRG** and **L-Katyusha** (Kovalev et al., 2020; Qian et al., 2021). These loopless versions stream-
 904 line the algorithms by eliminating the outer loop and introducing a biased coin-flip mechanism at
 905 each step. This simplification eases both the algorithms’ structure and their analyses, while preserv-
 906 ing their worst-case complexity bounds. **L-SVRG**, in particular, offers the advantage of setting the
 907 exit probability from the outer loop independently of the condition number, thus, enhancing both
 robustness and practical efficiency.

908 This technique of coin flipping allows to obtain variance reduction for the **CGD** algorithm. A rele-
 909 vant example is the **DIANA** algorithm proposed by Mishchenko et al. (2019). Its convergence was
 910 proved both in the convex and non-convex cases. Later, **MARINA** (Gorbunov et al., 2021) obtained
 911 the optimal convergence rate, improving in communication complexity compared to all previous
 912 first order methods. Finally, there is a line of work developing variance reduction in the federated
 913 setting using other methods and techniques (Chraïbi et al., 2019; Hanzely & Richtárik, 2020; Dinh
 914 et al., 2020; Peng et al., 2022).

915 Another method to obtain variance reduction is based on momentum. It was initially studied by
 916 Cutkosky & Orabona (2019), where they propose the **STORM** algorithm, which is a stochastic gra-
 917 dient descent algorithm with a momentum term for non-convex objectives. They obtain stationarity
 guarantees using adaptive stepsizes with optimal convergence rates. However, they require the vari-

918 ance of the stochastic gradient to be bounded by a constant, which is impractical. Using momentum
 919 for variance reduction has since been widely studied (Liu et al., 2020; Khanduri et al., 2020; Tran-
 920 Dinh et al., 2022; Li et al., 2022).

922 B BASIC FACTS

923 In this section, we present some basic facts along with their proofs that will be used later in the
 924 analysis.

925 **Fact 1.** For two matrices $\mathbf{A}, \mathbf{B} \in \mathbb{S}_+^d$, denote the i -th largest eigenvalues of \mathbf{A}, \mathbf{B} as $\lambda_i(\mathbf{A}), \lambda_i(\mathbf{B})$,
 926 if $\mathbf{A} \succeq \mathbf{B}$, then the following holds

$$927 \lambda_i(\mathbf{A}) \geq \lambda_i(\mathbf{B}). \quad (8)$$

928 *Proof.* According to the Courant-Fischer theorem, we write

$$929 \lambda_i(\mathbf{B}) = \max_{S: \dim S=i} \min_{x \in S \setminus \{0\}} \frac{x^\top \mathbf{B} x}{x^\top x}.$$

930 Let S_{\max}^i be a subspace of dimension i where the maximum is attained, we then have

$$931 \lambda_i(\mathbf{B}) = \min_{x \in S_{\max}^i \setminus \{0\}} \frac{x^\top \mathbf{B} x}{x^\top x}$$

$$932 \leq \min_{x \in S_{\max}^i \setminus \{0\}} \frac{x^\top \mathbf{A} x}{x^\top x} \leq \max_{S: \dim S=i} \min_{x \in S \setminus \{0\}} \frac{x^\top \mathbf{A} x}{x^\top x} = \lambda_i(\mathbf{A}).$$

933 \square

934 The following is a generalization of the bias-variance decomposition for the matrix norm.

935 **Fact 2.** (Variance Decomposition) Given a matrix $\mathbf{M} \in \mathbb{S}_{++}^d$, any vector $c \in \mathbb{R}^d$, and a random
 936 vector $x \in \mathbb{R}^d$ such that $\mathbb{E}[\|x\|] \leq +\infty$, the following bound holds

$$937 \mathbb{E} \left[\|x - \mathbb{E}[x]\|_{\mathbf{M}}^2 \right] = \mathbb{E} \left[\|x - c\|_{\mathbf{M}}^2 \right] - \|\mathbb{E}[x] - c\|_{\mathbf{M}}^2. \quad (9)$$

938 *Proof.* We have

$$939 \mathbb{E} \left[\|x - c\|_{\mathbf{M}}^2 \right] - \|\mathbb{E}[x] - c\|_{\mathbf{M}}^2$$

$$940 = \mathbb{E} [x^\top \mathbf{M} x] - 2\mathbb{E}[x]^\top \mathbf{M} c + c^\top \mathbf{M} c - \mathbb{E}[x]^\top \mathbf{M} \mathbb{E}[x] + 2\mathbb{E}[x]^\top \mathbf{M} c - c^\top \mathbf{M} c$$

$$941 = \mathbb{E} [x^\top \mathbf{M} x] - \mathbb{E}[x]^\top \mathbf{M} \mathbb{E}[x]$$

$$942 = \mathbb{E} [x^\top \mathbf{M} x] - 2 \cdot \mathbb{E}[x]^\top \mathbf{M} \mathbb{E}[x] + \mathbb{E}[x]^\top \mathbf{M} \mathbb{E}[x]$$

$$943 = \mathbb{E} \left[\|x - \mathbb{E}[x]\|_{\mathbf{M}}^2 \right].$$

944 This completes the proof. \square

945 **Fact 3.** The map $(\mathbf{A}, \mathbf{B}, \mathbf{X}) \mapsto \mathbf{A} - \mathbf{X} \mathbf{B}^{-1} \mathbf{X}$ is jointly concave on $\mathbb{S}_+^d \times \mathbb{S}_{++}^d \times \mathbb{S}^d$. It is also
 946 monotone increasing in variables \mathbf{A} and \mathbf{B} .

947 We refer the reader to Corollary 1.5.3 of Bhatia (2009) for the details and the proof. The following
 948 is a result of Fact 1 and Fact 3.

949 **Fact 4.** Suppose $\mathbf{L}_i \in \mathbb{S}_{++}^d$, for $i = 1, \dots, n$. Then, for every matrix $\mathbf{X} \in \mathbb{S}_{++}^d$, we define the
 950 following mapping

$$951 f(\mathbf{X}, \mathbf{L}_1, \dots, \mathbf{L}_n) = \frac{1}{n} \sum_{i=1}^n \lambda_{\max}(\mathbf{L}_i) \cdot \lambda_{\max}(\mathbf{L}_i \mathbf{X}^{-1}) \cdot \lambda_{\max}(\mathbf{X}^{-1}).$$

952 Then the above mapping is monotone decreasing in \mathbf{X} .

972 *Proof.* First we notice that from Fact 3 the mapping $\mathbf{X} \mapsto \mathbf{X}^{-1}$ is monotone decreasing. The latter
 973 means that if we have any $\mathbf{X}_1, \mathbf{X}_2 \in \mathbb{S}_{++}^d$ such that $\mathbf{X}_1 \succeq \mathbf{X}_2$, we have

$$974 \mathbf{X}_1^{-1} \preceq \mathbf{X}_2^{-1}.$$

975 Then it immediately follows, due to Fact 1, that

$$976 0 < \lambda_{\max}(\mathbf{X}_1^{-1}) \leq \lambda_{\max}(\mathbf{X}_2^{-1}).$$

977 We also notice that the relation $\lambda_{\max}(\mathbf{L}_i \mathbf{X}^{-1}) = \lambda_{\max}\left(\mathbf{L}_i^{\frac{1}{2}} \mathbf{X}^{-1} \mathbf{L}_i^{\frac{1}{2}}\right) = \lambda_{\max}(\mathbf{X}^{-1} \mathbf{L}_i)$, and
 978 that the mapping $\mathbf{X} \mapsto \mathbf{L}_i^{\frac{1}{2}} \mathbf{X}^{-1} \mathbf{L}_i^{\frac{1}{2}}$ is also monotone decreasing for every $i \in [n]$, so we have

$$979 0 < \lambda_{\max}(\mathbf{L}_i \mathbf{X}_1^{-1}) \leq \lambda_{\max}(\mathbf{L}_i \mathbf{X}_2^{-1}).$$

980 Since we have the coefficient $\lambda_{\max}(\mathbf{L}_i) > 0$, it follows that,

$$981 f(\mathbf{X}_1, \mathbf{L}_1, \dots, \mathbf{L}_n) \leq f(\mathbf{X}_2, \mathbf{L}_1, \dots, \mathbf{L}_n).$$

982 This means that $f(\mathbf{X})$ is monotone decreasing in \mathbf{X} . □

983 **Fact 5.** For any two matrices $\mathbf{A}, \mathbf{B} \in \mathbb{S}_{++}^d$, the following relation regarding their largest eigenvalue
 984 holds

$$985 \lambda_{\max}(\mathbf{AB}) \leq \lambda_{\max}(\mathbf{A}) \cdot \lambda_{\max}(\mathbf{B}). \quad (10)$$

986 *Proof.* Using the Courant-Fischer theorem, we can write

$$\begin{aligned} 987 \lambda_{\max}(\mathbf{AB}) &= \min_{S: \dim S=d} \max_{x \in S \setminus \{0\}} \frac{x^\top \mathbf{AB}x}{x^\top x} \\ 988 &= \max_{x \in \mathbb{R}^d \setminus \{0\}} \frac{x^\top \mathbf{AB}x}{x^\top x} \\ 989 &\leq \max_{x \in \mathbb{R}^d \setminus \{0\}} \frac{x^\top \mathbf{A}x}{x^\top x} \cdot \max_{x \in \mathbb{R}^d \setminus \{0\}} \frac{x^\top \mathbf{B}x}{x^\top x} \\ 1000 &= \lambda_{\max}(\mathbf{A}) \cdot \lambda_{\max}(\mathbf{B}). \end{aligned}$$

1001 □

1002 **Fact 6.** Given matrix $\mathbf{Q} \in \mathbb{S}_{++}^d$ and its matrix norm $\|\cdot\|_{\mathbf{Q}}$, its associated dual norm is $\|\cdot\|_{\mathbf{Q}^{-1}}$.

1003 *Proof.* Let us first recall the definition of the dual norm $\|\cdot\|_*$. For any vector $z \in \mathbb{R}^d$, it is defined as

$$1004 \|z\|_* := \sup\{z^\top x : \|x\|_{\mathbf{Q}} \leq 1\}.$$

1005 Solving this optimization problem is equivalent to solving $\sup\{z^\top x : \|x\|_{\mathbf{Q}}^2 = 1\}$. The Lagrange
 1006 function is given as

$$1007 f(x, \lambda) = z^\top x - \lambda \left(\|x\|_{\mathbf{Q}}^2 - 1 \right) = z^\top x - \lambda (x^\top \mathbf{Q}x - 1).$$

1008 Computing the derivatives we deduce that

$$1009 \frac{\partial f(x, \lambda)}{\partial x} = z - 2\lambda \cdot \mathbf{Q}x = 0, \quad \frac{\partial f(x, \lambda)}{\partial \lambda} = \|x\|_{\mathbf{Q}}^2 - 1 = 0.$$

1010 This leads to

$$1011 \lambda = \frac{\|z\|_{\mathbf{Q}^{-1}}}{2}, \quad x = \frac{\mathbf{Q}^{-1}z}{\|z\|_{\mathbf{Q}^{-1}}}.$$

1012 As a result, we have

$$\begin{aligned} 1013 \sup\{z^\top x : \|x\|_{\mathbf{Q}} \leq 1\} &= \sup\{z^\top x : \|x\|_{\mathbf{Q}}^2 = 1\} \\ 1014 &= z^\top z = \frac{z^\top \mathbf{Q}^{-1}z}{\|z\|_{\mathbf{Q}^{-1}}} = \|z\|_{\mathbf{Q}^{-1}}. \end{aligned}$$

1015 □

C PROPERTIES OF MATRIX SMOOTHNESS

C.1 THE MATRIX LIPSCHITZ-CONTINUOUS GRADIENT

In this section we describe some properties of matrix smoothness, matrix gradient Lipschitzness and their relations. The following proposition describes a sufficient condition for the matrix Lipschitz-continuity of the gradient.

Proposition 2. *Given twice continuously differentiable function $f : \mathbb{R}^d \mapsto \mathbb{R}$ with bounded Hessian,*

$$\nabla^2 f(x) \preceq \mathbf{L}, \quad (11)$$

where $\mathbf{L} \in \mathbb{S}_{++}^d$ and the generalized inequality holds for any $x \in \mathbb{R}^d$. Then f satisfies (3) with the matrix \mathbf{L} .

The below proposition is a variant of Proposition 1 and it characterizes the smoothness matrix of the objective function f , given the smoothness matrices of the component functions f_i .

Proposition 3. *Assume that f_i has \mathbf{L}_i -Lipschitz continuous gradient for every $i \in [n]$, then function f has \mathbf{L} -Lipschitz gradient with $\mathbf{L} \in \mathbb{S}_{++}^d$ satisfying*

$$\mathbf{L} \cdot \lambda_{\min}(\mathbf{L}) = \frac{1}{n} \sum_{i=1}^n \lambda_{\max}(\mathbf{L}_i) \cdot \mathbf{L}_i. \quad (12)$$

C.1.1 QUADRATICS

Given a matrix $\mathbf{A} \in \mathbb{S}_{++}^d$ and a vector $b \in \mathbb{R}^d$, consider the function $f(x) = \frac{1}{2}x^\top \mathbf{A}x + b^\top x + c$. Then its gradient is computed as $\nabla f(x) = \mathbf{A}x + b$ and $\nabla^2 f(x) = \mathbf{A}$. Inserting gradients formula into (3) we deduce

$$\sqrt{(x-y)^\top \mathbf{A} \mathbf{L}^{-1} \mathbf{A} (x-y)} \leq \sqrt{(x-y)^\top \mathbf{L} (x-y)},$$

for any $x, y \in \mathbb{R}^d$. This reduces to

$$\mathbf{A} \mathbf{L}^{-1} \mathbf{A} \preceq \mathbf{L}. \quad (13)$$

Since $\mathbf{A} \in \mathbb{S}_{++}^d$, we can also rewrite (13) as

$$\mathbf{A}^{\frac{1}{2}} \mathbf{L}^{-1} \mathbf{A}^{\frac{1}{2}} \preceq \mathbf{A}^{-\frac{1}{2}} \mathbf{L} \mathbf{A}^{-\frac{1}{2}},$$

which is equivalent to

$$\mathbf{A} \preceq \mathbf{L}. \quad (14)$$

Therefore, the ‘‘best’’ $\mathbf{L} \in \mathbb{S}_{++}^d$ that satisfies (3) is $\mathbf{L} = \mathbf{A} = \nabla^2 f(x)$, for every $x \in \mathbb{R}^d$. Now, let us look at a more general setting. Consider f given as follows,

$$f(x) = \sum_{i=1}^s \phi_i(\mathbf{M}_i x),$$

where $\mathbf{M}_i \in \mathbb{R}^{q_i \times d}$. Here $f : \mathbb{R}^d \mapsto \mathbb{R}$ is the sum of functions $\phi_i : \mathbb{R}^{q_i} \mapsto \mathbb{R}$. We assume that each function ϕ_i has matrix \mathbf{L}_i Lipschitz gradient. We have the following lemma regarding the matrix gradient Lipschitzness of f .

Proposition 4. *Assume that functions f and $\{\phi_i\}_{i=1}^s$ are described above. Then function f has \mathbf{L} -Lipschitz gradient, if the following condition is satisfied:*

$$\sum_{i=1}^s \lambda_{\max} \left(\mathbf{L}_i^{\frac{1}{2}} \mathbf{M}_i \mathbf{L}^{-1} \mathbf{M}_i^\top \mathbf{L}_i^{\frac{1}{2}} \right) = 1. \quad (15)$$

Note that Proposition 4 is a generalization of the previous case of quadratics, if we pick $s = 1$, $\mathbf{M}_i = \mathbf{A}^{\frac{1}{2}}$ and $\phi_1(x) = x^\top \mathbf{I}_d x$, the condition becomes $\mathbf{L} = \mathbf{A}$, which is exactly the solution given by (14). Thus we recover the result for quadratics. The linear term $bx + c$ is ignored in this case. In Proposition 4, we only intend to give a way of finding a matrix $\mathbf{L} \in \mathbb{S}_{++}^d$, so that f has \mathbf{L} -Lipschitz gradient. This does not mean, however, the \mathbf{L} here is optimal. The proof is deferred to Appendix C.3.

C.2 COMPARISON OF THE DIFFERENT SMOOTHNESS CONDITIONS

Let us recall the definition of matrix smoothness.

Definition 3. (\mathbf{L} -smoothness) Assume that $f : \mathbb{R}^d \rightarrow \mathbb{R}$ is a continuously differentiable function and matrix $\mathbf{L} \in \mathbb{S}_{++}^d$. We say that f is \mathbf{L} -smooth if for all $x, y \in \mathbb{R}^d$

$$f(y) \leq f(x) + \langle \nabla f(x), x - y \rangle + \frac{1}{2} \|x - y\|_{\mathbf{L}}^2. \quad (16)$$

We provide a proposition here which describes an equivalent form of stating \mathbf{L} -matrix smoothness of a function f . This proposition is used to illustrate the relation between matrix smoothness and matrix Lipschitz gradient.

Proposition 5. Let function $f : \mathbb{R}^d \rightarrow \mathbb{R}$ be continuously differentiable. Then the following statements are equivalent.

- (i) f is \mathbf{L} -matrix smooth.
- (ii) $\langle \nabla f(x) - \nabla f(y), x - y \rangle \leq \|x - y\|_{\mathbf{L}}^2$ for all $x, y \in \mathbb{R}^d$.

The two propositions, Proposition 6 and Proposition 7, formulated below illustrate the relation between matrix smoothness of f and matrix gradient Lipschitzness of f .

Proposition 6. Assume $f : \mathbb{R}^d \mapsto \mathbb{R}$ is a continuously differentiable function, and its gradient is \mathbf{L} -Lipschitz continuous with $\mathbf{L} \in \mathbb{S}_{++}^d$. Then function f is \mathbf{L} -matrix smooth.

Proposition 7. Assume $f : \mathbb{R}^d \rightarrow \mathbb{R}$ is a continuously differentiable function. Assume also that f is convex and \mathbf{L} -matrix smooth. Then ∇f is \mathbf{L} -Lipschitz continuous.

The next proposition shows that standard Lipschitzness of the gradient of a function is an immediate consequence of matrix Lipschitzness.

Proposition 8. Assume that the gradient of f is \mathbf{L} -Lipschitz continuous. Then ∇f is also \mathbf{L} -Lipschitz with $L = \lambda_{\max}(\mathbf{L})$.

C.3 PROOFS OF THE PROPOSITIONS REGARDING SMOOTHNESS

C.3.1 PROOF OF PROPOSITION 1

We start with the definition of \mathbf{L} -Lipschitz gradient of function f , and pick two arbitrary points $x, y \in \mathbb{R}^d$,

$$\|\nabla f(x) - \nabla f(y)\|_{\mathbf{L}^{-1}}^2 = \left\| \frac{1}{n} \sum_{i=1}^n (\nabla f_i(x) - \nabla f_i(y)) \right\|_{\mathbf{L}^{-1}}^2.$$

Applying the convexity of $\|\cdot\|_{\mathbf{L}^{-1}}^2$, we have

$$\|\nabla f(x) - \nabla f(y)\|_{\mathbf{L}^{-1}}^2 \leq \frac{1}{n} \sum_{i=1}^n \|\nabla f_i(x) - \nabla f_i(y)\|_{\mathbf{L}^{-1}}^2.$$

For each term within the summation, we use the definition of matrix norms and replace the matrix \mathbf{L}^{-1} with $\mathbf{L}_i^{-1/2} \mathbf{L}_i^{1/2} \mathbf{L}^{-1} \mathbf{L}_i^{1/2} \mathbf{L}_i^{-1/2}$, for every $i = 1, \dots, n$:

$$\begin{aligned} \|\nabla f(x) - \nabla f(y)\|_{\mathbf{L}^{-1}}^2 &= \frac{1}{n} \sum_{i=1}^n \left(\mathbf{L}_i^{-\frac{1}{2}} (\nabla f_i(x) - \nabla f_i(y)) \right)^\top \mathbf{L}_i^{\frac{1}{2}} \mathbf{L}^{-1} \mathbf{L}_i^{\frac{1}{2}} \left(\mathbf{L}_i^{-\frac{1}{2}} (\nabla f_i(x) - \nabla f_i(y)) \right) \\ &\leq \frac{1}{n} \sum_{i=1}^n \lambda_{\max} \left(\mathbf{L}_i^{\frac{1}{2}} \mathbf{L}^{-1} \mathbf{L}_i^{\frac{1}{2}} \right) \left\| \mathbf{L}_i^{-\frac{1}{2}} (\nabla f_i(x) - \nabla f_i(y)) \right\|^2 \\ &= \frac{1}{n} \sum_{i=1}^n \lambda_{\max} \left(\mathbf{L}_i^{\frac{1}{2}} \mathbf{L}^{-1} \mathbf{L}_i^{\frac{1}{2}} \right) \|\nabla f_i(x) - \nabla f_i(y)\|_{\mathbf{L}_i^{-1}}^2. \end{aligned}$$

Using the assumption that the gradient of each function f_i is \mathbf{L}_i -Lipschitz, we obtain,

$$\|\nabla f(x) - \nabla f(y)\|_{\mathbf{L}^{-1}}^2 \leq \frac{1}{n} \sum_{i=1}^n \lambda_{\max} \left(\mathbf{L}_i^{\frac{1}{2}} \mathbf{L}^{-1} \mathbf{L}_i^{\frac{1}{2}} \right) \|x - y\|_{\mathbf{L}_i}^2.$$

Replacing \mathbf{L}_i^{-1} with $\mathbf{L}^{-1/2} \mathbf{L}^{1/2} \mathbf{L}_i^{-1} \mathbf{L}^{1/2} \mathbf{L}^{-1/2}$

$$\begin{aligned} \|\nabla f(x) - \nabla f(y)\|_{\mathbf{L}^{-1}}^2 &= \frac{1}{n} \sum_{i=1}^n \lambda_{\max} \left(\mathbf{L}_i^{\frac{1}{2}} \mathbf{L}^{-1} \mathbf{L}_i^{\frac{1}{2}} \right) \cdot \left[(\mathbf{L}^{\frac{1}{2}}(x - y))^\top \mathbf{L}^{-\frac{1}{2}} \mathbf{L}_i \mathbf{L}^{-\frac{1}{2}} (\mathbf{L}^{\frac{1}{2}}(x - y)) \right] \\ &\leq \frac{1}{n} \sum_{i=1}^n \lambda_{\max} \left(\mathbf{L}_i^{\frac{1}{2}} \mathbf{L}^{-1} \mathbf{L}_i^{\frac{1}{2}} \right) \cdot \lambda_{\max} \left(\mathbf{L}^{-\frac{1}{2}} \mathbf{L}_i \mathbf{L}^{-\frac{1}{2}} \right) \left\| \mathbf{L}^{\frac{1}{2}}(x - y) \right\|^2. \end{aligned}$$

Using Fact 5, we are deduce the following bound,

$$\begin{aligned} \|\nabla f(x) - \nabla f(y)\|_{\mathbf{L}^{-1}}^2 &\leq \left(\frac{1}{n} \sum_{i=1}^n \lambda_{\max}(\mathbf{L}^{-1}) \cdot \lambda_{\max}(\mathbf{L}_i) \cdot \lambda_{\max}(\mathbf{L}_i \mathbf{L}^{-1}) \right) \cdot \|x - y\|_{\mathbf{L}}^2 \\ &= \|x - y\|_{\mathbf{L}}^2. \end{aligned}$$

C.3.2 PROOF OF PROPOSITION 2

We start with picking any two vector $x, y \in \mathbb{R}^d$. We have

$$\begin{aligned} \|\nabla f(x) - \nabla f(y)\|_{\mathbf{L}^{-1}}^2 &= \left\| \int_0^1 \nabla^2 f(\theta x + (1 - \theta)y)(x - y) d\theta \right\|_{\mathbf{L}^{-1}}^2 \\ &= (x - y)^\top \left(\int_0^1 \nabla^2 f(\theta x + (1 - \theta)y) d\theta \right)^\top \mathbf{L}^{-1} \left(\int_0^1 \nabla^2 f(\theta x + (1 - \theta)y) d\theta \right) (x - y). \end{aligned}$$

Denote $\mathbf{F} := \int_0^1 \nabla^2 f(\theta x + (1 - \theta)y) d\theta$, notice that \mathbf{F} is a symmetric matrix. Then, the previous identity becomes

$$\|\nabla f(x) - \nabla f(y)\|_{\mathbf{L}^{-1}}^2 = (x - y)^\top \mathbf{F}^\top \mathbf{L}^{-1} \mathbf{F} (x - y).$$

From the definition of \mathbf{F} and the bounded Hessian assumption, we have $\mathbf{F} \preceq \mathbf{L}$. Let us prove that $\mathbf{F} \mathbf{L}^{-1} \mathbf{F} \preceq \mathbf{L}$:

$$\begin{aligned} \mathbf{F} \mathbf{L}^{-1} \mathbf{F} \preceq \mathbf{L} &\iff \mathbf{L}^{-\frac{1}{2}} \mathbf{F} \mathbf{L} \mathbf{F} \mathbf{L}^{-\frac{1}{2}} \preceq \mathbf{I}_d \\ &\iff \mathbf{L}^{-\frac{1}{2}} \mathbf{F} \mathbf{L}^{-\frac{1}{2}} \cdot \mathbf{L}^{-\frac{1}{2}} \mathbf{F} \mathbf{L}^{-\frac{1}{2}} \preceq \mathbf{I}_d \\ &\iff \mathbf{L}^{-\frac{1}{2}} \mathbf{F} \mathbf{L}^{-\frac{1}{2}} \preceq \mathbf{I}_d \\ &\iff \mathbf{F} \preceq \mathbf{L}. \end{aligned}$$

This means that

$$\|\nabla f(x) - \nabla f(y)\|_{\mathbf{L}^{-1}}^2 \leq (x - y)^\top \mathbf{L} (x - y) = \|x - y\|_{\mathbf{L}}^2,$$

which completes the proof.

C.3.3 PROOF OF PROPOSITION 3

Suppose \mathbf{L} is a symmetric positive definite matrix satisfying (12). Let us now show that the function ∇f is \mathbf{L} -Lipschitz continuous. We start with picking any two points $x, y \in \mathbb{R}^d$, and notice that

$$\|\nabla f(x) - \nabla f(y)\|_{\mathbf{L}^{-1}}^2 = \left\| \frac{1}{n} \sum_{i=1}^n (\nabla f_i(x) - \nabla f_i(y)) \right\|_{\mathbf{L}^{-1}}^2.$$

Applying Jensen's inequality, we obtain

$$\|\nabla f(x) - \nabla f(y)\|_{\mathbf{L}^{-1}}^2 \leq \frac{1}{n} \sum_{i=1}^n \|\nabla f_i(x) - \nabla f_i(y)\|_{\mathbf{L}^{-1}}^2.$$

We then re-weight the norm appears in the summation individually,

$$\begin{aligned} \|\nabla f(x) - \nabla f(y)\|_{\mathbf{L}^{-1}}^2 &\leq \frac{1}{n} \sum_{i=1}^n (\nabla f_i(x) - \nabla f_i(y))^\top \mathbf{L}_i^{-\frac{1}{2}} \mathbf{L}_i^{\frac{1}{2}} \mathbf{L}^{-1} \mathbf{L}_i^{\frac{1}{2}} \mathbf{L}_i^{-\frac{1}{2}} (\nabla f_i(x) - \nabla f_i(y)) \\ &\leq \frac{1}{n} \sum_{i=1}^n \lambda_{\max}(\mathbf{L}_i) \cdot \lambda_{\max}(\mathbf{L}^{-1}) \cdot \|\nabla f_i(x) - \nabla f_i(y)\|_{\mathbf{L}_i^{-1}}^2. \end{aligned}$$

Utilizing the assumption that each f_i has \mathbf{L}_i Lipschitz gradient, we obtain

$$\begin{aligned} \|\nabla f(x) - \nabla f(y)\|_{\mathbf{L}^{-1}}^2 &\leq \frac{1}{n} \sum_{i=1}^n \lambda_{\max}(\mathbf{L}_i) \cdot \lambda_{\max}(\mathbf{L}^{-1}) \cdot \|x - y\|_{\mathbf{L}_i}^2 \\ &= \|x - y\|_{\lambda_{\max}(\mathbf{L}^{-1}) \cdot \frac{1}{n} \sum_{i=1}^n \lambda_{\max}(\mathbf{L}_i) \cdot \mathbf{L}_i}^2 \stackrel{(12)}{=} \|x - y\|_{\mathbf{L}}^2. \end{aligned}$$

C.3.4 PROOF OF PROPOSITION 4

For any x and y from \mathbb{R}^d , we have

$$\begin{aligned} &\|\nabla f(x) - \nabla f(y)\|_{\mathbf{L}^{-1}} \\ &= \left\| \sum_{i=1}^s \mathbf{M}_i^\top \nabla \phi_i(\mathbf{M}_i x) - \sum_{i=1}^s \mathbf{M}_i^\top \nabla \phi_i(\mathbf{M}_i y) \right\|_{\mathbf{L}^{-1}} \\ &= s \cdot \left\| \frac{1}{s} \sum_{i=1}^s \mathbf{M}_i^\top (\nabla \phi_i(\mathbf{M}_i x) - \nabla \phi_i(\mathbf{M}_i y)) \right\|_{\mathbf{L}^{-1}}. \end{aligned}$$

Applying the convexity of the norm $\|\cdot\|_{\mathbf{L}^{-1}}$,

$$\|\nabla f(x) - \nabla f(y)\|_{\mathbf{L}^{-1}} \leq s \cdot \frac{1}{s} \sum_{i=1}^s \|\mathbf{M}_i^\top (\nabla \phi_i(\mathbf{M}_i x) - \nabla \phi_i(\mathbf{M}_i y))\|_{\mathbf{L}^{-1}}.$$

Expanding the norm and applying the replacement trick for above \mathbf{L} and \mathbf{M}_i , we obtain

$$\begin{aligned} &\|\nabla f(x) - \nabla f(y)\|_{\mathbf{L}^{-1}} \\ &= \sum_{i=1}^s \sqrt{(\nabla \phi_i(\mathbf{M}_i x) - \nabla \phi_i(\mathbf{M}_i y))^\top \mathbf{M}_i \mathbf{L}^{-1} \mathbf{M}_i^\top (\nabla \phi_i(\mathbf{M}_i x) - \nabla \phi_i(\mathbf{M}_i y))} \\ &= \sum_{i=1}^s \sqrt{\mathbf{B}_i^\top \mathbf{L}_i^{\frac{1}{2}} \mathbf{M}_i \mathbf{L}^{-1} \mathbf{M}_i^\top \mathbf{L}_i^{\frac{1}{2}} \mathbf{B}_i} \\ &\leq \sum_{i=1}^s \sqrt{\lambda_{\max}(\mathbf{L}_i^{\frac{1}{2}} \mathbf{M}_i \mathbf{L}^{-1} \mathbf{M}_i^\top \mathbf{L}_i^{\frac{1}{2}})} \cdot \|\nabla \phi_i(\mathbf{M}_i x) - \nabla \phi_i(\mathbf{M}_i y)\|_{\mathbf{L}_i^{-1}}, \end{aligned}$$

where $\mathbf{B}_i := \mathbf{L}_i^{-\frac{1}{2}} (\nabla \phi_i(\mathbf{M}_i x) - \nabla \phi_i(\mathbf{M}_i y))$. Due to the assumption that the gradient of ϕ_i is \mathbf{L}_i -Lipschitz, we have

$$\begin{aligned} &\|\nabla f(x) - \nabla f(y)\|_{\mathbf{L}^{-1}} \\ &\leq \sum_{i=1}^s \sqrt{\lambda_{\max}(\mathbf{L}_i^{\frac{1}{2}} \mathbf{M}_i \mathbf{L}^{-1} \mathbf{M}_i^\top \mathbf{L}_i^{\frac{1}{2}})} \cdot \|\mathbf{M}_i(x - y)\|_{\mathbf{L}_i} \\ &= \sum_{i=1}^s \sqrt{\lambda_{\max}(\mathbf{L}_i^{\frac{1}{2}} \mathbf{M}_i \mathbf{L}^{-1} \mathbf{M}_i^\top \mathbf{L}_i^{\frac{1}{2}})} \cdot \sqrt{\left[\mathbf{L}^{\frac{1}{2}}(x - y)\right]^\top \mathbf{L}^{-\frac{1}{2}} \mathbf{M}_i^\top \mathbf{L}_i \mathbf{M}_i \mathbf{L}^{-\frac{1}{2}} \left[\mathbf{L}^{\frac{1}{2}}(x - y)\right]} \\ &\leq \sum_{i=1}^s \sqrt{\lambda_{\max}(\mathbf{L}_i^{\frac{1}{2}} \mathbf{M}_i \mathbf{L}^{-1} \mathbf{M}_i^\top \mathbf{L}_i^{\frac{1}{2}})} \cdot \lambda_{\max}(\mathbf{L}^{-\frac{1}{2}} \mathbf{M}_i^\top \mathbf{L}_i \mathbf{M}_i \mathbf{L}^{-\frac{1}{2}}) \cdot \|x - y\|_{\mathbf{L}} \\ &\leq \sum_{i=1}^s \lambda_{\max}(\mathbf{L}_i^{\frac{1}{2}} \mathbf{M}_i \mathbf{L}^{-1} \mathbf{M}_i^\top \mathbf{L}_i^{\frac{1}{2}}) \cdot \|x - y\|_{\mathbf{L}}, \end{aligned}$$

where the last inequality is due to the fact that,

$$\lambda_{\max} \left(\mathbf{L}_i^{\frac{1}{2}} \mathbf{M}_i \mathbf{L}^{-1} \mathbf{M}_i^\top \mathbf{L}_i^{\frac{1}{2}} \right) = \lambda_{\max} \left(\mathbf{L}^{-\frac{1}{2}} \mathbf{M}_i^\top \mathbf{L}_i \mathbf{M}_i \mathbf{L}^{-\frac{1}{2}} \right).$$

Recalling the condition of the proposition:

$$\sum_{i=1}^s \lambda_{\max} \left(\mathbf{L}_i^{\frac{1}{2}} \mathbf{M}_i \mathbf{L}^{-1} \mathbf{M}_i^\top \mathbf{L}_i^{\frac{1}{2}} \right) = 1,$$

we deduce

$$\|\nabla f(x) - \nabla f(y)\|_{\mathbf{L}^{-1}} \leq \|x - y\|_{\mathbf{L}}.$$

C.3.5 PROOF OF PROPOSITION 5

(i) \rightarrow (ii). If f is \mathbf{L} -matrix smooth, then for all $x, y \in \mathbb{R}^d$, we have

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

and

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

Summing up these two inequalities we get

$$\langle \nabla f(x) - \nabla f(y), x - y \rangle \leq \|x - y\|_{\mathbf{L}}^2.$$

(ii) \rightarrow (i). Choose any $x, y \in \mathbb{R}^d$, and define $z = x + t(y - x)$, then we have,

$$\begin{aligned} f(y) &= f(x) + \int_0^1 \langle \nabla f(x + t(y - x)), y - x \rangle dt \\ &= f(x) + \int_0^1 \langle \nabla f(z), y - x \rangle dt \\ &= f(x) + \langle \nabla f(x), y - x \rangle + \int_0^1 \langle \nabla f(z) - \nabla f(x), y - x \rangle dt \\ &= f(x) + \langle \nabla f(x), y - x \rangle + \int_0^1 \langle \nabla f(z) - \nabla f(x), z - x \rangle \cdot \frac{1}{t} dt. \end{aligned}$$

Using the assumption that for any $x, z \in \mathbb{R}^d$, we have

$$\langle \nabla f(z) - \nabla f(x), z - x \rangle \leq \|z - x\|_{\mathbf{L}}^2.$$

Plug this back into the previous identity, we obtain

$$\begin{aligned} f(y) &\leq f(x) + \langle \nabla f(x), y - x \rangle + \int_0^1 \|z - x\|_{\mathbf{L}}^2 \cdot \frac{1}{t} dt \\ &= f(x) + \langle \nabla f(x), y - x \rangle + \int_0^1 \|y - x\|_{\mathbf{L}}^2 \cdot t dt \\ &= f(x) + \langle \nabla f(x), y - x \rangle + \frac{1}{2} \|y - x\|_{\mathbf{L}}^2. \end{aligned}$$

C.3.6 PROOF OF PROPOSITION 6

We start with picking any two points $x, y \in \mathbb{R}^d$, using the generalized Cauchy-Schwarz inequality for dual norm, we have

$$\begin{aligned} \langle \nabla f(x) - \nabla f(y), x - y \rangle &\leq \|\nabla f(x) - \nabla f(y)\|_{\mathbf{L}^{-1}} \cdot \|x - y\|_{\mathbf{L}} \\ &\stackrel{(3)}{\leq} \|x - y\|_{\mathbf{L}} \cdot \|x - y\|_{\mathbf{L}} \\ &= \|x - y\|_{\mathbf{L}}^2 \end{aligned}$$

According to Proposition 5, this indicates that function f is \mathbf{L} -matrix smooth.

1296 C.3.7 PROOF OF PROPOSITION 7

1297 Using Proposition 5, we know that for any $x, y \in \mathbb{R}^d$, we have

$$1299 \quad \langle \nabla f(x) - \nabla f(y), x - y \rangle \leq \|x - y\|_{\mathbf{L}}^2. \quad (17)$$

1300 Now we pick any three points $x, y, z \in \mathbb{R}^d$. Using the \mathbf{L} -smoothness of f , we have

$$1302 \quad f(x + z) \geq f(x) + \langle \nabla f(x), z \rangle + \frac{1}{2} \|z\|_{\mathbf{L}}^2. \quad (18)$$

1303 Using the convexity of f we have

$$1304 \quad \langle \nabla f(y), x + z - y \rangle \leq f(x + z) - f(y). \quad (19)$$

1305 Combining (18) and (19), we obtain

$$1307 \quad \langle \nabla f(y), x + z - y \rangle \leq f(x) - f(y) + \langle \nabla f(x), z \rangle + \frac{1}{2} \|z\|_{\mathbf{L}}^2.$$

1308 Rearranging terms we get

$$1310 \quad \langle \nabla f(y) - \nabla f(x), z \rangle - \frac{1}{2} \|z\|_{\mathbf{L}}^2 \leq f(x) - f(y) - \langle \nabla f(y), x - y \rangle.$$

1311 The inequality holds for any z for fixed x and y , and the left hand side is maximized (w.r.t. z) when $z = \mathbf{L}^{-1} (\nabla f(y) - \nabla f(x))$. Plugging it in, we get

$$1314 \quad \frac{1}{2} \|\nabla f(x) - \nabla f(y)\|_{\mathbf{L}^{-1}}^2 \leq f(x) - f(y) - \langle \nabla f(y), x - y \rangle. \quad (20)$$

1315 By symmetry we can also obtain

$$1317 \quad \frac{1}{2} \|\nabla f(y) - \nabla f(x)\|_{\mathbf{L}^{-1}}^2 \leq f(y) - f(x) - \langle \nabla f(x), y - x \rangle.$$

1318 Adding (20) and its counterpart together, we get

$$1319 \quad \|\nabla f(x) - \nabla f(y)\|_{\mathbf{L}^{-1}}^2 \leq \langle \nabla f(x) - \nabla f(y), x - y \rangle. \quad (21)$$

1320 Combing (21) and (17), it follows

$$1322 \quad \|\nabla f(x) - \nabla f(y)\|_{\mathbf{L}^{-1}}^2 \leq \|x - y\|_{\mathbf{L}}^2.$$

1323 Note that \mathbf{L} and \mathbf{L}^{-1} are both positive definite matrices, so it is equivalent to

$$1324 \quad \|\nabla f(x) - \nabla f(y)\|_{\mathbf{L}^{-1}} \leq \|x - y\|_{\mathbf{L}}.$$

1325 This completes the proof.

1326 C.3.8 PROOF OF PROPOSITION 8

1327 Let us start with picking any two points $x, y \in \mathbb{R}^d$. With the matrix \mathbf{L} -Lipschitzness of the gradient of function f , we have

$$1329 \quad \|\nabla f(x) - \nabla f(y)\|_{\mathbf{L}^{-1}}^2 \leq \|x - y\|_{\mathbf{L}}^2.$$

1330 This implies

$$1331 \quad (x - y)^\top \mathbf{L}(x - y) - (\nabla f(x) - \nabla f(y))^\top \mathbf{L}^{-1} (\nabla f(x) - \nabla f(y)) \geq 0.$$

1332 Define function $f(\mathbf{X}) := a^\top \mathbf{X}a - b^\top \mathbf{X}^{-1}b$ for $\mathbf{X} \in \mathbb{S}_{++}^d$, where $a, b \in \mathbb{R}^d$ are fixed vectors. Then f is monotone increasing in \mathbf{X} . This can be shown in the following way, picking two matrices $\mathbf{X}_1, \mathbf{X}_2 \in \mathbb{S}_{++}^d$, where $\mathbf{X}_1 \succeq \mathbf{X}_2$. It is easy to see that $-\mathbf{X}_1^{-1} \succeq -\mathbf{X}_2^{-1}$, since from Fact 3 the map $\mathbf{X} \mapsto -\mathbf{X}^{-1}$ is monotone increasing for $\mathbf{X} \in \mathbb{S}_{++}^d$. Thus,

$$1333 \quad f(\mathbf{X}_1) - f(\mathbf{X}_2) = (x - y)^\top (\mathbf{X}_1 - \mathbf{X}_2)(x - y) \\ 1334 \quad + (\nabla f(x) - \nabla f(y))^\top (-\mathbf{X}_1^{-1} - (-\mathbf{X}_2^{-1})) (\nabla f(x) - \nabla f(y)) \geq 0.$$

1335 As a result, $f(\lambda_{\max}(\mathbf{L}) \cdot \mathbf{I}_d) \geq f(\mathbf{L}) \geq 0$, due to the fact that $\lambda_{\max}(\mathbf{L}) \cdot \mathbf{I}_d \succeq \mathbf{L}$. It remains to notice that

$$1336 \quad f(\lambda_{\max}(\mathbf{L}) \cdot \mathbf{I}_d) = \lambda_{\max}(\mathbf{L}) \|x - y\|^2 - \frac{1}{\lambda_{\max}(\mathbf{L})} \|\nabla f(x) - \nabla f(y)\|^2 \geq 0,$$

1337 which yields

$$1338 \quad \|\nabla f(x) - \nabla f(y)\|^2 \leq \lambda_{\max}^2(\mathbf{L}) \|x - y\|^2.$$

1339 Since we are working with $\mathbf{L} \in \mathbb{S}_{++}^d$, the above inequality implies

$$1340 \quad \|\nabla f(x) - \nabla f(y)\| \leq \lambda_{\max}(\mathbf{L}) \|x - y\|.$$

D ANALYSIS OF DET-MARINA

D.1 TECHNICAL LEMMAS

We first state some technical lemmas.

Lemma 1 (Descent lemma). *Assume that function f is L smooth, and $x^{k+1} = x^k - \mathbf{D} \cdot g^k$, where $\mathbf{D} \in \mathbb{S}_{++}^d$. Then we will have*

$$f(x^{k+1}) \leq f(x^k) - \frac{1}{2} \|\nabla f(x^k)\|_{\mathbf{D}}^2 + \frac{1}{2} \|g^k - \nabla f(x^k)\|_{\mathbf{D}}^2 - \frac{1}{2} \|x^{k+1} - x^k\|_{\mathbf{D}^{-1}-L}.$$

The following lemma is obtained for any sketch matrix $\mathbf{S} \in \mathbb{S}_+^d$ and any two positive definite matrices \mathbf{D} and \mathbf{L} .

Lemma 2 (Property of sketch matrix). *For any sketch matrix $\mathbf{S} \in \mathbb{S}_+^d$, a vector $t \in \mathbb{R}^d$, and matrices $\mathbf{D}, \mathbf{L} \in \mathbb{S}_{++}^d$, we have*

$$\mathbb{E} \left[\| \mathbf{S}t - t \|_{\mathbf{D}}^2 \right] \leq \lambda_{\max} \left(\mathbf{L}^{\frac{1}{2}} (\mathbb{E} [\mathbf{S} \mathbf{D} \mathbf{S}] - \mathbf{D}) \mathbf{L}^{\frac{1}{2}} \right) \cdot \|t\|_{\mathbf{L}^{-1}}^2. \quad (22)$$

Lemma 3. *Assume that Definition 1 holds and $h_i^0 = \nabla f_i(x^0)$, then for h_i^{k+1} from Algorithm 2, we have for any $\mathbf{D} \in \mathbb{S}_{++}^d$*

$$\|h^{k+1} - \nabla f(x^{k+1})\|_{\mathbf{D}}^2 = \|h_i^{k+1} - \nabla f_i(x^{k+1})\|_{\mathbf{D}}^2 = 0,$$

and

$$\|h_i^{k+1} - h_i^k\|_{\mathbf{L}_i^{-1}}^2 \leq \|x^{k+1} - x^k\|_{\mathbf{L}_i}^2.$$

The following lemmas describe the recurrence applied to terms in the Lyapunov function.

Lemma 4. *Suppose h^{k+1} and g^{k+1} are from Algorithm 2, then the following recurrence relation holds,*

$$\begin{aligned} & \mathbb{E} \left[\|g^{k+1} - h^{k+1}\|_{\mathbf{D}}^2 \right] \\ & \leq \frac{2\Lambda_{\mathbf{D},\mathbf{S}} \cdot \lambda_{\max}(\mathbf{D}^{-1}) \cdot \lambda_{\max}(\mathbf{D})}{n^2} \sum_{i=1}^n \lambda_{\max}(\mathbf{L}_i) \mathbb{E} \left[\|h_i^{k+1} - h_i^k\|_{\mathbf{L}_i^{-1}}^2 \right] \\ & \quad + \frac{2a^2 \Lambda_{\mathbf{D},\mathbf{S}} \cdot \lambda_{\max}(\mathbf{D}^{-1})}{n^2} \sum_{i=1}^n \mathbb{E} \left[\|g_i^k - h_i^k\|_{\mathbf{D}}^2 \right] + (1-a)^2 \mathbb{E} \left[\|g^k - h^k\|_{\mathbf{D}}^2 \right], \end{aligned} \quad (23)$$

where $\Lambda_{\mathbf{D},\mathbf{S}} = \lambda_{\max}(\mathbb{E}[\mathbf{S}_i^k \mathbf{D} \mathbf{S}_i^k] - \mathbf{D})$ for $\mathbf{D} \in \mathbb{S}_{++}^d$ and $\mathbf{S}_i^k \sim \mathbf{S}$.

Lemma 5. *Suppose h_i^{k+1} and g_i^{k+1} for $i \in [n]$ are from Algorithm 2, then the following recurrence holds,*

$$\begin{aligned} & \mathbb{E} \left[\|g_i^{k+1} - h_i^{k+1}\|_{\mathbf{D}}^2 \right] \\ & \leq (2a^2 \lambda_{\max}(\mathbf{D}^{-1}) \cdot \Lambda_{\mathbf{D},\mathbf{S}} + (1-a)^2) \cdot \mathbb{E} \left[\|g_i^k - h_i^k\|_{\mathbf{D}}^2 \right] \\ & \quad + 2\lambda_{\max}(\mathbf{D}^{-1}) \cdot \lambda_{\max}(\mathbf{D}) \cdot \Lambda_{\mathbf{D},\mathbf{S}} \cdot \lambda_{\max}(\mathbf{L}_i) \cdot \mathbb{E} \left[\|h_i^{k+1} - h_i^k\|_{\mathbf{L}_i^{-1}}^2 \right]. \end{aligned}$$

D.2 PROOF OF THEOREM 1

According to Lemma 1, we have

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

We then use the definition of g^{k+1} to derive an upper bound for $\mathbb{E} \left[\|g^{k+1} - \nabla f(x^{k+1})\|_{\mathcal{D}}^2 \right]$. Notice that,

$$g^{k+1} = \begin{cases} \nabla f(x^{k+1}) & \text{with probability } p, \\ g^k + \frac{1}{n} \sum_{i=1}^n \mathbf{S}_i^k (\nabla f_i(x^{k+1}) - \nabla f_i(x^k)) & \text{with probability } 1 - p. \end{cases}$$

As a result, from the tower property,

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

Using Fact 2, we have

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

Notice that the sketch matrix is unbiased, thus we have

$$\mathbb{E} [\mathbf{S}_i^k (\nabla f_i(x^{k+1}) - \nabla f_i(x^k)) \mid x^{k+1}, x^k] = \nabla f_i(x^{k+1}) - \nabla f_i(x^k),$$

and any two random vectors in the set $\{\mathbf{S}_i^k (\nabla f_i(x^{k+1}) - \nabla f_i(x^k))\}_{i=1}^n$ are independent from each other, if x^{k+1} and x^k are fixed. Therefore, we have

$$\begin{aligned} & \mathbb{E} \left[\|g^{k+1} - \nabla f(x^{k+1})\|_{\mathcal{D}}^2 \mid x^{k+1}, x^k \right] \\ &= \frac{1-p}{n^2} \sum_{i=1}^n \mathbb{E} \left[\left\| \mathbf{S}_i^k (\nabla f_i(x^{k+1}) - \nabla f_i(x^k)) - (\nabla f_i(x^{k+1}) - \nabla f_i(x^k)) \right\|_{\mathcal{D}}^2 \mid x^{k+1}, x^k \right] \\ &\quad + (1-p) \cdot \|g^k - \nabla f(x^k)\|_{\mathcal{D}}^2. \end{aligned} \tag{25}$$

Lemma 2 yields

$$\begin{aligned} & \mathbb{E} \left[\left\| \mathbf{S}_i^k (\nabla f_i(x^{k+1}) - \nabla f_i(x^k)) - (\nabla f_i(x^{k+1}) - \nabla f_i(x^k)) \right\|_{\mathcal{D}}^2 \mid x^{k+1}, x^k \right] \\ &\leq \lambda_{\max} \left(\mathbf{L}_i^{\frac{1}{2}} (\mathbb{E} [\mathbf{S}_i^k \mathbf{D} \mathbf{S}_i^k] - \mathbf{D}) \mathbf{L}_i^{\frac{1}{2}} \right) \|\nabla f_i(x^{k+1}) - \nabla f_i(x^k)\|_{\mathbf{L}_i^{-1}}^2. \end{aligned} \tag{26}$$

Assumption 2 implies

$$\begin{aligned} & \mathbb{E} \left[\left\| \mathbf{S}_i^k (\nabla f_i(x^{k+1}) - \nabla f_i(x^k)) - (\nabla f_i(x^{k+1}) - \nabla f_i(x^k)) \right\|_{\mathcal{D}}^2 \mid x^{k+1}, x^k \right] \\ &\leq \lambda_{\max} \left(\mathbf{L}_i^{\frac{1}{2}} (\mathbb{E} [\mathbf{S}_i^k \mathbf{D} \mathbf{S}_i^k] - \mathbf{D}) \mathbf{L}_i^{\frac{1}{2}} \right) \|x^{k+1} - x^k\|_{\mathbf{L}_i}^2. \end{aligned} \tag{27}$$

1458 Plugging (27) into (25), we deduce

$$1459 \mathbb{E} \left[\|g^{k+1} - \nabla f(x^{k+1})\|_D^2 \mid x^{k+1}, x^k \right]$$

$$1460 \leq \frac{1-p}{n^2} \sum_{i=1}^n \lambda_{\max} \left(L_i^{\frac{1}{2}} \left(\mathbb{E} [S_i^k D S_i^k] - D \right) L_i^{\frac{1}{2}} \right) \|x^{k+1} - x^k\|_{L_i}^2 + (1-p) \cdot \|g^k - \nabla f(x^k)\|_D^2.$$

1461

1462 Replacing L_i^{-1} with $L^{-1/2} L^{1/2} L_i^{-1} L^{1/2} L^{-1/2}$, we denote that

$$1463 \lambda_i := \lambda_{\max} \left(L_i^{\frac{1}{2}} \left(\mathbb{E} [S_i^k D S_i^k] - D \right) L_i^{\frac{1}{2}} \right),$$

1464 and rewrite the L_i -norm in the first term of RHS by the L -norm:

$$1465 \mathbb{E} \left[\|g^{k+1} - \nabla f(x^{k+1})\|_D^2 \mid x^{k+1}, x^k \right]$$

$$1466 = \frac{1-p}{n^2} \sum_{i=1}^n \lambda_i \cdot \left(L^{\frac{1}{2}} (x^{k+1} - x^k) \right)^\top L^{-\frac{1}{2}} L_i L^{-\frac{1}{2}} \left(L^{\frac{1}{2}} (x^{k+1} - x^k) \right)$$

$$1467 + (1-p) \|g^k - \nabla f(x^k)\|_D^2$$

$$1468 \leq \frac{1-p}{n^2} \sum_{i=1}^n \lambda_i \cdot \lambda_{\max} \left(L^{-\frac{1}{2}} L_i L^{-\frac{1}{2}} \right) \|x^{k+1} - x^k\|_L^2 + (1-p) \cdot \|g^k - \nabla f(x^k)\|_D^2.$$

1469 We further use Fact 5 to upper bound $\lambda_{\max} \left(L_i^{\frac{1}{2}} \left(\mathbb{E} [S_i^k D S_i^k] - D \right) L_i^{\frac{1}{2}} \right)$ by the product of

1470 $\lambda_{\max}(L_i)$ and $\lambda_{\max}(\mathbb{E}[S_i^k D S_i^k] - D)$. This allows us to simplify the expression since

1471 $\lambda_{\max}(\mathbb{E}[S_i^k D S_i^k] - D)$ is independent of the index i . Notice that we have already defined

$$1472 R(D, \mathcal{S}) = \frac{1}{n} \sum_{i=1}^n \lambda_{\max}(\mathbb{E}[S_i^k D S_i^k] - D) \cdot \lambda_{\max}(L_i) \cdot \lambda_{\max}(L^{-\frac{1}{2}} L_i L^{-\frac{1}{2}}).$$

1473 Taking expectation, using tower property and using the definition above, we deduce

$$1474 \mathbb{E} \left[\|g^{k+1} - \nabla f(x^{k+1})\|_D^2 \right]$$

$$1475 \leq \frac{(1-p) \cdot R(D, \mathcal{S})}{n} \mathbb{E} \left[\|x^{k+1} - x^k\|_L^2 \right] + (1-p) \mathbb{E} \left[\|g^k - \nabla f(x^k)\|_D^2 \right]. \quad (28)$$

1476 We construct the following Lyapunov function Φ_k ,

$$1477 \Phi_k = f(x^k) - f^* + \frac{1}{2p} \|g^k - \nabla f(x^k)\|_D^2. \quad (29)$$

1478 Using (24) and (28), we are able to get

$$1479 \mathbb{E} [\Phi_{k+1}] \leq \frac{1}{2p} \left[\frac{(1-p) \cdot R(D, \mathcal{S})}{n} \mathbb{E} \left[\|x^{k+1} - x^k\|_L^2 \right] + (1-p) \cdot \mathbb{E} \left[\|g^k - \nabla f(x^k)\|_D^2 \right] \right]$$

$$1480 + \mathbb{E} [f(x^k) - f^*] - \frac{1}{2} \mathbb{E} \left[\|\nabla f(x^k)\|_D^2 \right] + \frac{1}{2} \mathbb{E} \left[\|g^k - \nabla f(x^k)\|_D^2 \right]$$

$$1481 - \frac{1}{2} \mathbb{E} \left[\|x^{k+1} - x^k\|_{D^{-1}-L}^2 \right]$$

$$1482 = \mathbb{E} [\Phi_k] - \frac{1}{2} \mathbb{E} \left[\|\nabla f(x^k)\|_D^2 \right]$$

$$1483 + \left(\frac{(1-p) \cdot R(D, \mathcal{S})}{2np} \mathbb{E} \left[\|x^{k+1} - x^k\|_L^2 \right] - \frac{1}{2} \mathbb{E} \left[\|x^{k+1} - x^k\|_{D^{-1}-L}^2 \right] \right)$$

$$1484 = \mathbb{E} [\Phi_k] - \frac{1}{2} \mathbb{E} \left[\|\nabla f(x^k)\|_D^2 \right]$$

$$1485 + \frac{1}{2} \left(\frac{(1-p) \cdot R(D, \mathcal{S})}{np} \mathbb{E} \left[\|x^{k+1} - x^k\|_L^2 \right] - \mathbb{E} \left[\|x^{k+1} - x^k\|_{D^{-1}-L}^2 \right] \right).$$

We can rewrite the last term as

$$\mathbb{E} \left[(x^{k+1} - x^k)^\top \left[\frac{(1-p) \cdot R(\mathbf{D}, \mathcal{S})}{np} \mathbf{L} + \mathbf{L} - \mathbf{D}^{-1} \right] (x^{k+1} - x^k) \right]. \quad (30)$$

We require the matrix in between to be negative semi-definite, which is

$$\mathbf{D}^{-1} \succeq \left(\frac{(1-p) \cdot R(\mathbf{D}, \mathcal{S})}{np} + 1 \right) \mathbf{L}.$$

This leads to the result that the expression (30) is always non-positive. After dropping the last term, the relation between $\mathbb{E} [\Phi_{k+1}]$ and $\mathbb{E} [\Phi_k]$ becomes

$$\mathbb{E} [\Phi_{k+1}] \leq \mathbb{E} [\Phi_k] - \frac{1}{2} \mathbb{E} \left[\|\nabla f(x^k)\|_{\mathbf{D}}^2 \right].$$

Unrolling this recurrence, we get

$$\frac{1}{K} \sum_{k=0}^{K-1} \mathbb{E} \left[\|\nabla f(x^k)\|_{\mathbf{D}}^2 \right] \leq \frac{2(\mathbb{E} [\Phi_0] - \mathbb{E} [\Phi_K])}{K}. \quad (31)$$

The left hand side can be viewed as $\mathbb{E} \left[\|\nabla f(\tilde{x}^K)\|_{\mathbf{D}}^2 \right]$, where \tilde{x}^K is drawn uniformly at random from $\{x_k\}_{k=0}^{K-1}$. From $\Phi_K > 0$, we obtain

$$\begin{aligned} \frac{2(\mathbb{E} [\Phi_0] - \mathbb{E} [\Phi_K])}{K} &\leq \frac{2\Phi_0}{K} \\ &= \frac{2 \left(f(x^0) - f^* + \frac{1}{2p} \|g^0 - \nabla f(x^0)\|_{\mathbf{D}}^2 \right)}{K} \\ &= \frac{2(f(x^0) - f^*)}{K}. \end{aligned}$$

Plugging in the simplified result into (31), and performing determinant normalization, we get

$$\mathbb{E} \left[\|\nabla f(\tilde{x}^K)\|_{\frac{\mathbf{D}}{\det(\mathbf{D})^{1/d}}}^2 \right] \leq \frac{2(f(x^0) - f^*)}{\det(\mathbf{D})^{1/d} K}. \quad (32)$$

Remark 7. We can achieve a slightly more refined stepsize condition than (4) for *det-MARINA*, which is given as follows

$$\mathbf{D} \succeq \left(\frac{(1-p) \cdot \tilde{R}(\mathbf{D}, \mathcal{S})}{np} + 1 \right) \mathbf{L}, \quad (33)$$

where

$$\tilde{R}(\mathbf{D}, \mathcal{S}) := \frac{1}{n} \sum_{i=1}^n \lambda_{\max} \left(\mathbf{L}_i^{\frac{1}{2}} \left(\mathbb{E} [\mathbf{S}_i^k \mathbf{D} \mathbf{S}_i^k] - \mathbf{D} \right) \mathbf{L}_i^{\frac{1}{2}} \right) \cdot \lambda_{\max} \left(\mathbf{L}^{-\frac{1}{2}} \mathbf{L}_i \mathbf{L}^{-\frac{1}{2}} \right).$$

This is obtained if we do not use Fact 5 to upper bound $\lambda_{\max} \left(\mathbf{L}_i^{\frac{1}{2}} \left(\mathbb{E} [\mathbf{S}_i^k \mathbf{D} \mathbf{S}_i^k] - \mathbf{D} \right) \mathbf{L}_i^{\frac{1}{2}} \right)$ by the product of $\lambda_{\max}(\mathbf{L}_i)$ and $\lambda_{\max}(\mathbb{E}[\mathbf{S}_i^k \mathbf{D} \mathbf{S}_i^k] - \mathbf{D})$. However, (33) results in a condition that is much harder to solve even if we assume $\mathbf{D} = \gamma \cdot \mathbf{W}$. So instead of using the more refined condition (33), we turn to (4). Notice that both of the two conditions (33) and (4) reduce to the stepsize condition for *MARINA* in the scalar setting.

D.3 COMPARISON OF DIFFERENT STEPSIZES

In Corollary 3, we focus on the special stepsize where we fix $\mathbf{W} = \mathbf{L}^{-1}$, and show that in this case *det-MARINA* always beats *MARINA* in terms of both iteration and communication complexities. However, other choices for \mathbf{W} are also possible. Specifically, we consider the cases where $\mathbf{W} = \text{diag}^{-1}(\mathbf{L})$ and $\mathbf{W} = \mathbf{I}_d$.

1566 D.3.1 THE DIAGONAL CASE

1567 We consider $\mathbf{W} = \text{diag}^{-1}(\mathbf{L})$. The following corollary describes the optimal stepsize and the
1568 iteration complexity.
1569

1570 **Corollary 8.** *If we take $\mathbf{W} = \text{diag}^{-1}(\mathbf{L})$ in Corollary 1, then the optimal stepsize satisfies*
1571

$$1572 \mathbf{D}_{\text{diag}^{-1}(\mathbf{L})}^* = \frac{2}{1 + \sqrt{1 + 4\alpha\beta \cdot \Lambda_{\text{diag}^{-1}(\mathbf{L}),\mathcal{S}}}} \cdot \text{diag}^{-1}(\mathbf{L}). \quad (34)$$

1573 *This stepsize results in a better iteration complexity of **det-MARINA** compared to scalar **MARINA**.*
1574

1575 From this corollary we know that **det-MARINA** has a better iteration complexity when $\mathbf{W} =$
1576 $\text{diag}^{-1}(\mathbf{L})$. And since the same sketch is used for **MARINA** and **det-MARINA**, the communi-
1577 cation complexity is improved as well. However, in general there is no clear relation between the
1578 iteration complexity of $\mathbf{W} = \mathbf{L}^{-1}$ case and $\mathbf{W} = \text{diag}^{-1}(\mathbf{L})$ case. This is also confirmed by one
1579 of our experiments, see Figure 6 to see the comparison of **det-MARINA** using optimal stepsizes in
1580 different cases.
1581

1583 D.3.2 THE IDENTITY CASE

1584 In this setting, \mathbf{W} is the d -dimensional identity matrix \mathbf{I}_d . Then the stepsize of our algorithm reduces
1585 to a scalar γ , where γ is determined through Corollary 1. Notice that in this case we do not reduce
1586 to the standard **MARINA** case because we are still using the matrix Lipschitz gradient assumption
1587 with $\mathbf{L} \in \mathbb{S}_{++}^d$.
1588

1589 **Corollary 9.** *If we take $\mathbf{W} = \mathbf{I}_d$, the optimal stepsize is given by*
1590

$$1591 \mathbf{D}_{\mathbf{I}_d}^* = \frac{2}{1 + \sqrt{1 + 4\alpha\beta \frac{1}{\lambda_{\max}(\mathbf{L})} \cdot \omega}} \cdot \frac{\mathbf{I}_d}{\lambda_{\max}(\mathbf{L})}. \quad (35)$$

1592 *This stepsize results in a better iteration complexity of **det-MARINA** compared to scalar **MARINA**.*
1593

1594 The result in this corollary tells us that using scalar stepsize with matrix Lipschitz gradient assump-
1595 tion alone can result in acceleration of **MARINA**. However, the use of matrix stepsize allows us
1596 to also take into consideration the "structure" of the stepsize, thus allows more flexibility. When
1597 the structure of the stepsize is chosen properly, combining matrix gradient Lipschitzness and matrix
1598 stepsize can result in a faster rate, as it can also be observed from the experiments in Figure 6. The
1599 choices of \mathbf{W} we consider here are in some sense inspired by the matrix stepsize **GD**, where the
1600 optimal stepsize is \mathbf{L}^{-1} . In general, how to identify the best structure for the matrix stepsize remains
1601 a open problem.
1602

1604 D.4 PROOFS OF THE COROLLARIES

1605 D.4.1 PROOF OF COROLLARY 1

1606 We start with rewriting (4) as
1607

$$1608 \left(\frac{1-p}{np} \cdot R(\mathbf{D}, \mathcal{S}) + 1 \right) \mathbf{D}^{\frac{1}{2}} \mathbf{L} \mathbf{D}^{\frac{1}{2}} \preceq \mathbf{I}_d.$$

1609 Plugging in the definition of $R(\mathbf{D}, \mathcal{S})$ and $\mathbf{D} = \gamma \mathbf{W}$, we get
1610

$$1611 \gamma \left(\frac{1-p}{np} \cdot \frac{1}{n} \sum_{i=1}^n \lambda_{\max}(\mathbf{L}_i) \lambda_{\max}(\mathbf{L}^{-1} \mathbf{L}_i) \cdot \lambda_{\max}(\mathbb{E}[\mathbf{S}_i^k \mathbf{W} \mathbf{S}_i^k] - \mathbf{W}) \cdot \gamma + 1 \right) \mathbf{W}^{\frac{1}{2}} \mathbf{L} \mathbf{W}^{\frac{1}{2}} \preceq \mathbf{I}_d.$$

1612 This generalized inequality is equivalent to the following inequality,
1613

$$1614 \gamma \left(\frac{1-p}{np} \cdot \frac{1}{n} \sum_{i=1}^n \lambda_{\max}(\mathbf{L}_i) \lambda_{\max}(\mathbf{L}^{-1} \mathbf{L}_i) \cdot \lambda_{\max}(\mathbb{E}[\mathbf{S}_i^k \mathbf{W} \mathbf{S}_i^k] - \mathbf{W}) \cdot \gamma + 1 \right) \cdot \lambda_{\max}(\mathbf{W}^{\frac{1}{2}} \mathbf{L} \mathbf{W}^{\frac{1}{2}}) \leq 1,$$

1620 which is a quadratic inequality on γ . Notice that we have already defined

$$1621 \alpha = \frac{1-p}{np}; \quad \beta = \frac{1}{n} \sum_{i=1}^n \lambda_{\max}(\mathbf{L}_i) \cdot \lambda_{\max}(\mathbf{L}^{-1}\mathbf{L}_i);$$

$$1622 \Lambda_{\mathbf{W},\mathcal{S}} = \lambda_{\max}(\mathbb{E}[\mathbf{S}_i^k \mathbf{W} \mathbf{S}_i^k] - \mathbf{W}); \quad \lambda_{\mathbf{W}} = \lambda_{\max}^{-1}(\mathbf{W}^{\frac{1}{2}} \mathbf{L} \mathbf{W}^{\frac{1}{2}}).$$

1623 As a result, the above inequality can be written equivalently as

$$1624 \alpha \beta \Lambda_{\mathbf{W},\mathcal{S}} \cdot \gamma^2 + \gamma - \lambda_{\mathbf{W}} \leq 0,$$

1625 which yields the upper bound on γ

$$1626 \gamma \leq \frac{\sqrt{1 + 4\alpha\beta \cdot \Lambda_{\mathbf{W},\mathcal{S}} \lambda_{\mathbf{W}}} - 1}{2\alpha\beta \cdot \Lambda_{\mathbf{W},\mathcal{S}}}.$$

1627 Since $\sqrt{1 + 4\alpha\beta \cdot \Lambda_{\mathbf{W},\mathcal{S}} \lambda_{\mathbf{W}}} + 1 > 0$, we can simplify the result as

$$1628 \gamma \leq \frac{2\lambda_{\mathbf{W}}}{1 + \sqrt{1 + 4\alpha\beta \cdot \Lambda_{\mathbf{W},\mathcal{S}} \lambda_{\mathbf{W}}}}.$$

1629 D.4.2 PROOF OF COROLLARY 3

1630 It is obvious that (7) directly follows from plugging $\mathbf{W} = \mathbf{L}^{-1}$ into (6). The optimal stepsize is
1631 obtained as the product of γ and \mathbf{L}^{-1} . The iteration complexity of **MARINA**, according to Gorbunov
1632 et al. (2021), is

$$1633 K \geq K_1 = \mathcal{O}\left(\frac{\Delta_0 L}{\varepsilon^2} \left(1 + \sqrt{\frac{(1-p)\omega}{pn}}\right)\right). \quad (36)$$

1634 On the other hand,

$$1635 \det(\mathbf{L})^{\frac{1}{d}} \leq \lambda_{\max}(\mathbf{L}) = L. \quad (37)$$

1636 In addition, using the inequality

$$1637 \sqrt{1+4t} \leq 1 + 2\sqrt{t}, \quad (38)$$

1638 which holds for any $t \geq 0$, we have the following bound

$$1639 \frac{(1 + \sqrt{1 + 4\alpha\beta \cdot \Lambda_{\mathbf{L}^{-1},\mathcal{S}}})}{2} \leq 1 + \sqrt{\alpha\beta \cdot \Lambda_{\mathbf{L}^{-1},\mathcal{S}}}.$$

1640 Next we prove that

$$1641 1 + \sqrt{\alpha\beta \cdot \Lambda_{\mathbf{L}^{-1},\mathcal{S}}} \leq 1 + \sqrt{\frac{(1-p)}{pn}} \cdot \omega, \quad (39)$$

1642 which is equivalent to proving

$$1643 \frac{1}{n} \sum_{i=1}^n \lambda_{\max}(\mathbf{L}_i) \lambda_{\max}(\mathbf{L}_i \mathbf{L}^{-1}) \cdot \lambda_{\max}(\mathbb{E}[\mathbf{S}_i^k \mathbf{L}^{-1} \mathbf{S}_i^k] - \mathbf{L}^{-1}) \leq \omega.$$

1644 The left hand side can be upper bounded by,

$$1645 \frac{1}{n} \sum_{i=1}^n \lambda_{\max}(\mathbf{L}_i) \lambda_{\max}(\mathbf{L}^{-1} \mathbf{L}_i) \cdot \lambda_{\max}(\mathbf{L}^{-1}) \cdot \frac{\lambda_{\max}(\mathbb{E}[\mathbf{S}_i^k \mathbf{L}^{-1} \mathbf{S}_i^k] - \mathbf{L}^{-1})}{\lambda_{\max}(\mathbf{L}^{-1})}$$

$$1646 \leq \frac{\lambda_{\max}(\mathbb{E}[\mathbf{S}_i^k \mathbf{L}^{-1} \mathbf{S}_i^k] - \mathbf{L}^{-1})}{\lambda_{\max}(\mathbf{L}^{-1})},$$

1647 where the inequality is a consequence of Proposition 1. We further bound the last term with

$$1648 \frac{\lambda_{\max}(\mathbb{E}[\mathbf{S}_i^k \mathbf{L}^{-1} \mathbf{S}_i^k] - \mathbf{L}^{-1})}{\lambda_{\max}(\mathbf{L}^{-1})} = \lambda_{\max}\left(\mathbb{E}\left[\mathbf{S}_i^k \cdot \frac{\mathbf{L}^{-1}}{\lambda_{\max}(\mathbf{L}^{-1})} \cdot \mathbf{S}_i^k\right] - \frac{\mathbf{L}^{-1}}{\lambda_{\max}(\mathbf{L}^{-1})}\right)$$

$$1649 \leq \lambda_{\max}(\mathbb{E}[\mathbf{S}_i^k \mathbf{S}_i^k] - \mathbf{I}_d) =: \omega.$$

Here, the last inequality is due to the monotonicity of the mapping $\mathbf{X} \mapsto \lambda_{\max}(\mathbb{E}[\mathbf{S}_i^k \mathbf{X} \mathbf{S}_i^k] - \mathbf{X})$ with $\mathbf{X} \in \mathbb{S}_{++}^d$, which can be shown as follows, let us pick any $\mathbf{X}_1, \mathbf{X}_2 \in \mathbb{S}_{++}^d$ and $\mathbf{X}_1 \preceq \mathbf{X}_2$,

$$(\mathbb{E}[\mathbf{S}_i^k \mathbf{X}_2 \mathbf{S}_i^k] - \mathbf{X}_2) - (\mathbb{E}[\mathbf{S}_i^k \mathbf{X}_1 \mathbf{S}_i^k] - \mathbf{X}_1) = \mathbb{E}[\mathbf{S}_i^k (\mathbf{X}_2 - \mathbf{X}_1) \mathbf{S}_i^k] - (\mathbf{X}_2 - \mathbf{X}_1) \succeq \mathbf{O}_d.$$

The above inequality is due to the convexity of the mapping $\mathbf{S}_i^k \mapsto \mathbf{S}_i^k \mathbf{X} \mathbf{S}_i^k$. As a result, we have

$$\lambda_{\max}(\mathbb{E}[\mathbf{S}_i^k \mathbf{X}_2 \mathbf{S}_i^k] - \mathbf{X}_2) \geq \lambda_{\max}(\mathbb{E}[\mathbf{S}_i^k \mathbf{X}_1 \mathbf{S}_i^k] - \mathbf{X}_1),$$

whenever $\mathbf{X}_2 \succeq \mathbf{X}_1$. Due to the fact that

$$\frac{\mathbf{L}^{-1}}{\lambda_{\max}(\mathbf{L}^{-1})} \preceq \mathbf{I}_d,$$

we have

$$\lambda_{\max}\left(\mathbb{E}\left[\mathbf{S}_i^k \cdot \frac{\mathbf{L}^{-1}}{\lambda_{\max}(\mathbf{L}^{-1})} \cdot \mathbf{S}_i^k\right] - \frac{\mathbf{L}^{-1}}{\lambda_{\max}(\mathbf{L}^{-1})}\right) \leq \lambda_{\max}(\mathbb{E}[\mathbf{S}_i^k \cdot \mathbf{I}_d \cdot \mathbf{S}_i^k] - \mathbf{I}_d) = \omega.$$

Combining (37) and (39), we know that the iteration complexity of **det-MARINA** is always better than that of **MARINA**.

D.4.3 PROOF OF COROLLARY 4

The number of bits sent in expectation is

$$\mathcal{O}(d + K(pd + (1-p)\zeta_S)) = \mathcal{O}((Kp+1)d + (1-p)K\zeta_S).$$

The special case where we choose $p = \zeta_S/d$ indicates that

$$\alpha = \frac{1-p}{np} = \frac{1}{n} \left(\frac{d}{\zeta_S} - 1 \right).$$

In order to reach an error of ε^2 , we need

$$K = \mathcal{O}\left(\frac{\Delta_0 \cdot \det(\mathbf{L})^{\frac{1}{d}}}{\varepsilon^2} \cdot \left(1 + \sqrt{1 + \frac{4\beta}{n} \left(\frac{d}{\zeta_S} - 1\right) \cdot \Lambda_{\mathbf{L}^{-1}, \mathcal{S}}}\right)\right),$$

which is the iteration complexity. Applying once again (38) and using the fact that $p = \zeta_S/d$, the communication complexity in this case is given by

$$\begin{aligned} & \mathcal{O}\left(d + \frac{\Delta_0 \cdot \det(\mathbf{L})^{\frac{1}{d}}}{\varepsilon^2} \cdot \left(1 + \sqrt{1 + \frac{4\beta}{n} \left(\frac{d}{\zeta_S} - 1\right) \cdot \Lambda_{\mathbf{L}^{-1}, \mathcal{S}}}\right) \cdot (pd + (1-p)\zeta_S)\right) \\ & \leq \mathcal{O}\left(d + \frac{2\Delta_0 \cdot \det(\mathbf{L})^{\frac{1}{d}}}{\varepsilon^2} \cdot \left(1 + \sqrt{\frac{\beta}{n} \left(\frac{d}{\zeta_S} - 1\right) \cdot \Lambda_{\mathbf{L}^{-1}, \mathcal{S}}}\right) \cdot (pd + (1-p)\zeta_S)\right) \\ & \leq \mathcal{O}\left(d + \frac{4\Delta_0 \cdot \det(\mathbf{L})^{\frac{1}{d}}}{\varepsilon^2} \cdot \left(\zeta_S + \sqrt{\frac{\beta \cdot \Lambda_{\mathbf{L}^{-1}, \mathcal{S}}}{n} \cdot \zeta_S(d - \zeta_S)}\right)\right). \end{aligned}$$

Ignoring the coefficient we get

$$\mathcal{O}\left(d + \frac{\Delta_0 \cdot \det(\mathbf{L})^{\frac{1}{d}}}{\varepsilon^2} \cdot \left(\zeta_S + \sqrt{\frac{\beta \cdot \Lambda_{\mathbf{L}^{-1}, \mathcal{S}}}{n} \cdot \zeta_S(d - \zeta_S)}\right)\right).$$

D.4.4 PROOF OF COROLLARY 8

Applying Corollary 1, notice that in this case

$$\lambda_{\text{diag}^{-1}(\mathbf{L})} = \lambda_{\max}^{-1}\left(\text{diag}^{-\frac{1}{2}}(\mathbf{L}) \mathbf{L} \text{diag}^{-\frac{1}{2}}(\mathbf{L})\right) = 1,$$

we obtain $\mathbf{D}_{\text{diag}^{-1}(\mathbf{L})}^*$. The iteration complexity is given by

$$\mathcal{O}\left(\frac{\det(\text{diag}(\mathbf{L}))^{\frac{1}{d}} \cdot \Delta_0}{\varepsilon^2} \cdot \left(\frac{1 + \sqrt{1 + 4\alpha\beta\Lambda_{\text{diag}^{-1}(\mathbf{L}), \mathcal{S}}}}{2}\right)\right).$$

We now compare it to the iteration complexity of **MARINA**, which is given in (36). We know that each diagonal element \mathbf{L}_{jj} satisfies $\mathbf{L}_{jj} \leq \lambda_{\max}(\mathbf{L}) = L$ for $j = 1, \dots, d$. As a result,

$$\det(\text{diag}(\mathbf{L}))^{\frac{1}{d}} \leq L. \quad (40)$$

From (38), we deduce

$$\frac{1 + \sqrt{1 + 4\alpha\beta \cdot \Lambda_{\text{diag}^{-1}(\mathbf{L}), \mathcal{S}}}}{2} \leq 1 + \sqrt{\alpha\beta \cdot \Lambda_{\text{diag}^{-1}(\mathbf{L}), \mathcal{S}}}.$$

Now, let us prove the below inequality

$$1 + \sqrt{\alpha\beta \cdot \Lambda_{\text{diag}^{-1}(\mathbf{L}), \mathcal{S}}} \leq 1 + \sqrt{\frac{(1-p)}{pn}} \cdot \omega. \quad (41)$$

The latter is equivalent to

$$\beta \cdot \Lambda_{\text{diag}^{-1}(\mathbf{L}), \mathcal{S}} \leq \omega.$$

Plugging in the definition of β , ω and $\Lambda_{\text{diag}^{-1}(\mathbf{L}), \mathcal{S}}$ and using the relation given in Proposition 1, we obtain,

$$\lambda_{\max} \left(\mathbb{E} \left[\mathbf{S}_i^k \frac{\text{diag}^{-1}(\mathbf{L})}{\lambda_{\max}(\mathbf{L}^{-1})} \mathbf{S}_i^k - \frac{\text{diag}^{-1}(\mathbf{L})}{\lambda_{\max}(\mathbf{L}^{-1})} \right] \right) \leq \lambda_{\max} \left(\mathbb{E} [\mathbf{S}_i^k \mathbf{I}_d \mathbf{S}_i^k] - \mathbf{I}_d \right).$$

Thus, it is enough to prove that

$$\frac{\text{diag}^{-1}(\mathbf{L})}{\lambda_{\max}(\mathbf{L}^{-1})} \preceq \mathbf{I}_d.$$

We can further simplify the above inequality as

$$\lambda_{\min}(\mathbf{L}) \leq \lambda_{\min}(\text{diag}(\mathbf{L})),$$

which is always true for any $\mathbf{L} \in \mathbb{S}_{++}^d$. Combining (40) and (41) we conclude the proof.

D.4.5 PROOF OF COROLLARY 9

Using the explicit formula for the optimal stepsize $D_{I_d}^*$, we deduce the following iteration complexity for

$$\mathcal{O} \left(\frac{\lambda_{\max}(\mathbf{L}) \Delta_0}{\varepsilon^2} \cdot \left(\frac{1 + \sqrt{1 + 4\alpha\beta \frac{\omega}{\lambda_{\max}(\mathbf{L})}}}{2} \right) \right). \quad (42)$$

Recall that $\lambda_{\max}(\mathbf{L}) = L$, we obtain using (38) that

$$\frac{1 + \sqrt{1 + 4\alpha\beta \frac{\omega}{\lambda_{\max}(\mathbf{L})}}}{2} \leq 1 + \sqrt{\alpha\beta \frac{\omega}{\lambda_{\max}(\mathbf{L})}}.$$

The comparison of two iteration complexities, given in (42) and (36) reduces to

$$1 + \sqrt{\alpha\beta \frac{\omega}{\lambda_{\max}(\mathbf{L})}} \leq 1 + \sqrt{\frac{1-p}{np}} \omega.$$

This is equivalent to

$$\beta \cdot \frac{1}{\lambda_{\max}(\mathbf{L})} \leq 1.$$

Utilizing Proposition 1, the above inequality can be rewritten as

$$\frac{1}{\lambda_{\max}(\mathbf{L}^{-1}) \cdot \lambda_{\max}(\mathbf{L})} \leq 1,$$

which is exactly

$$\lambda_{\min}(\mathbf{L}) \leq \lambda_{\max}(\mathbf{L}).$$

1782 E ANALYSIS OF DET-DASHA

1783 E.1 PROOF OF THEOREM 2

1784 Using Lemma 1 and taking expectations, we are able to obtain

$$\begin{aligned}
1785 & \mathbb{E} [f(x^{k+1})] \\
1786 & \leq \mathbb{E} [f(x^k)] - \frac{1}{2} \mathbb{E} [\|\nabla f(x^k)\|_{\mathcal{D}}^2] - \frac{1}{2} \mathbb{E} [\|x^{k+1} - x^k\|_{\mathcal{D}^{-1}-\mathcal{L}}^2] + \frac{1}{2} \mathbb{E} [\|g^k - \nabla f(x^k)\|_{\mathcal{D}}^2] \\
1787 & \leq \mathbb{E} [f(x^k)] - \frac{1}{2} \mathbb{E} [\|\nabla f(x^k)\|_{\mathcal{D}}^2] - \frac{1}{2} \mathbb{E} [\|x^{k+1} - x^k\|_{\mathcal{D}^{-1}-\mathcal{L}}^2] \\
1788 & \quad + \mathbb{E} \left[\frac{1}{2} \|g^k - h^k + h^k - \nabla f(x^k)\|_{\mathcal{D}}^2 \right] \\
1789 & \leq \mathbb{E} [f(x^k)] - \frac{1}{2} \mathbb{E} [\|\nabla f(x^k)\|_{\mathcal{D}}^2] - \frac{1}{2} \mathbb{E} [\|x^{k+1} - x^k\|_{\mathcal{D}^{-1}-\mathcal{L}}^2] \\
1790 & \quad + \mathbb{E} [\|g^k - h^k\|_{\mathcal{D}}^2 + \|h^k - \nabla f(x^k)\|_{\mathcal{D}}^2], \tag{43}
\end{aligned}$$

1791 where the last step is due to the convexity of the norm. Using Lemma 4, we obtain

$$\begin{aligned}
1792 & \mathbb{E} [\|g^{k+1} - h^{k+1}\|_{\mathcal{D}}^2] \leq \frac{2\omega_{\mathcal{D}} \cdot \lambda_{\max}(\mathcal{D})}{n^2} \sum_{i=1}^n \lambda_{\max}(\mathbf{L}_i) \mathbb{E} [\|h_i^{k+1} - h_i^k\|_{\mathcal{L}_i^{-1}}^2] \\
1793 & \quad + \frac{2a^2\omega_{\mathcal{D}}}{n^2} \sum_{i=1}^n \mathbb{E} [\|g_i^k - h_i^k\|_{\mathcal{D}}^2] + (1-a)^2 \mathbb{E} [\|g^k - h^k\|_{\mathcal{D}}^2]. \tag{44}
\end{aligned}$$

1794 Using Lemma 5, we get

$$\begin{aligned}
1795 & \mathbb{E} [\|g_i^{k+1} - h_i^{k+1}\|_{\mathcal{D}}^2] \leq (2a^2\omega_{\mathcal{D}} + (1-a)^2) \cdot \mathbb{E} [\|g_i^k - h_i^k\|_{\mathcal{D}}^2] \\
1796 & \quad + 2\omega_{\mathcal{D}} \cdot \lambda_{\max}(\mathcal{D}) \cdot \lambda_{\max}(\mathbf{L}_i) \cdot \mathbb{E} [\|h_i^{k+1} - h_i^k\|_{\mathcal{L}_i^{-1}}^2]. \tag{45}
\end{aligned}$$

1797 Now let us fix $\kappa \in [0, +\infty)$, $\eta \in [0, +\infty)$ which we will determine later, and construct the following Lyapunov function Φ_k

$$\Phi_k = \mathbb{E} [f(x^k) - f^*] + \kappa \cdot \mathbb{E} [\|g^k - h^k\|_{\mathcal{D}}^2] + \eta \cdot \mathbb{E} \left[\frac{1}{n} \sum_{i=1}^n \|g_i^k - h_i^k\|_{\mathcal{D}}^2 \right]. \tag{46}$$

1798 Combining (43), (44) and (45), we get

$$\begin{aligned}
1799 & \Phi_{k+1} \\
1800 & \leq \mathbb{E} \left[f(x^k) - f^* - \frac{1}{2} \|\nabla f(x^k)\|_{\mathcal{D}}^2 \right] \\
1801 & \quad + \mathbb{E} \left[-\frac{1}{2} \|x^{k+1} - x^k\|_{\mathcal{D}^{-1}-\mathcal{L}}^2 + \|g^k - h^k\|_{\mathcal{D}}^2 + \|h^k - \nabla f(x^k)\|_{\mathcal{D}}^2 \right] \\
1802 & \quad + \kappa(1-a)^2 \mathbb{E} [\|g^k - h^k\|_{\mathcal{D}}^2] + \frac{2\kappa \cdot \omega_{\mathcal{D}} \lambda_{\max}(\mathcal{D})}{n} \cdot \frac{1}{n} \sum_{i=1}^n \lambda_{\max}(\mathbf{L}_i) \mathbb{E} [\|h_i^{k+1} - h_i^k\|_{\mathcal{L}_i^{-1}}^2] \\
1803 & \quad + \frac{2a^2\omega_{\mathcal{D}} \cdot \kappa}{n} \cdot \frac{1}{n} \sum_{i=1}^n \mathbb{E} [\|g_i^k - h_i^k\|_{\mathcal{D}}^2] + \eta (2a^2\omega_{\mathcal{D}} + (1-a)^2) \cdot \frac{1}{n} \sum_{i=1}^n \mathbb{E} [\|g_i^k - h_i^k\|_{\mathcal{D}}^2] \\
1804 & \quad + 2\eta \cdot \omega_{\mathcal{D}} \cdot \lambda_{\max}(\mathcal{D}) \cdot \frac{1}{n} \sum_{i=1}^n \lambda_{\max}(\mathbf{L}_i) \cdot \mathbb{E} [\|h_i^{k+1} - h_i^k\|_{\mathcal{L}_i^{-1}}^2].
\end{aligned}$$

Rearranging terms, and notice that $\|h^k - \nabla f(x^k)\|_{\mathbf{D}}^2 = 0$,

$$\begin{aligned}
\Phi_{k+1} &\leq \mathbb{E} [f(x^k) - f^*] - \frac{1}{2} \mathbb{E} [\|\nabla f(x^k)\|_{\mathbf{D}}^2] \\
&\quad - \frac{1}{2} \mathbb{E} [\|x^{k+1} - x^k\|_{\mathbf{D}^{-1}-\mathbf{L}}^2] + (1 + \kappa(1-a)^2) \mathbb{E} [\|g^k - h^k\|_{\mathbf{D}}^2] \\
&\quad + \left(\frac{2a^2\omega_{\mathbf{D}} \cdot \kappa}{n} + \eta(2a^2\omega_{\mathbf{D}} + (1-a)^2) \right) \cdot \frac{1}{n} \sum_{i=1}^n \mathbb{E} [\|g_i^k - h_i^k\|_{\mathbf{D}}^2] \\
&\quad + \left(\frac{2\kappa \cdot \omega_{\mathbf{D}} \lambda_{\max}(\mathbf{D})}{n} + 2\eta \cdot \omega_{\mathbf{D}} \cdot \lambda_{\max}(\mathbf{D}) \right) \cdot \frac{1}{n} \sum_{i=1}^n \lambda_{\max}(\mathbf{L}_i) \cdot \mathbb{E} [\|h_i^{k+1} - h_i^k\|_{\mathbf{L}_i^{-1}}^2].
\end{aligned}$$

In order to proceed, we consider the choice of κ and η , for κ ,

$$1 + \kappa(1-a)^2 \leq \kappa. \quad (47)$$

It is then clear that the choice of $\kappa = \frac{1}{a}$ satisfies the condition. On the other hand, we look at the terms involving $\mathbb{E} [\|g_i^k - h_i^k\|_{\mathbf{D}}^2]$, we can rewrite as

$$T_1 := \left(\frac{2a^2\omega_{\mathbf{D}} \cdot \kappa}{n} + \eta(2a^2\omega_{\mathbf{D}} + (1-a)^2) \right) \cdot \frac{1}{n} \sum_{i=1}^n \mathbb{E} [\|g_i^k - h_i^k\|_{\mathbf{D}}^2].$$

Picking $\kappa = \frac{1}{a}$ and $a = \frac{1}{2\omega_{\mathbf{D}}+1}$, the T_1 can be simplified as

$$T_1 = \left(\frac{2\omega_{\mathbf{D}}}{n \cdot (2\omega_{\mathbf{D}} + 1)} + \eta \cdot \frac{4\omega_{\mathbf{D}}^2 + 2\omega_{\mathbf{D}}}{(2\omega_{\mathbf{D}} + 1)^2} \right) \cdot \frac{1}{n} \sum_{i=1}^n \mathbb{E} [\|g_i^k - h_i^k\|_{\mathbf{D}}^2].$$

We pick η so that it satisfies

$$\left(\frac{2\omega_{\mathbf{D}}}{n \cdot (2\omega_{\mathbf{D}} + 1)} + \eta \cdot \frac{4\omega_{\mathbf{D}}^2 + 2\omega_{\mathbf{D}}}{(2\omega_{\mathbf{D}} + 1)^2} \right) \leq \eta. \quad (48)$$

Taking $\eta = \frac{2\omega_{\mathbf{D}}}{n}$, which is the minimum value satisfying (48), we conclude that

$$T_1 \leq \eta \cdot \frac{1}{n} \sum_{i=1}^n \mathbb{E} [\|g_i^k - h_i^k\|_{\mathbf{D}}^2]. \quad (49)$$

Combining (47) and (49), we are able to conclude that

$$\begin{aligned}
\Phi_{k+1} &\leq \mathbb{E} [f(x^k) - f^*] + \kappa \cdot \mathbb{E} [\|g^k - h^k\|_{\mathbf{D}}^2] + \eta \cdot \frac{1}{n} \sum_{i=1}^n \mathbb{E} [\|g_i^k - h_i^k\|_{\mathbf{D}}^2] \\
&\quad - \frac{1}{2} \mathbb{E} [\|\nabla f(x^k)\|_{\mathbf{D}}^2] - \frac{1}{2} \mathbb{E} [\|x^{k+1} - x^k\|_{\mathbf{D}^{-1}-\mathbf{L}}^2] \\
&\quad + \left(\frac{2\kappa \cdot \omega_{\mathbf{D}} \lambda_{\max}(\mathbf{D})}{n} + 2\eta \cdot \omega_{\mathbf{D}} \cdot \lambda_{\max}(\mathbf{D}) \right) \cdot \frac{1}{n} \sum_{i=1}^n \lambda_{\max}(\mathbf{L}_i) \cdot \mathbb{E} [\|h_i^{k+1} - h_i^k\|_{\mathbf{L}_i^{-1}}^2].
\end{aligned}$$

Using the definition of Φ_k and Lemma 3, we obtain

$$\begin{aligned}
\Phi_{k+1} &\leq \Phi_k - \frac{1}{2} \mathbb{E} [\|\nabla f(x^k)\|_{\mathbf{D}}^2] - \frac{1}{2} \mathbb{E} [\|x^{k+1} - x^k\|_{\mathbf{D}^{-1}-\mathbf{L}}^2] \\
&\quad \left(\frac{2\kappa \cdot \omega_{\mathbf{D}} \lambda_{\max}(\mathbf{D})}{n} + 2\eta \cdot \omega_{\mathbf{D}} \cdot \lambda_{\max}(\mathbf{D}) \right) \cdot \frac{1}{n} \sum_{i=1}^n \lambda_{\max}(\mathbf{L}_i) \cdot \mathbb{E} [\|x^{k+1} - x^k\|_{\mathbf{L}_i}^2] \\
&= \Phi_k - \frac{1}{2} \mathbb{E} [\|\nabla f(x^k)\|_{\mathbf{D}}^2] + \mathbb{E} [\|x^{k+1} - x^k\|_{\mathbf{N}}^2],
\end{aligned}$$

where $\mathbf{N} \in \mathbb{S}^d$ is defined as

$$\mathbf{N} := \left(\frac{2\kappa \cdot \omega_{\mathbf{D}} \lambda_{\max}(\mathbf{D})}{n} + 2\eta \cdot \omega_{\mathbf{D}} \cdot \lambda_{\max}(\mathbf{D}) \right) \cdot \frac{1}{n} \sum_{i=1}^n \lambda_{\max}(\mathbf{L}_i) \cdot \mathbf{L}_i - \frac{1}{2} \mathbf{D}^{-1} + \frac{1}{2} \mathbf{L}.$$

We require $\mathbf{N} \preceq \mathbf{O}_d$, which leads to the condition on \mathbf{D} :

$$\mathbf{D}^{-1} - \mathbf{L} - \frac{4\lambda_{\max}(\mathbf{D}) \cdot \omega_{\mathbf{D}} \cdot (4\omega_{\mathbf{D}} + 1)}{n} \cdot \frac{1}{n} \sum_{i=1}^n \lambda_{\max}(\mathbf{L}_i) \cdot \mathbf{L}_i \succeq \mathbf{O}_d.$$

Given the above condition is satisfied, we have the recurrence

$$\frac{1}{2} \mathbb{E} \left[\|\nabla f(x^k)\|_{\mathbf{D}}^2 \right] \leq \Phi_k - \Phi_{k+1}$$

Summing up for $k = 0 \dots K-1$, we obtain

$$\sum_{k=0}^{K-1} \mathbb{E} \left[\|\nabla f(x^k)\|_{\mathbf{D}}^2 \right] \leq 2(\Phi_0 - \Phi_K). \quad (50)$$

Notice that we also have

$$\begin{aligned} \Phi_0 &= f(x^0) - f^* + (2\omega_{\mathbf{D}} + 1) \|g^0 - h^0\|_{\mathbf{D}}^2 + \frac{2\omega_{\mathbf{D}}}{n} \sum_{i=1}^n \|g_i^0 - h_i^0\|^2 \\ &= f(x^0) - f^*, \end{aligned}$$

We divide both sides of (50) by K , and perform determinant normalization,

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

This is to say

$$\mathbb{E} \left[\|\nabla f(\tilde{x}^K)\|_{\frac{\mathbf{D}}{\det(\mathbf{D})^{1/d}}}^2 \right] \leq \frac{2(f(x^0) - f^*)}{\det(\mathbf{D})^{1/d} \cdot K},$$

where \tilde{x}^K is chosen uniformly randomly from the first K iterates of the algorithm.

E.2 PROOFS OF THE COROLLARIES

E.2.1 PROOF OF COROLLARY 2

Plug $\mathbf{D} = \gamma_{\mathbf{W}} \cdot \mathbf{W}$ into the stepsize condition in Theorem 2, we obtain

$$\frac{1}{\gamma_{\mathbf{W}}} \cdot \mathbf{W}^{-1} - \mathbf{L} - \frac{4\gamma_{\mathbf{W}} \cdot \lambda_{\max}(\mathbf{W}) \cdot \omega_{\mathbf{W}} (4\omega_{\mathbf{W}} + 1)}{n} \cdot \frac{1}{n} \sum_{i=1}^n \lambda_{\max}(\mathbf{L}_i) \cdot \mathbf{L}_i \succeq \mathbf{O}_d.$$

We then simplify the above condition as

$$\begin{aligned} &\frac{1}{\gamma_{\mathbf{W}}} \cdot \mathbf{L}^{-\frac{1}{2}} \mathbf{W}^{-1} \mathbf{L}^{-\frac{1}{2}} \\ &\succeq \mathbf{I}_d + \frac{4\gamma_{\mathbf{W}} \cdot \lambda_{\max}(\mathbf{W}) \cdot \omega_{\mathbf{W}} (4\omega_{\mathbf{W}} + 1)}{n} \cdot \mathbf{L}^{-\frac{1}{2}} \left(\frac{1}{n} \sum_{i=1}^n \lambda_{\max}(\mathbf{L}_i) \cdot \mathbf{L}_i \right) \mathbf{L}^{-\frac{1}{2}}. \end{aligned}$$

Using Proposition 3, we have

$$\frac{1}{\gamma_{\mathbf{W}}} \cdot \mathbf{L}^{-\frac{1}{2}} \mathbf{W}^{-1} \mathbf{L}^{-\frac{1}{2}} - \frac{4\gamma_{\mathbf{W}} \cdot \lambda_{\max}(\mathbf{W}) \cdot \omega_{\mathbf{W}} (4\omega_{\mathbf{W}} + 1)}{n} \cdot \lambda_{\min}(\mathbf{L}) \cdot \mathbf{I}_d \succeq \mathbf{I}_d.$$

Taking the minimum eigenvalue of both sides, we obtain that,

$$\frac{1}{\gamma_{\mathbf{W}}} \cdot \lambda_{\min} \left(\mathbf{L}^{-\frac{1}{2}} \mathbf{W}^{-1} \mathbf{L}^{-\frac{1}{2}} \right) - \frac{4\gamma_{\mathbf{W}} \cdot \lambda_{\max}(\mathbf{W}) \cdot \omega_{\mathbf{W}} (4\omega_{\mathbf{W}} + 1)}{n} \cdot \lambda_{\min}(\mathbf{L}) \geq 1,$$

If we denote $C_{\mathbf{W}} := \frac{\lambda_{\max}(\mathbf{W}) \cdot \omega_{\mathbf{W}} (4\omega_{\mathbf{W}} + 1)}{n} > 0$, and $\lambda_{\mathbf{W}} := \lambda_{\max}^{-1} \left(\mathbf{L}^{\frac{1}{2}} \mathbf{W} \mathbf{L}^{\frac{1}{2}} \right)$, we can write

$$4 \cdot C_{\mathbf{W}} \cdot \lambda_{\min}(\mathbf{L}) \cdot \gamma_{\mathbf{W}}^2 + \gamma_{\mathbf{W}} - \lambda_{\mathbf{W}} \leq 0.$$

The solution is given by

$$\gamma_{\mathbf{W}} \leq \frac{2\lambda_{\mathbf{W}}}{1 + \sqrt{1 + 16C_{\mathbf{W}} \lambda_{\min}(\mathbf{L}) \cdot \lambda_{\mathbf{W}}}}.$$

1944 E.2.2 PROOF OF COROLLARY 5

1945 The best scaling factor in this case is given as, according to Corollary 2,

$$1946 \gamma_{L^{-1}} = \frac{2}{1 + \sqrt{1 + 16C_{L^{-1}} \cdot \lambda_{\min}(\mathbf{L})}}.$$

1947 In order to reach a ε^2 stationary point, we need

$$1948 K \geq \frac{\det(\mathbf{L})^{\frac{1}{d}} (f(x^0) - f^*)}{\varepsilon^2} \cdot \left(1 + \sqrt{1 + 16C_{L^{-1}} \cdot \lambda_{\min}(\mathbf{L})}\right).$$

1954 E.2.3 PROOF OF COROLLARY 6

1955 The iteration complexity of **det-DASHA** is given by, according to, Corollary 5,

$$1956 \mathcal{O}\left(\frac{f(x^0) - f^*}{\varepsilon^2} \cdot \left(1 + \sqrt{1 + 16C_{L^{-1}} \cdot \lambda_{\min}(\mathbf{L})}\right) \cdot \det(\mathbf{L})^{\frac{1}{d}}\right).$$

1957 Using the inequality $\sqrt{1+t} \leq 1 + \sqrt{t}$ for $t > 0$, and leaving out the coefficients, we obtain

$$1958 \mathcal{O}\left(\frac{f(x^0) - f^*}{\varepsilon^2} \cdot \left(1 + \sqrt{C_{L^{-1}} \cdot \lambda_{\min}(\mathbf{L})}\right) \cdot \det(\mathbf{L})^{\frac{1}{d}}\right).$$

1959 Notice that

$$1960 C_{L^{-1}} \cdot \lambda_{\min}(\mathbf{L}) = \lambda_{\max}(\mathbf{L}^{-1}) \cdot \frac{\omega_{L^{-1}}(4\omega_{L^{-1}} + 1)}{n} \cdot \lambda_{\min}(\mathbf{L}) = \frac{\omega_{L^{-1}}(4\omega_{L^{-1}} + 1)}{n}.$$

1961 As a result, the iteration complexity can be further simplified as

$$1962 \mathcal{O}\left(\frac{f(x^0) - f^*}{\varepsilon^2} \cdot \left(1 + \frac{\omega_{L^{-1}}}{\sqrt{n}}\right) \cdot \det(\mathbf{L})^{\frac{1}{d}}\right).$$

1963 The iteration complexity of **DASHA** is, according to Tyurin & Richtárik (2024, Corollary 6.2)

$$1964 \mathcal{O}\left(\frac{1}{\varepsilon^2} \cdot (f(x^0) - f^*) \left(L + \frac{\omega}{\sqrt{n}} \widehat{L}\right)\right),$$

1965 where $\widehat{L} = \sqrt{\frac{1}{n} \sum_{i=1}^n L_i^2}$. Since $\det(\mathbf{L})^{\frac{1}{d}} \leq \lambda_{\max}(\mathbf{L}) = L$, and $L \leq \widehat{L}$, it is easy to see that compared to **DASHA**, **det-DASHA** has a better iteration complexity when the momentum is the same. Notice that those two algorithms use the same sketch, thus, it also indicates that the communication complexity of the two algorithms are the same.

1980 E.2.4 PROOF OF COROLLARY 7

1981 The iteration complexity of **det-MARINA** is given by

$$1982 \mathcal{O}\left(\frac{f(x^0) - f^*}{\varepsilon^2} \cdot \det(\mathbf{L})^{\frac{1}{d}} \cdot \left(1 + \sqrt{\alpha\beta\Lambda_{L^{-1},S}}\right)\right),$$

1983 after removing logarithmic factors. Plugging in the definitions we obtain in the case of $\omega_{L^{-1}} + 1 = \frac{1}{p}$, we have

$$1984 \mathcal{O}\left(\frac{f(x^0) - f^*}{\varepsilon^2} \cdot \det(\mathbf{L})^{\frac{1}{d}} \cdot \left(1 + \frac{\omega_{L^{-1}}}{n}\right)\right).$$

1985 From the proof of Corollary 6, we know that the iteration complexity of **det-DASHA** is

$$1986 \mathcal{O}\left(\frac{1}{\varepsilon^2} \cdot (f(x^0) - f^*) \left(L + \frac{\omega}{\sqrt{n}} \widehat{L}\right)\right).$$

1987 It is easy to see that in this case the two algorithms have the same iteration complexity asymptotically. Notice that the communication complexity is the product of bytes sent per iteration and the number of iterations. **det-DASHA** clearly sends less bytes per iteration because it always sent the compressed gradient differences, which means that it has a better communication complexity than **det-MARINA**.

F DISTRIBUTED DET-CGD

This section is a brief summary of the distributed **det-CGD** algorithm and its theoretical analysis. The details can be found in (Li et al., 2024b). The algorithm follows the standard FL paradigm. See the pseudocode in Algorithm 3.

Algorithm 3 Distributed **det-CGD**

```

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  $g^k = \frac{1}{n} \sum_{i=1}^n S_i^k \nabla f_i(x^k)$ ;
9:   computes  $x^{k+1} = x^k - Dg^k$ ;
10:  broadcasts  $x^{k+1}$ .
11: end for
12: Return:  $x^K$ 

```

Below is the main convergence result for the algorithm.

Theorem 3. *Suppose that f is L -smooth. Under the Assumptions 1,2, if the stepsize satisfies*

$$DL D \preceq D, \quad (51)$$

then the following convergence bound is true for the iteration of Algorithm 3:

$$\min_{0 \leq k \leq K-1} \mathbb{E} \left[\left\| \nabla f(x^k) \right\|_{\frac{D}{\det(D)^{1/d}}}^2 \right] \leq \frac{2(1 + \frac{\lambda_D}{n})^K (f(x^0) - f^*)}{\det(D)^{1/d} K} + \frac{2\lambda_D \Delta^*}{\det(D)^{1/d} n}, \quad (52)$$

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

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

Remark 8. *On the right hand side of (52) we observe that increasing K will only reduce the first term, that corresponds to the convergence error. Whereas, the second term, which does not depend on K , will remain constant, if the other parameters of the algorithm are fixed. This testifies to the neighborhood phenomenon which we discussed in Section 2.*

Remark 9. *If the stepsize satisfies the below conditions,*

$$DL D \preceq D, \quad \lambda_D \leq \min \left\{ \frac{n}{K}, \frac{n\varepsilon^2}{4\Delta^*} \det(D)^{1/d} \right\}, \quad K \geq \frac{12(f(x^0) - f^*)}{\det(D)^{1/d} \varepsilon^2}, \quad (53)$$

then we obtain ε -stationary point.

One can see that in the convergence guarantee of **det-CGD** in the distributed case, the result (52) is not variance-reduced. Because of this limitation, in order to reach a ε stationary point, the stepsize condition in (53) is restrictive.

G EXTENSION OF DET-CGD2 IN MARINA FORM

In this section we want to extend **det-CGD2** into its variance reduced counterpart in **MARINA** form.

G.1 EXTENSION OF DET-CGD2 TO ITS VARIANCE REDUCED COUNTERPART

We call **det-MARINA** as the extension of **det-CGD1**, and Algorithm 4 as the extension of **det-CGD2** due to the difference in the order of applying sketches and stepsize matrices. The key difference

Algorithm 4 det-CGD2-VR

```

1: Input: starting point  $x^0$ , stepsize matrix  $\mathbf{D}$ , probability  $p \in (0, 1]$ , number of iterations  $K$ 
2: Initialize  $g^0 = \mathbf{D} \cdot \nabla f(x^0)$ 
3: for  $k = 0, 1, \dots, K - 1$  do
4:   Sample  $c_k \sim \text{Be}(p)$ 
5:   Broadcast  $g^k$  to all workers
6:   for  $i = 1, 2, \dots$  in parallel do
7:      $x^{k+1} = x^k - g^k$ 
8:     Set  $g_i^{k+1} = \begin{cases} \mathbf{D} \cdot \nabla f_i(x^{k+1}) & \text{if } c_k = 1 \\ g^k + \mathbf{T}_i^k \mathbf{D} (\nabla f_i(x^{k+1}) - \nabla f_i(x^k)) & \text{if } c_k = 0 \end{cases}$ 
9:   end for
10:   $g^{k+1} = \frac{1}{n} \sum_{i=1}^n g_i^{k+1}$ 
11: end for
12: Return:  $\tilde{x}^K$  chosen uniformly at random from  $\{x^k\}_{k=0}^{K-1}$ 

```

between **det-CGD1** and **det-CGD2** is that in **det-CGD1** the gradient is sketched first and then multiplied by the stepsize, while for **det-CGD2**, the gradient is multiplied by the stepsize first after which the product is sketched. The convergence for Algorithm 4 can be proved in a similar manner as Theorem 1.

Theorem 4. *Let Assumptions 1 and 2 hold, with the gradient of f being L -Lipschitz. If the stepsize matrix $\mathbf{D} \in \mathbb{S}_{++}^d$ satisfies*

$$\mathbf{D}^{-1} \succeq \left(\frac{(1-p) \cdot R'(\mathbf{D}, \mathcal{S})}{np} + 1 \right) \mathbf{L},$$

where

$$R'(\mathbf{D}, \mathcal{S}) = \frac{1}{n} \sum_{i=1}^n \lambda_{\max} \left(\mathbf{D} \mathbb{E} [\mathbf{T}_i^k \mathbf{D}^{-1} \mathbf{T}_i^k] \mathbf{D} \mathbf{L}_i^{\frac{1}{2}} - \mathbf{L}_i^{\frac{1}{2}} \mathbf{D} \right) \cdot \lambda_{\max}(\mathbf{L}_i) \cdot \lambda_{\max} \left(\mathbf{L}^{-\frac{1}{2}} \mathbf{L}_i \mathbf{L}^{-\frac{1}{2}} \right).$$

Then after K iterations of Algorithm 4, we have

$$\mathbb{E} \left[\left\| \nabla f(\tilde{x}^K) \right\|_{\frac{\mathbf{D}}{\det(\mathbf{D})^{1/d}}}^2 \right] \leq \frac{2(f(x^0) - f^*)}{\det(\mathbf{D})^{1/d} \cdot K}.$$

This is to say that in order to reach a ε -stationary point, we require

$$K \geq \frac{2(f(x^0) - f^*)}{\det(\mathbf{D})^{1/d} \cdot \varepsilon^2}.$$

If we look at the scalar case where $\mathbf{D} = \gamma \cdot \mathbf{I}_d$, $\mathbf{L}_i = L_i \cdot \mathbf{I}_d$ and $\mathbf{L} = L \cdot \mathbf{I}_d$, then the condition in Theorem 4 reduces to

$$\frac{(1-p)\omega L^2}{np} + L - \frac{1}{\gamma} \leq 0. \quad (54)$$

Notice that here $\omega = \lambda_{\max} \left(\mathbb{E} \left[(\mathbf{T}_i^k)^2 \right] \right) - 1$, and we have $L^2 = \frac{1}{n} \sum_{i=1}^n L_i^2$, which is due to the relation given in Proposition 5. This condition coincides with the condition for convergence of **MARINA**. One may also check that, the update rule in Algorithm 4, is the same as **MARINA** in the scalar case. However, the condition given in Theorem 4 is not simpler than Theorem 1, contrary to the single-node case. We emphasize that Algorithm 4 is not suitable for the federated learning setting where the clients have limited resources. In order to perform the update, each client is required to store the stepsize matrix \mathbf{D} which is of size $d \times d$. In the over-parameterized regime, the dataset size is $m \times d$ where m is the number of data samples, and we have $d > m$. This means that the stepsize matrix each client needs to store is even larger than the dataset itself, which is unacceptable given the limited resources each client has.

2106 G.2 ANALYSIS OF ALGORITHM 4
2107

2108 We first present two lemmas which are necessary for the proofs of Theorem 4.

2109 **Lemma 6.** Assume that function f is L -smooth, and $x^{k+1} = x^k - g^k$, and matrix $\mathbf{D} \in \mathbb{S}_{++}^d$. Then
2110 we will have

2111
$$f(x^{k+1}) \leq f(x^k) - \frac{1}{2} \|\nabla f(x^k)\|_{\mathbf{D}}^2 + \frac{1}{2} \|\mathbf{D} \cdot \nabla f(x^k) - g^k\|_{\mathbf{D}^{-1}}^2 - \frac{1}{2} \|x^{k+1} - x^k\|_{\mathbf{D}^{-1} - \mathbf{L}}^2. \quad (55)$$

2112

2113 This lemma is formulated in a different way from Lemma 1 on purpose.

2114 **Lemma 7.** For any sketch matrix $\mathbf{T} \in \mathbb{S}_+^d$, vector $t \in \mathbb{R}^d$, matrix $\mathbf{D} \in \mathbb{S}_{++}^d$ and matrix $\mathbf{L} \in \mathbb{S}_{++}^d$,
2115 we have

2116
$$\mathbb{E} \left[\|\mathbf{T}\mathbf{D}t - \mathbf{D}t\|_{\mathbf{D}^{-1}}^2 \right] \leq \lambda_{\max} \left(\mathbf{L}^{\frac{1}{2}} \mathbf{D} \mathbb{E} [\mathbf{T}\mathbf{D}^{-1}\mathbf{T}] \mathbf{D}\mathbf{L}^{\frac{1}{2}} - \mathbf{L}^{\frac{1}{2}} \mathbf{D}\mathbf{L}^{\frac{1}{2}} \right) \|t\|_{\mathbf{L}^{-1}}^2. \quad (56)$$

2117
2118

2119 G.3 PROOF OF THEOREM 4
2120

2121 We start with Lemma 6,

2122
$$\mathbb{E} [f(x^{k+1})] \leq \mathbb{E} [f(x^k)] - \mathbb{E} \left[\frac{1}{2} \|\nabla f(x^k)\|_{\mathbf{D}}^2 \right]$$

2123
$$+ \mathbb{E} \left[\frac{1}{2} \|\mathbf{D} \cdot \nabla f(x^k) - g^k\|_{\mathbf{D}^{-1}}^2 \right] - \mathbb{E} \left[\frac{1}{2} \|x^{k+1} - x^k\|_{\mathbf{D}^{-1} - \mathbf{L}}^2 \right]. \quad (57)$$

2124
2125
2126

2127 Now we do the same as Theorem 1 and look at the term $\mathbb{E} \left[\|\mathbf{D} \cdot \nabla f(x^{k+1}) - g^{k+1}\|_{\mathbf{D}^{-1}}^2 \right]$. Recall
2128 that g^k here is given by

2129
$$g^{k+1} = \begin{cases} \mathbf{D} \cdot \nabla f(x^{k+1}) & \text{with probability } p \\ g^k + \frac{1}{n} \sum_{i=1}^n \mathbf{T}_i^k \mathbf{D} (\nabla f_i(x^{k+1}) - \nabla f_i(x^k)) & \text{with probability } 1 - p. \end{cases}$$

2130
2131

2132 As a result, we have

2133
$$\mathbb{E} \left[\|g^{k+1} - \mathbf{D}\nabla f(x^{k+1})\|_{\mathbf{D}^{-1}}^2 \mid x^{k+1}, x^k \right]$$

2134
$$= \mathbb{E} \left[\mathbb{E} \left[\|g^{k+1} - \mathbf{D}\nabla f(x^{k+1})\|_{\mathbf{D}^{-1}}^2 \mid x^{k+1}, x^k, c_k \right] \right]$$

2135
$$= p \cdot \|\mathbf{D}\nabla f(x^{k+1}) - \mathbf{D}\nabla f(x^{k+1})\|_{\mathbf{D}^{-1}}^2$$

2136
$$+ (1-p) \cdot \mathbb{E} \left[\left\| g^k + \frac{1}{n} \sum_{i=1}^n \mathbf{T}_i^k \mathbf{D} (\nabla f_i(x^{k+1}) - \nabla f_i(x^k)) - \mathbf{D}\nabla f(x^{k+1}) \right\|_{\mathbf{D}^{-1}}^2 \mid x^{k+1}, x^k \right]$$

2137
2138
$$= (1-p) \cdot \mathbb{E} \left[\left\| g^k + \frac{1}{n} \sum_{i=1}^n \mathbf{T}_i^k \mathbf{D} (\nabla f_i(x^{k+1}) - \nabla f_i(x^k)) - \mathbf{D}\nabla f(x^{k+1}) \right\|_{\mathbf{D}^{-1}}^2 \mid x^{k+1}, x^k \right].$$

2139
2140
2141
2142
2143
2144

2145 For the sake of presentation, we use $\mathbb{E}_k[\cdot]$ to denote the conditional expectation $\mathbb{E}[\cdot \mid x_k, x_{k+1}]$ on
2146 x_k, x_{k+1} . Using Fact 2 with $x = \frac{1}{n} \sum_{i=1}^n \mathbf{T}_i^k \mathbf{D} (\nabla f_i(x^{k+1}) - \nabla f_i(x^k))$, $c = \mathbf{D}\nabla f(x^{k+1}) - g^k$,
2147 we are able to obtain that,

2148
$$(1-p)\mathbb{E}_k \left[\left\| g^k + \frac{1}{n} \sum_{i=1}^n \mathbf{T}_i^k \mathbf{D} (\nabla f_i(x^{k+1}) - \nabla f_i(x^k)) - \mathbf{D}\nabla f(x^{k+1}) \right\|_{\mathbf{D}^{-1}}^2 \right]$$

2149
$$= (1-p)\mathbb{E}_k \left[\left\| \frac{1}{n} \sum_{i=1}^n \mathbf{T}_i^k \mathbf{D} (\nabla f_i(x^{k+1}) - \nabla f_i(x^k)) - \mathbf{D} (\nabla f(x^{k+1}) - \nabla f(x^k)) \right\|_{\mathbf{D}^{-1}}^2 \right]$$

2150
$$+ (1-p) \|g^k - \nabla f(x^k)\|_{\mathbf{D}^{-1}}^2$$

2151
$$= (1-p)\mathbb{E}_k \left[\left\| \frac{1}{n} \sum_{i=1}^n [\mathbf{T}_i^k \mathbf{D} (\nabla f_i(x^{k+1}) - \nabla f_i(x^k)) - \mathbf{D} (\nabla f_i(x^{k+1}) - \nabla f_i(x^k))] \right\|_{\mathbf{D}^{-1}}^2 \right]$$

2152
$$+ (1-p) \|g^k - \nabla f(x^k)\|_{\mathbf{D}^{-1}}^2.$$

2153
2154
2155
2156
2157
2158
2159

It is not hard to notice that for the sketch matrices we pick, the following identity holds due to the unbiasedness,

$$\mathbb{E}_k [\mathbf{T}_i^k \mathbf{D} (\nabla f_i(x^{k+1}) - \nabla f_i(x^k))] = \mathbf{D} (\nabla f_i(x^{k+1}) - \nabla f_i(x^k)),$$

and any two random vectors in the set $\{\mathbf{T}_i^k \mathbf{D} (\nabla f_i(x^{k+1}) - \nabla f_i(x^k))\}_{i=1}^n$ are independent if x^{k+1}, x^k are fixed. As a result

$$\begin{aligned} & \mathbb{E}_k \left[\|g^{k+1} - \mathbf{D} \nabla f(x^{k+1})\|_{\mathbf{D}^{-1}}^2 \right] \\ &= \frac{1-p}{n^2} \sum_{i=1}^n \mathbb{E}_k \left[\|\mathbf{T}_i^k (\mathbf{D} \nabla f_i(x^{k+1}) - \mathbf{D} \nabla f_i(x^k)) - (\mathbf{D} \nabla f_i(x^{k+1}) - \mathbf{D} \nabla f_i(x^k))\|_{\mathbf{D}^{-1}}^2 \right] \\ & \quad + (1-p) \cdot \|g^k - \mathbf{D} \nabla f(x^k)\|_{\mathbf{D}^{-1}}^2. \end{aligned} \quad (58)$$

For each term within the summation, we can further upper bound it using Lemma 7

$$\begin{aligned} & \mathbb{E}_k \left[\|\mathbf{T}_i^k (\mathbf{D} \nabla f_i(x^{k+1}) - \mathbf{D} \nabla f_i(x^k)) - (\mathbf{D} \nabla f_i(x^{k+1}) - \mathbf{D} \nabla f_i(x^k))\|_{\mathbf{D}^{-1}}^2 \right] \\ & \leq \lambda_{\max} \left(\mathbf{L}_i^{\frac{1}{2}} \mathbf{D} \mathbb{E} [\mathbf{T}_i^k \mathbf{D}^{-1} \mathbf{T}_i^k] \mathbf{D} \mathbf{L}_i^{\frac{1}{2}} - \mathbf{L}_i^{\frac{1}{2}} \mathbf{D} \mathbf{L}_i^{\frac{1}{2}} \right) \|\nabla f_i(x^{k+1}) - \nabla f_i(x^k)\|_{\mathbf{L}_i^{-1}}^2 \\ & \leq \lambda_{\max} \left(\mathbf{L}_i^{\frac{1}{2}} \mathbf{D} \mathbb{E} [\mathbf{T}_i^k \mathbf{D}^{-1} \mathbf{T}_i^k] \mathbf{D} \mathbf{L}_i^{\frac{1}{2}} - \mathbf{L}_i^{\frac{1}{2}} \mathbf{D} \mathbf{L}_i^{\frac{1}{2}} \right) \|x^{k+1} - x^k\|_{\mathbf{L}_i}^2. \end{aligned}$$

Where the last inequality is due to Assumption 2. Plugging back into (58), we get

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

Applying the replacement trick from the proof of Theorem 1, we obtain

$$\begin{aligned} & \mathbb{E}_k \left[\|g^{k+1} - \mathbf{D} \nabla f(x^{k+1})\|_{\mathbf{D}^{-1}}^2 \right] \\ & \leq \frac{1-p}{n^2} \sum_{i=1}^n \lambda_{\max} \left(\mathbf{L}_i^{\frac{1}{2}} \mathbf{D} \mathbb{E} [\mathbf{T}_i^k \mathbf{D}^{-1} \mathbf{T}_i^k] \mathbf{D} \mathbf{L}_i^{\frac{1}{2}} - \mathbf{L}_i^{\frac{1}{2}} \mathbf{D} \mathbf{L}_i^{\frac{1}{2}} \right) \\ & \quad \times \left\langle \mathbf{L}^{\frac{1}{2}} (x^{k+1} - x^k), \left(\mathbf{L}^{-\frac{1}{2}} \mathbf{L}_i \mathbf{L}^{-\frac{1}{2}} \right) \cdot \mathbf{L}^{\frac{1}{2}} (x^{k+1} - x^k) \right\rangle + (1-p) \cdot \|g^k - \mathbf{D} \nabla f(x^k)\|_{\mathbf{D}^{-1}}^2 \\ & \leq \frac{1-p}{n^2} \sum_{i=1}^n \lambda_{\max} \left(\mathbf{L}_i^{\frac{1}{2}} (\mathbf{D} \mathbb{E} [\mathbf{T}_i^k \mathbf{D}^{-1} \mathbf{T}_i^k] \mathbf{D} - \mathbf{D}) \mathbf{L}_i^{\frac{1}{2}} \right) \cdot \lambda_{\max} \left(\mathbf{L}^{-\frac{1}{2}} \mathbf{L}_i \mathbf{L}^{-\frac{1}{2}} \right) \|x^{k+1} - x^k\|_{\mathbf{L}}^2 \\ & \quad + (1-p) \cdot \|g^k - \mathbf{D} \nabla f(x^k)\|_{\mathbf{D}^{-1}}^2. \end{aligned}$$

Applying Fact 5, we obtain

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

Recalling the definition of $R'(\mathbf{D}, \mathcal{S})$, we further simplify it to

$$\begin{aligned} & \mathbb{E}_k \left[\|g^{k+1} - \mathbf{D} \nabla f(x^{k+1})\|_{\mathbf{D}^{-1}}^2 \right] \\ & \leq \frac{(1-p) \cdot R'(\mathbf{D}, \mathcal{S})}{n} \|x^{k+1} - x^k\|_{\mathbf{L}}^2 + (1-p) \cdot \|g^k - \mathbf{D} \nabla f(x^k)\|_{\mathbf{D}^{-1}}^2. \end{aligned}$$

2214 Taking expectation again and using the tower property, we get

$$2215 \mathbb{E} \left[\|g^{k+1} - \mathbf{D}\nabla f(x^{k+1})\|_{\mathbf{D}^{-1}}^2 \right] \quad (59)$$

$$2216 \leq \frac{(1-p) \cdot R'(\mathbf{D}, \mathcal{S})}{n} \mathbb{E} \left[\|x^{k+1} - x^k\|_{\mathbf{L}}^2 \right] + (1-p) \cdot \mathbb{E} \left[\|g^k - \mathbf{D}\nabla f(x^k)\|_{\mathbf{D}^{-1}}^2 \right]. \quad (60)$$

2217
2218
2219 Construct the Lyapunov function Φ_k as follows,

$$2220 \Phi_k = f(x^k) - f^* + \frac{1}{2p} \|g^k - \mathbf{D}\nabla f(x^k)\|_{\mathbf{D}^{-1}}^2.$$

2221 Utilizing (57) and (59), we are able to get

$$\begin{aligned} 2222 \mathbb{E} [\Phi_{k+1}] &\leq \mathbb{E} [f(x^k) - f^*] - \frac{1}{2} \mathbb{E} \left[\|\nabla f(x^k)\|_{\mathbf{D}}^2 \right] \\ 2223 &\quad + \frac{1}{2} \mathbb{E} \left[\|g^k - \mathbf{D}\nabla f(x^k)\|_{\mathbf{D}^{-1}}^2 \right] - \frac{1}{2} \mathbb{E} \left[\|x^{k+1} - x^k\|_{\mathbf{D}^{-1}-\mathbf{L}}^2 \right] \\ 2224 &\quad + \frac{1}{2p} \cdot \frac{(1-p)R'(\mathbf{D}, \mathcal{S})}{n} \mathbb{E} \left[\|x^{k+1} - x^k\|_{\mathbf{L}}^2 \right] + \frac{1-p}{2p} \mathbb{E} \left[\|g^k - \mathbf{D}\nabla f(x^k)\|_{\mathbf{D}^{-1}}^2 \right] \\ 2225 &= \mathbb{E} [\Phi_k] - \frac{1}{2} \mathbb{E} \left[\|\nabla f(x^k)\|_{\mathbf{D}}^2 \right] \\ 2226 &\quad + \frac{1}{2} \left(\frac{(1-p)R'(\mathbf{D}, \mathcal{S})}{np} \mathbb{E} \left[\|x^{k+1} - x^k\|_{\mathbf{L}}^2 \right] - \mathbb{E} \left[\|x^{k+1} - x^k\|_{\mathbf{D}^{-1}-\mathbf{L}}^2 \right] \right). \end{aligned}$$

2227 Now, notice that the last term in the above inequality is non-positive as guaranteed by the condition

$$2228 \mathbf{D}^{-1} \succeq \left(\frac{(1-p)R'(\mathbf{D}, \mathcal{S})}{np} + 1 \right) \mathbf{L}.$$

2229 This leads to the recurrence after ignoring the last term,

$$2230 \mathbb{E} [\Phi_{k+1}] \leq \mathbb{E} [\Phi_k] - \frac{1}{2} \mathbb{E} \left[\|\nabla f(x^k)\|_{\mathbf{D}}^2 \right].$$

2231 Unrolling this recurrence, we get

$$2232 \frac{1}{K} \sum_{k=0}^{K-1} \mathbb{E} \left[\|\nabla f(x^k)\|_{\mathbf{D}}^2 \right] \leq \frac{2(\mathbb{E} [\Phi_0] - \mathbb{E} [\Phi_K])}{K}.$$

2233 The left hand side can be viewed as average over \tilde{x}^K , which is drawn uniformly at random from $\{x_k\}_{k=0}^{K-1}$, while the right hand side can be simplified as

$$2234 \frac{2(\mathbb{E} [\Phi_0] - \mathbb{E} [\Phi_K])}{K} \leq \frac{2\Phi_0}{K} = \frac{2 \left(f(x^0) - f^* + \frac{1}{2p} \|g^0 - \nabla f(x^0)\|_{\mathbf{D}}^2 \right)}{K}.$$

2235 Recalling that $g^0 = \nabla f(x^0)$ and performing determinant normalization as Li et al. (2024b), we get

$$2236 \mathbb{E} \left[\|\nabla f(\tilde{x}^K)\|_{\frac{\mathbf{D}}{\det(\mathbf{D})^{1/d}}}^2 \right] \leq \frac{2(f(x^0) - f^*)}{\det(\mathbf{D})^{1/d} K}.$$

2237 H PROOFS OF THE TECHNICAL LEMMAS

2238 H.1 PROOF OF LEMMA 1

2239 Let $\bar{x}^{k+1} := x^k - \mathbf{D} \cdot \nabla f(x^k)$. Since f has a matrix \mathbf{L} -Lipschitz gradient, f is also \mathbf{L} -smooth. From the \mathbf{L} -smoothness of f , we have

$$\begin{aligned} 2240 f(x^{k+1}) &\leq f(x^k) + \langle \nabla f(x^k), x^{k+1} - x^k \rangle + \frac{1}{2} \langle x^{k+1} - x^k, \mathbf{L}(x^{k+1} - x^k) \rangle \\ 2241 &= f(x^k) + \langle \nabla f(x^k) - g^k, x^{k+1} - x^k \rangle + \langle g^k, x^{k+1} - x^k \rangle + \frac{1}{2} \langle x^{k+1} - x^k, \mathbf{L}(x^{k+1} - x^k) \rangle. \end{aligned}$$

2268
2269
2270
2271
2272
2273
2274
2275

We can merge the last two terms and obtain,

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

2276
2277

We add and subtract $\langle \nabla f(x^k) - g^k, \mathbf{D} \cdot g^k \rangle$,

2278
2279
2280
2281
2282
2283
2284
2285

$$\begin{aligned} f(x^{k+1}) &\leq f(x^k) + \langle \nabla f(x^k) - g^k, \mathbf{D}(\nabla f(x^k) - g^k) \rangle - \langle \nabla f(x^k) - g^k, \mathbf{D} \cdot \nabla f(x^k) \rangle \\ &\quad - \left\langle x^{k+1} - x^k, \left(\mathbf{D}^{-1} - \frac{1}{2} \mathbf{L} \right) (x^{k+1} - x^k) \right\rangle \\ &= f(x^k) + \|\nabla f(x^k) - g^k\|_{\mathbf{D}}^2 - \langle x^{k+1} - \bar{x}^{k+1}, \mathbf{D}^{-1}(x^k - \bar{x}^{k+1}) \rangle \\ &\quad - \left\langle x^{k+1} - x^k, \left(\mathbf{D}^{-1} - \frac{1}{2} \mathbf{L} \right) (x^{k+1} - x^k) \right\rangle. \end{aligned}$$

2286
2287

Decomposing the term $\langle x^{k+1} - \bar{x}^{k+1}, \mathbf{D}^{-1}(x^k - \bar{x}^{k+1}) \rangle$, we get

2288
2289
2290
2291
2292

$$\begin{aligned} f(x^{k+1}) &\leq f(x^k) + \|\nabla f(x^k) - g^k\|_{\mathbf{D}}^2 - \left\langle x^{k+1} - x^k, \left(\mathbf{D}^{-1} - \frac{1}{2} \mathbf{L} \right) (x^{k+1} - x^k) \right\rangle \\ &\quad - \frac{1}{2} \left(\|x^{k+1} - \bar{x}^{k+1}\|_{\mathbf{D}^{-1}}^2 + \|x^k - \bar{x}^{k+1}\|_{\mathbf{D}^{-1}}^2 - \|x^{k+1} - x^k\|_{\mathbf{D}^{-1}}^2 \right). \end{aligned}$$

2293

Plugging in the definition of x^{k+1}, \bar{x}^{k+1} , we get

2294
2295
2296
2297
2298
2299
2300
2301

$$\begin{aligned} f(x^{k+1}) &\leq f(x^k) + \|\nabla f(x^k) - g^k\|_{\mathbf{D}}^2 - \|x^{k+1} - x^k\|_{\mathbf{D}^{-1} - \frac{1}{2} \mathbf{L}}^2 \\ &\quad - \frac{1}{2} \left(\|\mathbf{D}(\nabla f(x^k) - g^k)\|_{\mathbf{D}^{-1}}^2 + \|\mathbf{D} \cdot \nabla f(x^k)\|_{\mathbf{D}^{-1}}^2 - \|x^{k+1} - x^k\|_{\mathbf{D}^{-1}}^2 \right) \\ &= f(x^k) + \|\nabla f(x^k) - g^k\|_{\mathbf{D}}^2 - \|x^{k+1} - x^k\|_{\mathbf{D}^{-1} - \frac{1}{2} \mathbf{L}}^2 \\ &\quad - \frac{1}{2} \left(\|\nabla f(x^k) - g^k\|_{\mathbf{D}}^2 + \|\nabla f(x^k)\|_{\mathbf{D}}^2 - \|x^{k+1} - x^k\|_{\mathbf{D}^{-1}}^2 \right). \end{aligned}$$

2302

Rearranging terms we get,

2303
2304
2305
2306
2307
2308

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

2309

H.2 PROOF OF LEMMA 2

2310

The definition of the weighted norm yields

2311
2312
2313
2314
2315
2316
2317
2318
2319
2320
2321

$$\begin{aligned} \mathbb{E} \left[\|S t - t\|_{\mathbf{D}}^2 \right] &= \mathbb{E} [\langle t, (\mathbf{S} - \mathbf{I}_d) \mathbf{D} (\mathbf{S} - \mathbf{I}_d) t \rangle] \\ &= \langle t, \mathbb{E} [(\mathbf{S} - \mathbf{I}_d) \mathbf{D} (\mathbf{S} - \mathbf{I}_d)] t \rangle \\ &= \left\langle t, \mathbf{L}^{-\frac{1}{2}} \cdot \mathbb{E} \left[\mathbf{L}^{\frac{1}{2}} (\mathbf{S} - \mathbf{I}_d) \mathbf{D} (\mathbf{S} - \mathbf{I}_d) \mathbf{L}^{\frac{1}{2}} \right] \cdot \mathbf{L}^{-\frac{1}{2}} t \right\rangle \\ &= \left\langle \mathbf{L}^{-\frac{1}{2}} t, \mathbb{E} \left[\mathbf{L}^{\frac{1}{2}} (\mathbf{S} - \mathbf{I}_d) \mathbf{D} (\mathbf{S} - \mathbf{I}_d) \mathbf{L}^{\frac{1}{2}} \right] \cdot \mathbf{L}^{-\frac{1}{2}} t \right\rangle \\ &\leq \lambda_{\max} \left(\mathbb{E} \left[\mathbf{L}^{\frac{1}{2}} (\mathbf{S} - \mathbf{I}_d) \mathbf{D} (\mathbf{S} - \mathbf{I}_d) \mathbf{L}^{\frac{1}{2}} \right] \right) \left\| \mathbf{L}^{-\frac{1}{2}} t \right\|^2 \\ &= \lambda_{\max} \left(\mathbf{L}^{\frac{1}{2}} (\mathbb{E} [\mathbf{S} \mathbf{D} \mathbf{S}] - \mathbf{D}) \mathbf{L}^{\frac{1}{2}} \right) \cdot \|t\|_{\mathbf{L}^{-1}}^2. \end{aligned}$$

2322 H.3 PROOF OF LEMMA 4
2323

2324 Throughout the following proof, we denote $\mathbb{E}_{\mathcal{S}}[\cdot]$ as taking expectation with respect to the
2325 randomness contained within the sketch sampled from distribution \mathcal{S} . We estimate the term
2326 $\mathbb{E}_{\mathcal{S}}\left[\|g^{k+1} - h^{k+1}\|_{\mathcal{D}}^2\right]$ in order to construct the Lyapunov function. For $\mathbb{E}_{\mathcal{S}}\left[\|g^{k+1} - h^{k+1}\|_{\mathcal{D}}^2\right]$,
2327 we have

$$\begin{aligned} 2328 \mathbb{E}_{\mathcal{S}}\left[\|g^{k+1} - h^{k+1}\|_{\mathcal{D}}^2\right] &= \mathbb{E}_{\mathcal{S}}\left[\left\|g^k + \frac{1}{n} \sum_{i=1}^n m_i^{k+1} - h^{k+1}\right\|_{\mathcal{D}}^2\right] \\ 2329 &= \mathbb{E}_{\mathcal{S}}\left[\left\|g^k + \frac{1}{n} \sum_{i=1}^n \mathbf{S}_i^k (h_i^{k+1} - h_i^k - a(g_i^k - h_i^k)) - h^{k+1}\right\|_{\mathcal{D}}^2\right] \end{aligned}$$

2330 Using Fact 3, we obtain

$$\begin{aligned} 2331 \mathbb{E}_{\mathcal{S}}\left[\|g^{k+1} - h^{k+1}\|_{\mathcal{D}}^2\right] &= \mathbb{E}_{\mathcal{S}}\left[\left\|\frac{1}{n} \sum_{i=1}^n \mathbf{S}_i^k (h_i^{k+1} - h_i^k - a(g_i^k - h_i^k)) - (h^{k+1} - h^k - a(g^k - h^k))\right\|_{\mathcal{D}}^2\right] \\ 2332 &\quad + (1-a)^2 \|h^k - g^k\|_{\mathcal{D}}^2 \\ 2333 &= \mathbb{E}_{\mathcal{S}}\left[\left\|\frac{1}{n} \sum_{i=1}^n \mathbf{S}_i^k (h_i^{k+1} - h_i^k - a(g_i^k - h_i^k)) - \frac{1}{n} \sum_{i=1}^n (h_i^{k+1} - h_i^k - a(g_i^k - h_i^k))\right\|_{\mathcal{D}}^2\right] \\ 2334 &\quad + (1-a)^2 \|h^k - g^k\|_{\mathcal{D}}^2 \\ 2335 &= \frac{1}{n^2} \sum_{i=1}^n \mathbb{E}_{\mathcal{S}}\left[\left\|\mathbf{S}_i^k (h_i^{k+1} - h_i^k - a(g_i^k - h_i^k)) - (h_i^{k+1} - h_i^k - a(g_i^k - h_i^k))\right\|_{\mathcal{D}}^2\right] \\ 2336 &\quad + (1-a)^2 \|h^k - g^k\|_{\mathcal{D}}^2. \end{aligned}$$

2337 Here, the last identity is obtained from the unbiasedness of the sketches:

$$2338 \mathbb{E}_{\mathcal{S}}[\mathbf{S}_i^k (h_i^{k+1} - h_i^k - a(g_i^k - h_i^k))] = h_i^{k+1} - h_i^k - a(g_i^k - h_i^k).$$

2339 We can further use Lemma 2, and obtain

$$\begin{aligned} 2340 \mathbb{E}_{\mathcal{S}}\left[\|g^{k+1} - h^{k+1}\|_{\mathcal{D}}^2\right] &\leq \frac{1}{n^2} \sum_{i=1}^n \lambda_{\max}\left(\mathbf{D}^{-\frac{1}{2}} (\mathbb{E}[\mathbf{S}_i^k \mathbf{D} \mathbf{S}_i^k] - \mathbf{D}) \mathbf{D}^{-\frac{1}{2}}\right) \|h_i^{k+1} - h_i^k - a(g_i^k - h_i^k)\|_{\mathcal{D}}^2 \\ 2341 &\quad + (1-a)^2 \|g^k - h^k\|_{\mathcal{D}}^2 \\ 2342 &\leq \frac{1}{n^2} \sum_{i=1}^n \lambda_{\max}(\mathbf{D}^{-1}) \cdot \lambda_{\max}(\mathbb{E}[\mathbf{S}_i^k \mathbf{D} \mathbf{S}_i^k] - \mathbf{D}) \|h_i^{k+1} - h_i^k - a(g_i^k - h_i^k)\|_{\mathcal{D}}^2 \\ 2343 &\quad + (1-a)^2 \|g^k - h^k\|_{\mathcal{D}}^2. \end{aligned}$$

2344 We can rewrite the above bound, after applying Jensen's inequality as

$$\begin{aligned} 2345 \mathbb{E}_{\mathcal{S}}\left[\|g^{k+1} - h^{k+1}\|_{\mathcal{D}}^2\right] &\leq \frac{2\Lambda_{\mathcal{D},\mathcal{S}} \cdot \lambda_{\max}(\mathbf{D}^{-1})}{n^2} \sum_{i=1}^n \|h_i^{k+1} - h_i^k\|_{\mathcal{D}}^2 \\ 2346 &\quad + \frac{2a^2 \Lambda_{\mathcal{D},\mathcal{S}} \cdot \lambda_{\max}(\mathbf{D}^{-1})}{n^2} \sum_{i=1}^n \|g_i^k - h_i^k\|_{\mathcal{D}}^2 \\ 2347 &\quad + (1-a)^2 \|g^k - h^k\|_{\mathcal{D}}^2. \end{aligned}$$

2376 Notice that we have

$$2377 \quad \|h_i^{k+1} - h_i^k\|_{\mathbf{D}}^2 \leq \lambda_{\max}(\mathbf{D}) \cdot \lambda_{\max}(\mathbf{L}_i) \cdot \|h_i^{k+1} - h_i^k\|_{\mathbf{L}_i^{-1}}^2.$$

2379 Thus, it is not hard to see that

$$2380 \quad \mathbb{E}_{\mathcal{S}} \left[\|g^{k+1} - h^{k+1}\|_{\mathbf{D}}^2 \right] \leq \frac{2\Lambda_{\mathbf{D},\mathcal{S}} \cdot \lambda_{\max}(\mathbf{D}^{-1}) \cdot \lambda_{\max}(\mathbf{D})}{n^2} \sum_{i=1}^n \lambda_{\max}(\mathbf{L}_i) \|h_i^{k+1} - h_i^k\|_{\mathbf{L}_i^{-1}}^2$$

$$2381 \quad + \frac{2a^2\Lambda_{\mathbf{D},\mathcal{S}} \cdot \lambda_{\max}(\mathbf{D}^{-1})}{n^2} \sum_{i=1}^n \|g_i^k - h_i^k\|_{\mathbf{D}}^2$$

$$2382 \quad + (1-a)^2 \|g^k - h^k\|_{\mathbf{D}}^2.$$

2383 We obtain the inequality in the lemma after taking expectation again and applying tower property.

2390 H.4 PROOF OF LEMMA 5

2391 Similarly, we then try to bound the terms $\mathbb{E}_{\mathcal{S}} \left[\|g_i^{k+1} - h_i^{k+1}\|_{\mathbf{D}}^2 \right]$. We start with

$$2392 \quad \mathbb{E}_{\mathcal{S}} \left[\|g_i^{k+1} - h_i^{k+1}\|_{\mathbf{D}}^2 \right]$$

$$2393 \quad = \mathbb{E}_{\mathcal{S}} \left[\|g_i^k + \mathbf{S}_i^k (h_i^{k+1} - h_i^k - a(g_i^k - h_i^k)) - h_i^{k+1}\|_{\mathbf{D}}^2 \right]$$

$$2394 \quad = \mathbb{E}_{\mathcal{S}} \left[\|\mathbf{S}_i^k (h_i^{k+1} - h_i^k - a(g_i^k - h_i^k)) - (h_i^{k+1} - h_i^k - a(g_i^k - h_i^k)) + (1-a)(h_i^k - g_i^k)\|_{\mathbf{D}}^2 \right].$$

2399 Using Fact 3,

$$2400 \quad \mathbb{E}_{\mathcal{S}} \left[\|g_i^{k+1} - h_i^{k+1}\|_{\mathbf{D}}^2 \right]$$

$$2401 \quad = \mathbb{E}_{\mathcal{S}} \left[\|\mathbf{S}_i^k (h_i^{k+1} - h_i^k - a(g_i^k - h_i^k)) - (h_i^{k+1} - h_i^k - a(g_i^k - h_i^k))\|_{\mathbf{D}}^2 \right]$$

$$2402 \quad + (1-a)^2 \|h_i^k - g_i^k\|_{\mathbf{D}}^2.$$

2406 Using Lemma 2

$$2407 \quad \mathbb{E}_{\mathcal{S}} \left[\|g_i^{k+1} - h_i^{k+1}\|_{\mathbf{D}}^2 \right]$$

$$2408 \quad \stackrel{(22)}{\leq} \lambda_{\max}(\mathbf{D}^{-\frac{1}{2}} (\mathbb{E} [\mathbf{S}_i^k \mathbf{D} \mathbf{S}_i^k] - \mathbf{D}) \mathbf{D}^{-\frac{1}{2}}) \|h_i^{k+1} - h_i^k - a(g_i^k - h_i^k)\|_{\mathbf{D}}^2$$

$$2409 \quad + (1-a)^2 \|g_i^k - h_i^k\|_{\mathbf{D}}^2$$

$$2410 \quad \leq \lambda_{\max}(\mathbf{D}^{-1}) \cdot \Lambda_{\mathbf{D},\mathcal{S}} \|h_i^{k+1} - h_i^k - a(g_i^k - h_i^k)\|_{\mathbf{D}}^2 + (1-a)^2 \|g_i^k - h_i^k\|_{\mathbf{D}}^2$$

$$2411 \quad \leq 2\lambda_{\max}(\mathbf{D}^{-1}) \cdot \Lambda_{\mathbf{D},\mathcal{S}} \|h_i^{k+1} - h_i^k\|_{\mathbf{D}}^2 + 2a^2\lambda_{\max}(\mathbf{D}^{-1}) \cdot \Lambda_{\mathbf{D},\mathcal{S}} \|g_i^k - h_i^k\|_{\mathbf{D}}^2$$

$$2412 \quad + (1-a)^2 \|g_i^k - h_i^k\|_{\mathbf{D}}^2$$

$$2413 \quad \leq 2\lambda_{\max}(\mathbf{D}^{-1}) \cdot \lambda_{\max}(\mathbf{D}) \cdot \Lambda_{\mathbf{D},\mathcal{S}} \cdot \lambda_{\max}(\mathbf{L}_i) \cdot \|h_i^{k+1} - h_i^k\|_{\mathbf{L}_i^{-1}}^2$$

$$2414 \quad + 2a^2\lambda_{\max}(\mathbf{D}^{-1}) \cdot \Lambda_{\mathbf{D},\mathcal{S}} \|g_i^k - h_i^k\|_{\mathbf{D}}^2 + (1-a)^2 \|g_i^k - h_i^k\|_{\mathbf{D}}^2$$

$$2415 \quad = (2a^2\lambda_{\max}(\mathbf{D}^{-1}) \cdot \Lambda_{\mathbf{D},\mathcal{S}} + (1-a)^2) \|g_i^k - h_i^k\|_{\mathbf{D}}^2$$

$$2416 \quad + 2\lambda_{\max}(\mathbf{D}^{-1}) \cdot \lambda_{\max}(\mathbf{D}) \cdot \Lambda_{\mathbf{D},\mathcal{S}} \cdot \lambda_{\max}(\mathbf{L}_i) \cdot \|h_i^{k+1} - h_i^k\|_{\mathbf{L}_i^{-1}}^2.$$

2424 Taking expectation again, and using tower property, we are able to obtain,

$$2425 \quad \mathbb{E} \left[\|g_i^{k+1} - h_i^{k+1}\|_{\mathbf{D}}^2 \right]$$

$$2426 \quad \leq (2a^2\lambda_{\max}(\mathbf{D}^{-1}) \cdot \Lambda_{\mathbf{D},\mathcal{S}} + (1-a)^2) \mathbb{E} \left[\|g_i^k - h_i^k\|_{\mathbf{D}}^2 \right]$$

$$2427 \quad + 2\lambda_{\max}(\mathbf{D}^{-1}) \cdot \lambda_{\max}(\mathbf{D}) \cdot \Lambda_{\mathbf{D},\mathcal{S}} \cdot \lambda_{\max}(\mathbf{L}_i) \cdot \mathbb{E} \left[\|h_i^{k+1} - h_i^k\|_{\mathbf{L}_i^{-1}}^2 \right].$$

2430
2431
2432
2433
2434
2435
2436
2437
2438
2439
2440
2441
2442
2443
2444
2445
2446
2447
2448
2449
2450
2451
2452
2453
2454
2455
2456
2457
2458
2459
2460
2461
2462
2463
2464
2465
2466
2467
2468
2469
2470
2471
2472
2473
2474
2475
2476
2477
2478
2479
2480
2481
2482
2483

H.5 PROOF OF LEMMA 6

From Proposition 5, we know that the objective is L -smooth. Let $\bar{x}^{k+1} = x^k - \mathbf{D} \cdot \nabla f(x^k)$, then L -smoothness yields

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

Simplifying the above inner-products we have,

$$f(x^{k+1}) \leq f(x^k) + \langle \nabla f(x^k) - \mathbf{D}^{-1} \cdot g^k, -g^k \rangle - \left\langle x^{k+1} - x^k, \left(\mathbf{D}^{-1} - \frac{1}{2} \mathbf{L} \right) (x^{k+1} - x^k) \right\rangle.$$

We then add and subtract $\langle \nabla f(x^k) - \mathbf{D}^{-1} \cdot g^k, \mathbf{D} \cdot \nabla f(x^k) - g^k \rangle$, which

$$\begin{aligned} f(x^{k+1}) &\leq f(x^k) + \langle \nabla f(x^k) - \mathbf{D}^{-1} \cdot g^k, \mathbf{D} \cdot \nabla f(x^k) - g^k \rangle - \langle \nabla f(x^k) - \mathbf{D}^{-1} \cdot g^k, \mathbf{D} \cdot \nabla f(x^k) \rangle \\ &\quad - \left\langle x^{k+1} - x^k, \left(\mathbf{D}^{-1} - \frac{1}{2} \mathbf{L} \right) (x^{k+1} - x^k) \right\rangle \\ &= f(x^k) + \|\nabla f(x^k) - \mathbf{D}^{-1} \cdot g^k\|_{\mathbf{D}}^2 - \langle \mathbf{D}^{-1}(x^{k+1} - \bar{x}^{k+1}), x^k - \bar{x}^{k+1} \rangle \\ &\quad - \left\langle x^{k+1} - x^k, \left(\mathbf{D}^{-1} - \frac{1}{2} \mathbf{L} \right) (x^{k+1} - x^k) \right\rangle. \end{aligned}$$

Decomposing the inner product term we deduce,

$$\begin{aligned} f(x^{k+1}) &\leq f(x^k) + \|\mathbf{D}^{-1} (\mathbf{D} \cdot \nabla f(x^k) - g^k)\|_{\mathbf{D}}^2 - \left\langle x^{k+1} - x^k, \left(\mathbf{D}^{-1} - \frac{1}{2} \mathbf{L} \right) (x^{k+1} - x^k) \right\rangle \\ &\quad - \frac{1}{2} \left(\|x^{k+1} - \bar{x}^{k+1}\|_{\mathbf{D}^{-1}}^2 + \|x^k - \bar{x}^{k+1}\|_{\mathbf{D}^{-1}}^2 - \|x^{k+1} - x^k\|_{\mathbf{D}^{-1}}^2 \right) \\ &= f(x^k) + \|\mathbf{D} \cdot \nabla f(x^k) - g^k\|_{\mathbf{D}^{-1}}^2 - \|x^{k+1} - x^k\|_{\mathbf{D}^{-1} - \frac{1}{2} \mathbf{L}}^2 \\ &\quad - \frac{1}{2} \left(\|\mathbf{D} \cdot \nabla f(x^k) - g^k\|_{\mathbf{D}^{-1}}^2 + \|\mathbf{D} \cdot \nabla f(x^k)\|_{\mathbf{D}^{-1}}^2 - \|x^{k+1} - x^k\|_{\mathbf{D}^{-1}}^2 \right). \end{aligned}$$

Therefore,

$$f(x^{k+1}) \leq f(x^k) + \frac{1}{2} \|\mathbf{D} \nabla f(x^k) - g^k\|_{\mathbf{D}^{-1}}^2 - \frac{1}{2} \|\nabla f(x^k)\|_{\mathbf{D}}^2 - \frac{1}{2} \|x^{k+1} - x^k\|_{\mathbf{D}^{-1} - \mathbf{L}}^2.$$

H.6 PROOF OF LEMMA 7

We start with

$$\begin{aligned} \mathbb{E} \left[\|\mathbf{T} \mathbf{D} t - \mathbf{D} t\|_{\mathbf{D}^{-1}}^2 \right] &= \mathbb{E} \left[\|\mathbf{T} - \mathbf{I}_d\|_{\mathbf{D}^{-1}}^2 \|t\|_{\mathbf{D}^{-1}}^2 \right] \\ &= \langle t, \mathbb{E} [\mathbf{D}(\mathbf{T} - \mathbf{I}_d) \mathbf{D}^{-1} (\mathbf{T} - \mathbf{I}_d) \mathbf{D}] \cdot t \rangle \\ &= \langle t, \mathbf{D} (\mathbb{E} [\mathbf{T} \mathbf{D}^{-1} \mathbf{T}] - \mathbf{D}^{-1}) \mathbf{D} \cdot t \rangle \\ &= \left\langle \mathbf{L}^{-\frac{1}{2}} t, \mathbf{L}^{\frac{1}{2}} \mathbf{D} (\mathbb{E} [\mathbf{T} \mathbf{D}^{-1} \mathbf{T}] - \mathbf{D}^{-1}) \mathbf{D} \mathbf{L}^{\frac{1}{2}} \cdot \mathbf{L}^{-\frac{1}{2}} t \right\rangle \\ &\leq \lambda_{\max} \left(\mathbf{L}^{\frac{1}{2}} \mathbf{D} \mathbb{E} [\mathbf{T} \mathbf{D}^{-1} \mathbf{T}] \mathbf{D} \mathbf{L}^{\frac{1}{2}} - \mathbf{L}^{\frac{1}{2}} \mathbf{D} \mathbf{L}^{\frac{1}{2}} \right) \cdot \|\mathbf{L}^{-\frac{1}{2}} t\|^2 \\ &= \lambda_{\max} \left(\mathbf{L}^{\frac{1}{2}} \mathbf{D} \mathbb{E} [\mathbf{T} \mathbf{D}^{-1} \mathbf{T}] \mathbf{D} \mathbf{L}^{\frac{1}{2}} - \mathbf{L}^{\frac{1}{2}} \mathbf{D} \mathbf{L}^{\frac{1}{2}} \right) \cdot \|t\|_{\mathbf{L}^{-1}}^2 \end{aligned}$$

This completes the proof.

I EXPERIMENTS

In this section, we conduct numerical experiments to back up the theoretical results for **det-MARINA** and **det-DASHA**. The code for the experiments can be found in <https://anonymous.4open.science/r/detCGD-VR-Code-865B>. All the codes for the experiments are written in Python 3.11 with NumPy and SciPy package. The code was run on a machine with AMD Ryzen 9 5900HX Radeon Graphics @ 3.3 GHz and 8 cores 16 threads. The datasets in LibSVM are typically non-IID real world datasets, and it is randomly distributed across all the clients.

I.1 THE SETTING

We first state the experiment setting. We are interested in the following logistic regression problem with a non-convex regularizer. The objective is given as

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

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 whose size is m_i . The constant $\lambda > 0$ is the coefficient of the regularizer. Larger λ means the model is more regular. For each function f_i , its Hessian can be upper bounded by

$$\mathbf{L}_i = \frac{1}{m_i} \sum_{i=1}^{m_i} \frac{a_i a_i^\top}{4} + 2\lambda \cdot \mathbf{I}_d;$$

and, therefore, the Hessian of f is bounded by

$$\mathbf{L} = \frac{1}{\sum_{i=1}^n m_i} \sum_{i=1}^n \sum_{j=1}^{m_i} \frac{a_i a_i^\top}{4} + 2\lambda \cdot \mathbf{I}_d.$$

Due to Proposition 2, it immediately follows that f_i and f satisfy Definition 1 with $\mathbf{L}_i \in \mathbb{S}_{++}^d$ and $\mathbf{L} \in \mathbb{S}_{++}^d$, respectively.

In the following subsections, we perform several numerical experiments comparing the performance of **DCGD**, **det-CGD**, **MARINA**, **DASHA**, **det-MARINA** and **det-DASHA**. The datasets we used are from the LibSVM repository (Chang & Lin, 2011).

I.2 COMPARISON OF ALL THE METHODS

In this section, we present several plots which compare all relevant methods to the **det-MARINA** and **det-DASHA**. The methods are the following: (i) **DCGD** with scalar stepsize γ_2 , (ii) **det-CGD** with matrix stepsize \mathbf{D}_3^* , (iii) **MARINA** with scalar stepsize γ_1 , (iv) **DASHA** with scalar stepsize γ_4 , (v) **det-MARINA** with \mathbf{D}_{L-1}^* , (vi) **det-DASHA** with \mathbf{D}_{L-1}^{**} . Throughout the experiment, $\varepsilon = 0.01$, and $\lambda = 0.9$, we are using the same Rand- τ sketch for all the algorithms, and we run all the algorithms for a fixed number of iteration $K = 10000$.

It can be seen in Figure 2, the performance in terms of communication complexity of **det-DASHA** and **det-MARINA** is better than their scalar counterpart **DASHA** and **MARINA** respectively. This validates the efficiency of using a matrix stepsize over a scalar stepsize. Furthermore, we notice that **det-DASHA** and **det-MARINA** have better communication complexity in this case, compared to **det-CGD**. In addition, we observe variance reduction.

Notice that the optimal stepsizes of **det-CGD** and **DCGD** require information of function value differences at x^* . Furthermore, the stepsizes are also constrained by the number of iterations K and the error ε^2 . Meanwhile, for the variance reduced methods, we do not require such considerations, which is much more practical in general.

I.3 IMPROVEMENTS OVER **MARINA**

The purpose of this experiment is to compare the iteration complexity of **MARINA**, with **det-MARINA** using Rand- τ sketches, thus showing improvements of **det-MARINA** upon **MARINA**.

2538
 2539
 2540
 2541
 2542
 2543
 2544
 2545
 2546
 2547
 2548
 2549
 2550
 2551
 2552
 2553
 2554
 2555
 2556
 2557
 2558
 2559
 2560
 2561
 2562
 2563
 2564
 2565
 2566
 2567
 2568
 2569
 2570
 2571
 2572
 2573
 2574
 2575
 2576
 2577
 2578
 2579
 2580
 2581
 2582
 2583
 2584
 2585
 2586
 2587
 2588
 2589
 2590
 2591

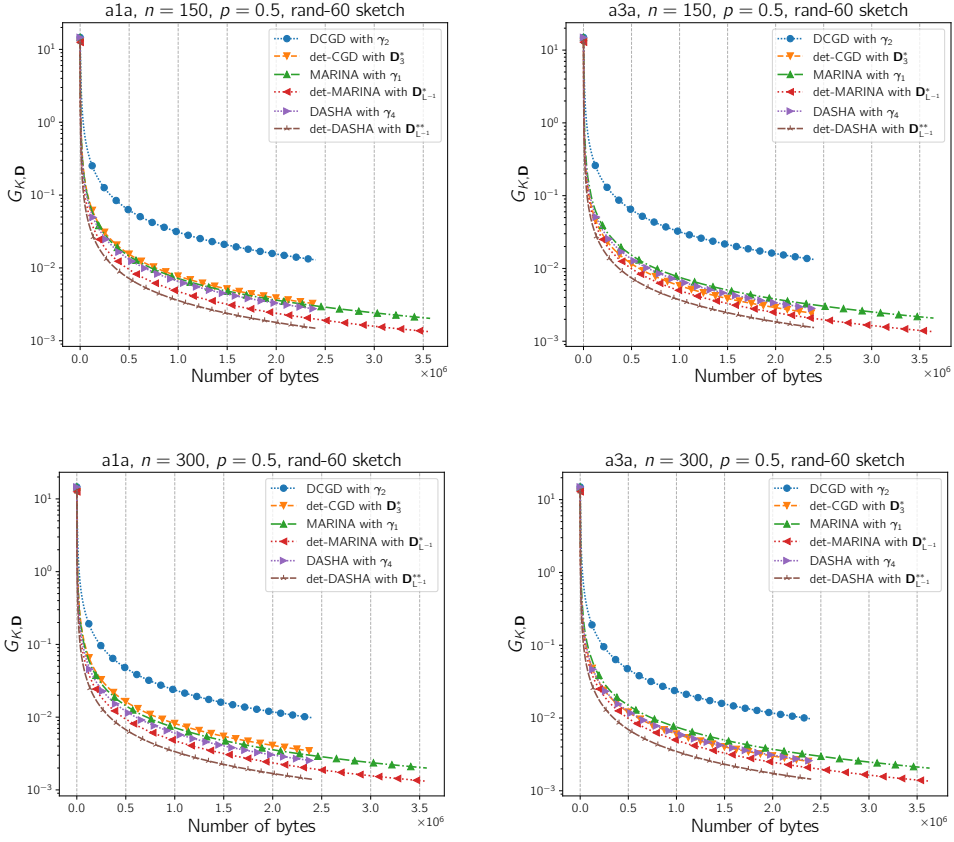


Figure 2: Comparison of **DCGD** with optimal scalar stepsize, **det-CGD** with matrix stepsize D_3^* , **MARINA** with optimal scalar stepsize, **DASHA** with optimal scalar stepsize, **det-MARINA** with optimal stepsize D_{L-1}^* and **det-DASHA** with optimal stepsize D_{L-1}^{**} . Throughout the experiment, we are using Rand- τ sketch with $\tau = 60$, and each algorithm is run for a fixed number of iterations $K = 10000$. The momentum of **DASHA** is set as $1/2\omega+1$ and **det-DASHA** is $1/2\omega_D+1$. The notation n in the title stands for the number of clients in each case, and p stands for the probability used by **MARINA** and **det-MARINA**.

Using Theorem C.1 from (Gorbunov et al., 2021), we deduce the optimal stepsize for **MARINA**, is

$$\gamma_1 = \frac{1}{L \left(1 + \sqrt{\frac{(1-p)\omega}{pn}} \right)}, \quad (61)$$

where ω is the quantization coefficient. In particular, for the Rand- τ compressor $\omega = \frac{d}{\tau} - 1$. For the full definition see Section 1.3 of (Gorbunov et al., 2021). The stepsize for **det-MARINA** is determined through Corollary 1. We use the notation D_W^* to denote the optimal stepsize for each choice of W , here we list some of the optimal stepsizes for different W , which are used in the experiment section. We have

$$\begin{aligned} D_{I_d}^* &= \frac{2}{1 + \sqrt{1 + 4\alpha\beta \frac{1}{\lambda_{\max}(L)} \cdot \omega}} \cdot \frac{I_d}{\lambda_{\max}(L)}, \\ D_{L^{-1}}^* &= \frac{2}{1 + \sqrt{1 + 4\alpha\beta \cdot \lambda_{\max}(\mathbb{E}[S_i^k L^{-1} S_i^k] - L^{-1})}} \cdot L^{-1}, \\ D_{\text{diag}^{-1}(L)}^* &= \frac{2}{1 + \sqrt{1 + 4\alpha\beta \cdot \lambda_{\max}(\mathbb{E}[S_i^k \text{diag}^{-1}(L) S_i^k] - \text{diag}^{-1}(L))}} \cdot \text{diag}^{-1}(L) \end{aligned} \quad (62)$$

In this experiment, we aim to compare **det-MARINA** with stepsize $\mathbf{D}_{L^{-1}}^*$ to the standard **MARINA** with the optimal scalar stepsize. Rand- τ compressor is used in the comparison. Throughout the experiments, λ is fixed at 0.3. We set the x -axis to be the number of iterations, while y -axis to be the expectation of the corresponding matrix norm of the gradient of the function, which is defined as

$$G_{K,D} = \mathbb{E} \left[\left\| \nabla f(\bar{x}^K) \right\|_{\mathbf{D}/\det(\mathbf{D})^{1/d}}^2 \right]. \quad (63)$$

Notice that this criterion is comparable to the standard Euclidean norm Li et al. (2024b), and for a fixed \mathbf{D} , we have

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

As it is illustrated in Figure 3, **det-MARINA** always has a faster convergence rate compared to **MARINA** if they use the same sketch, this justifies the result we have in Corollary 3. Notice that in some cases, **det-MARINA** with Rand-1 sketch even outperforms standard **MARINA** with Rand-80 sketch. This further demonstrates the superiority of matrix stepsizes and smoothness over the standard scalar setting.

I.4 IMPROVEMENTS ON NON VARIANCE REDUCED METHODS

In this section, we compare two non-variance reduced methods, distributed compressed gradient descent (**DCGD**) and distributed **det-CGD**, with two variance reduced methods, **MARINA**, and **det-MARINA**. Rand-1 sketch is used throughout this experiment for all the algorithms, for non variance reduced method ε^2 is fixed at 0.01 in order to determine the optimal stepsize. The purpose of this experiment is to show the advantages of variance reduced methods over non variance reduced methods. **DCGD** was initially proposed in (Khirirat et al., 2018). Later on **DIANA** was proposed in (Mishchenko et al., 2019) and then combined with variance reduction technique. Recently Shulgin & Richtárik (2022) proposed shifted **DCGD**, which is a shifted version of **DCGD** and proved its convergence in the (strongly) convex setting. A general analysis on **SGD** type methods in the non-convex world is provided by Khaled & Richtárik (2023), including **DCGD** and shifted **DCGD**. In our case, in order to determine the optimal scalar stepsize for **DCGD**, one can simply use Proposition 4 in (Khaled & Richtárik, 2023). One can check that in order to satisfy $\min_{0 \leq k \leq K-1} \mathbb{E} \left[\left\| \nabla f(x^k) \right\|^2 \right] \leq \varepsilon^2$ the stepsize condition for **DCGD** in the non-convex case reduces to

$$\gamma_2 \leq \min \left\{ \frac{1}{L}, \sqrt{\frac{n}{\omega L L_{\max} K}}, \frac{n\varepsilon^2}{4L L_{\max} \omega \cdot \Delta^*} \right\},$$

where L is the smoothness constant for f , L_i is the smoothness constant for f_i , $L_{\max} = \max_i L_i$, K is the total number of iterations, $\Delta^* = f(x^*) - \frac{1}{n} \sum_{i=1}^n f_i(x^*)$. The constant ω is associated with the compressor used in the algorithm, for Rand- τ sketch, it is $\frac{d}{\tau} - 1$. For distributed **det-CGD** according to Li et al. (2024b), the stepsize condition in order to satisfy $\min_{0 \leq k \leq K-1} \mathbb{E} \left[\left\| \nabla f(x) \right\|_{\mathbf{D}/\det(\mathbf{D})^{1/d}}^2 \right] \leq \varepsilon^2$ is

$$\mathbf{D} \mathbf{L} \mathbf{D} \preceq \mathbf{D}, \quad \lambda_{\mathbf{D}} \leq \min \left\{ \frac{n}{K}, \frac{n\varepsilon^2}{4\Delta^*} \det(\mathbf{D})^{1/d} \right\}, \quad (64)$$

where $\lambda_{\mathbf{D}}$ is defined as

$$\lambda_{\mathbf{D}} = \max_i \left\{ \lambda_{\max} \left(\mathbb{E} \left[\mathbf{L}_i^{\frac{1}{2}} (\mathbf{S}_i^k - \mathbf{I}_d) \mathbf{D} \mathbf{L} \mathbf{D} (\mathbf{S}_i^k - \mathbf{I}_d) \mathbf{L}_i^{\frac{1}{2}} \right] \right) \right\}. \quad (65)$$

In general cases, there is no easy way to find a optimal stepsize matrix \mathbf{D} satisfying (64), alternatively, we choose the optimal diagonal stepsize \mathbf{D}_3^* similarly to (Li et al., 2024b). The stepsize condition for **MARINA** has already been described by (61). Note that we only consider **MARINA**, but not **DIANA** or shifted **DCGD**, because **DIANA** and shifted **DCGD** offer suboptimal rates compared to **MARINA** in the non-convex setting. For **det-MARINA**, we fix $\mathbf{W} = \mathbf{L}^{-1}$, and use $\mathbf{D}_{L^{-1}}^*$ as the stepsize matrix. In theory, **det-MARINA** in this case should always out perform **MARINA** in terms of iteration complexity.

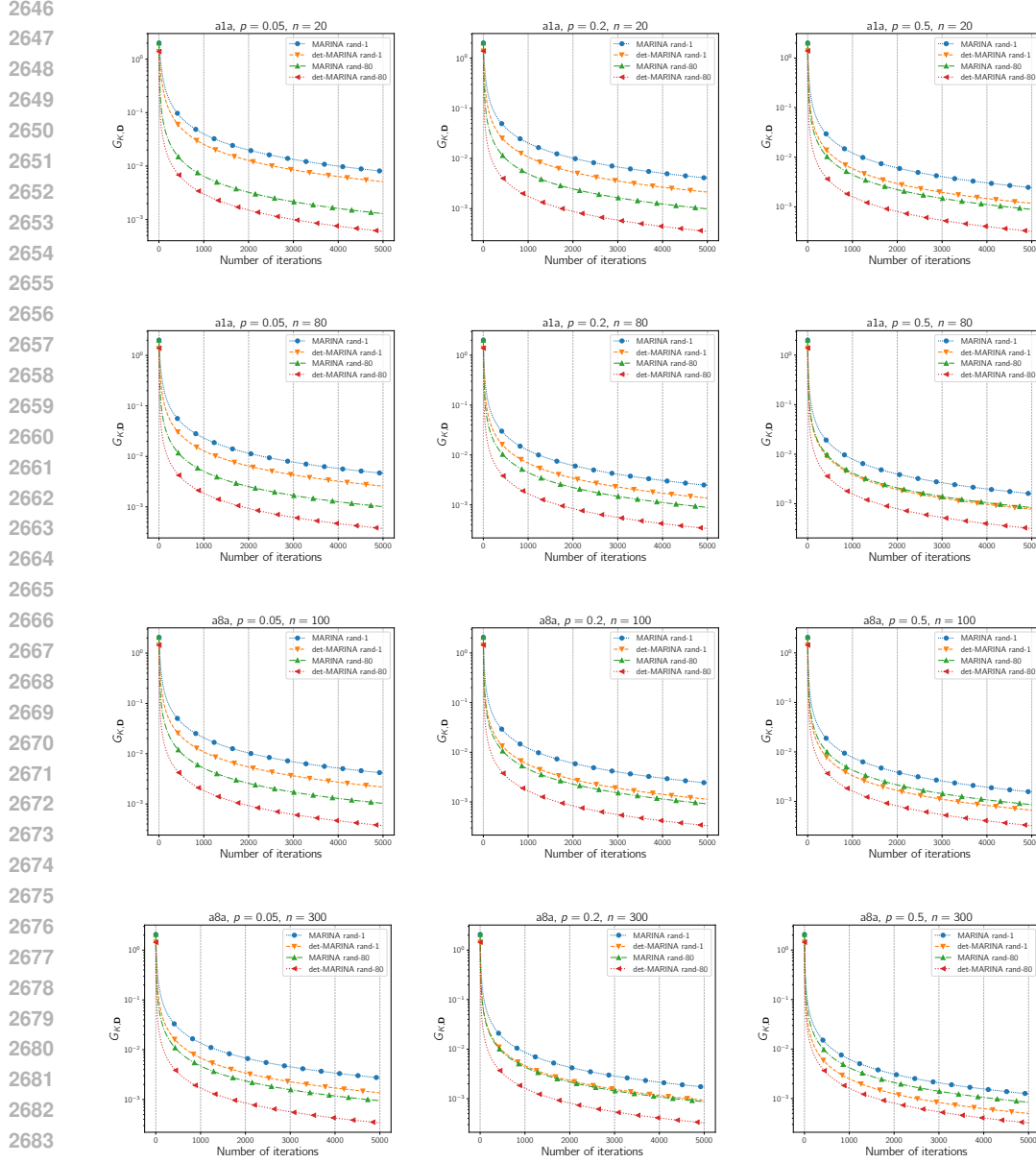


Figure 3: In this experiment, we aim to compare det-MARINA with stepsize D_{L-1}^* to the standard MARINA with the optimal scalar stepsize. Rand- τ compressor is used in the comparison. Throughout the experiments, λ is fixed at 0.3. Optimal stepsize is calculated in each case with respect to the sketch used. The x -axis denotes the number of iterations while the notation $G_{K,D}$ for the y -axis is defined in (63), which is the averaged matrix norm of the gradient. The notation p in the title denotes the probability used in the two algorithms, n denotes the number of clients in each setting.

In Figure 4, in each plot, we observe that det-MARINA outperforms MARINA and the rest of the non-variance reduced methods. This is expected, since our theory confirms that det-MARINA indeed has a better rate compared to MARINA , and the stepsizes of the non-variance reduced methods are negatively affected by the neighborhood. When p is reasonably large, the variance reduced methods considered here outperform the non-variance reduced methods. In this experiment we consider only the comparison involving det-CGD .

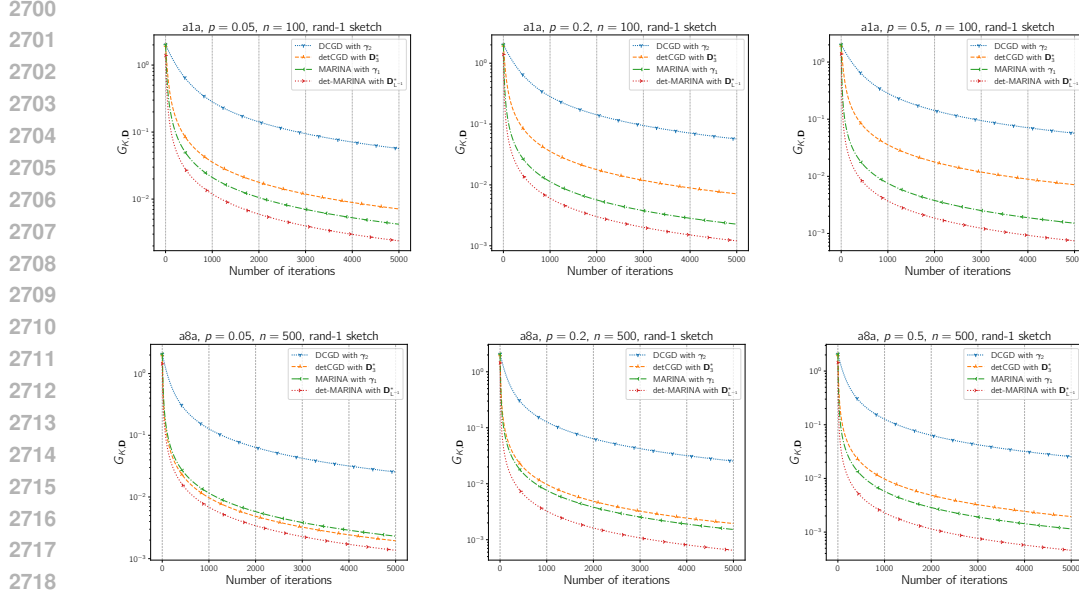


Figure 4: Comparison of **DCGD** with optimal scalar stepsize γ_2 , **det-CGD** with optimal diagonal stepsize D_3^* , **MARINA** with optimal scalar stepsize γ_1 , and **det-MARINA** with optimal stepsize D_{L-1}^* with respect to $\mathbf{W} = \mathbf{L}^{-1}$. In each case, probability p is chosen the set $\{0.05, 0.2, 0.5\}$ for **MARINA** and **det-MARINA**. $\lambda = 0.3$ is fixed throughout the experiment. The notation n in the title indicates the number of clients in each case.

I.5 IMPROVEMENTS OVER **det-CGD**

In this section, we compare **det-CGD** in the distributed case with **det-MARINA**, which are both algorithms using matrix stepsizes and matrix smoothness. The purpose of this experiment is to show that **det-MARINA** improves on the current state of the art matrix stepsize compressed gradient method when the objective function is non-convex. Throughout the experiment, $\lambda = 0.3$ is fixed, and for **det-CGD**, $\varepsilon^2 = 0.01$ is fixed in order to determine its stepsize. For a thorough comparison, we select the stepsize for **det-CGD** in the following way. Let us denote the stepsize as $\mathbf{D} = \gamma_{\mathbf{W}} \cdot \mathbf{W}$, where $\gamma_{\mathbf{W}} \in \mathbb{R}_{++}$, $\mathbf{W} \in \mathbb{S}_{++}^d$. We first fix a matrix \mathbf{W} , in this case, we pick \mathbf{W} from the set $\{\mathbf{L}^{-1}, \text{diag}^{-1}(\mathbf{L}), \mathbf{I}_d\}$, and then we determine the optimal scaling $\gamma_{\mathbf{W}}$ for each case using the condition given in (Li et al., 2024b) (see (64) and (65)). Then, we denote the matrix stepsizes for **det-CGD**

$$\mathbf{D}_1 = \gamma_{\mathbf{I}_d} \cdot \mathbf{I}_d, \quad \mathbf{D}_2 = \gamma_{\text{diag}^{-1}(\mathbf{L})} \cdot \text{diag}^{-1}(\mathbf{L}), \quad \mathbf{D}_3 = \gamma_{\mathbf{L}^{-1}} \cdot \mathbf{L}^{-1}. \quad (66)$$

For **det-MARINA**, we use the stepsize D_{L-1}^* , which is described in (62). In this experiment, we compare **det-CGD** using three stepsizes $\mathbf{D}_1, \mathbf{D}_2, \mathbf{D}_3$ with **det-MARINA** using stepsize D_{L-1}^* .

From Figure 5, it is clear that **det-MARINA** outperforms **det-CGD** with all matrix optimal stepsizes with respect to a fixed \mathbf{W} considered here. This is expected, since the convergence rate of non-variance reduced methods are affected by its neighborhood. This experiment demonstrates the advantages of **det-MARINA** over **det-CGD**, and is also supported by our theory. Notice that though different \mathbf{W} are considered for **det-CGD**, their convergence rates are similar, which is also mentioned by Li et al. (2024b).

I.6 COMPARING DIFFERENT STEPSIZE CHOICES

This experiment is designed to see the how **det-MARINA** works under different stepsize choices. As it is mentioned in Appendix I.3, for each choice of $\mathbf{W} \in \mathbb{S}_{++}^d$, an optimal stepsize $D_{\mathbf{W}}^*$ can be determined. Here we compare **det-MARINA** using three different stepsize choices $D_{L-1}^*, D_{\text{diag}^{-1}(\mathbf{L})}^*$ and $D_{\mathbf{I}_d}^*$. These stepsizes are explicitly defined in (62). Throughout the experiment, we fix $\lambda = 0.3$, Rand-1 sketch is used in all cases.

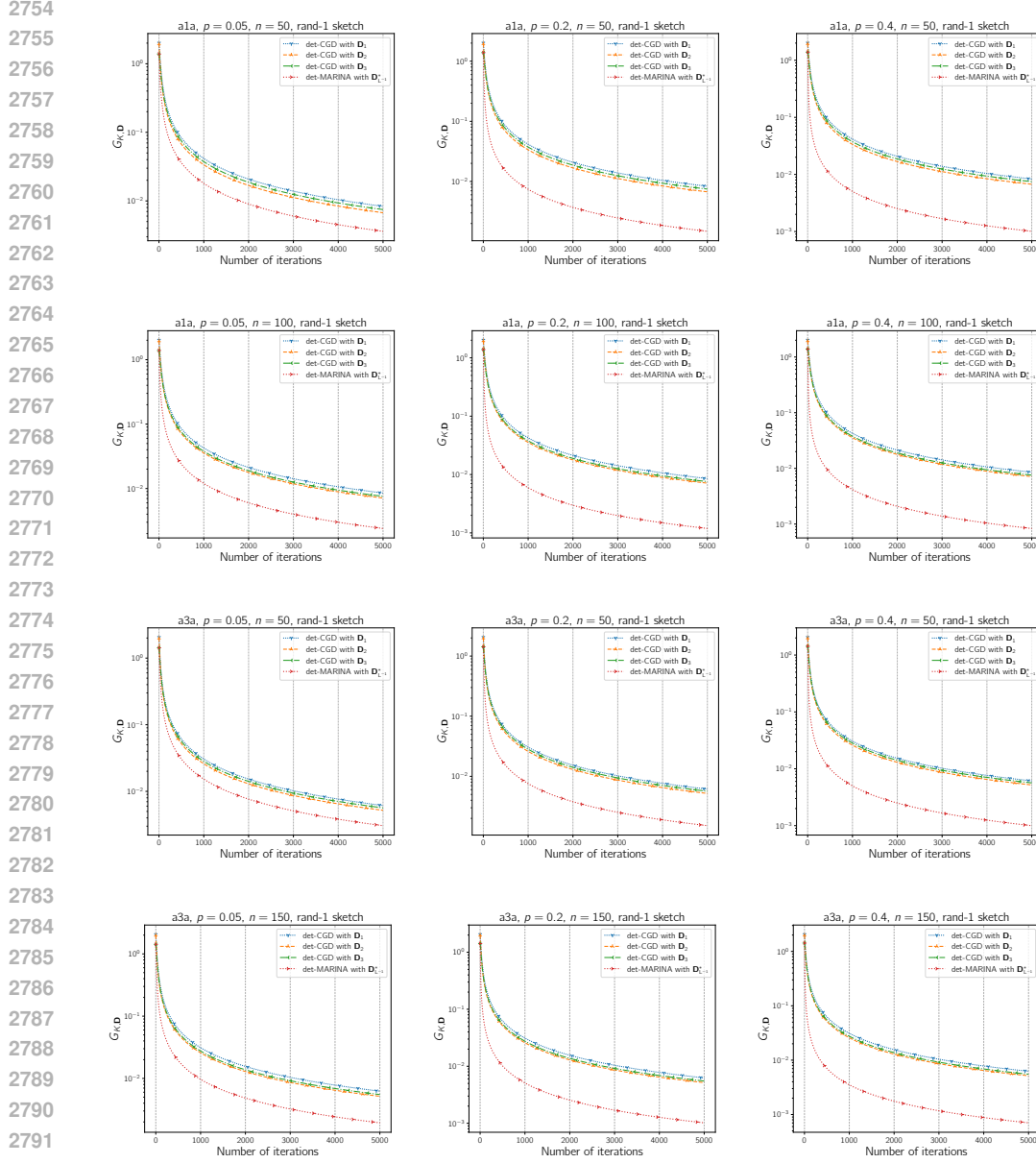


Figure 5: Comparison of **det-CGD** with matrix stepsize D_1 , D_2 and D_3 and **det-MARINA** with optimal matrix stepsize with respect to $W = L^{-1}$. The stepsizes $\{D_i\}_{i=1}^3$ are described in (66). Throughout the experiment ε^2 is fixed at 0.01, the notation p in the title refers to the probability for **det-MARINA**, n denotes the number of clients considered, Rand-1 sketch is used in all cases for all the algorithms.

We can observe from Figure 6 that, in almost all cases **det-MARINA** with stepsize $D_{\text{diag}^{-1}(L)}^*$ and D_{L-1}^* outperforms **det-MARINA** with $D_{I_d}^*$. As **det-MARINA** with $D_{I_d}^*$ can be viewed as **MARINA** using scalar stepsize but under matrix Lipschitz gradient assumption, this demonstrates the effectiveness of using a matrix stepsize over the scalar stepsize. However, in Figure 6, there are cases where **det-MARINA** with $D_{\text{diag}^{-1}(L)}^*$ outperforms D_{L-1}^* . This tells us the two stepsizes are perhaps incomparable in general cases. This is similar to **det-CGD**, where optimal stepsizes with respect to a subspace associated with a fixed W^{-1} are incomparable.

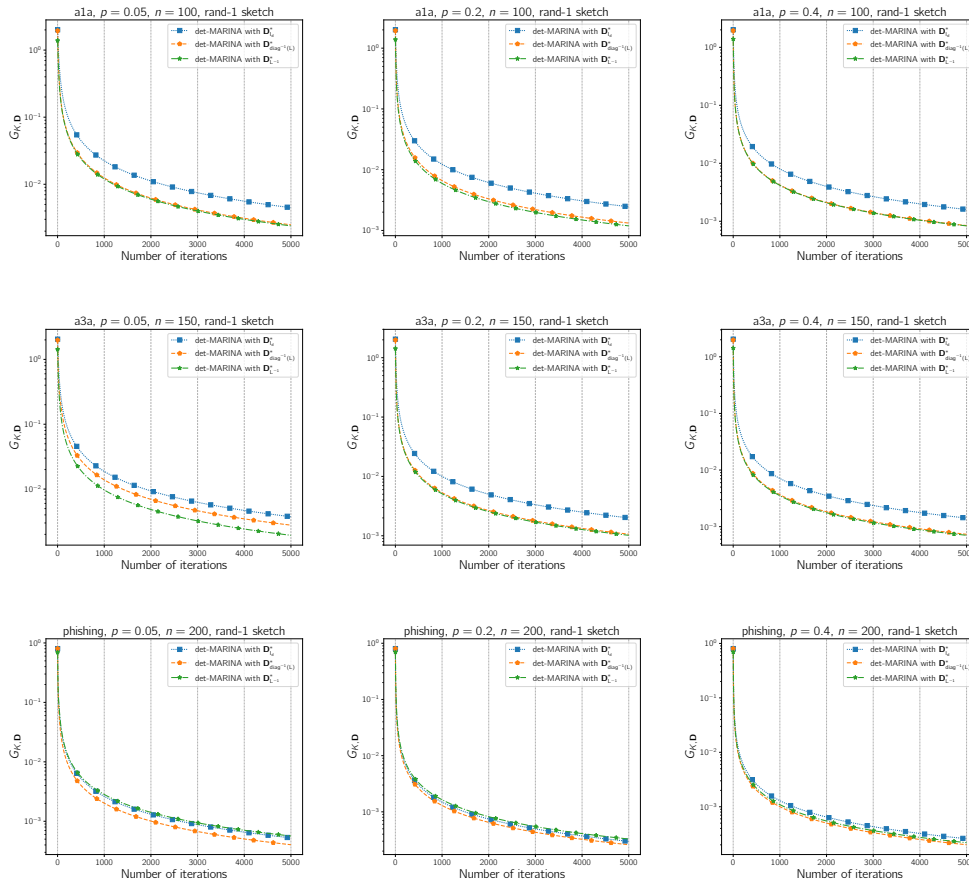


Figure 6: Comparison of **det-MARINA** with matrix stepsize $D_{I_d}^*$, $D_{\text{diag}^{-1}(L)}^*$ and D_{L-1}^* . The stepsizes are defined in (62). Throughout the experiment, $\lambda = 0.3$ is fixed, Rand-1 sketch is used in all cases. The notation p in the title indicates the probability of sending the true gradient for **det-MARINA**, n denotes the number of clients considered.

I.7 COMPARING COMMUNICATION COMPLEXITY

In this section, we perform an experiment on how different probabilities p will affect the overall communication complexity of **det-MARINA**. We use D_{L-1}^* as the stepsize, which is determined with respect to the sketch used. Rand- τ sketches are used in these experiments, and we vary the minibatch size τ to provide a more comprehensive comparison. For Rand- τ sketch S and any $A \in \mathbb{S}_{++}^d$, one can show that

$$\mathbb{E}[SAS] = \frac{d}{\tau} \left(\frac{d-\tau}{d-1} \text{diag}(A) + \frac{\tau-1}{d-1} A \right). \quad (67)$$

Combining (67) and (62), we can find out the corresponding matrix stepsize easily. In the experiment, a fixed number of iterations ($K = 5000$) is performed for each **det-MARINA** with the corresponding stepsize.

As it can be observed from Figure 7, in each dataset, the communication complexity tends to increase with the increase of probability p . However, when the number of iteration is fixed, a larger p often means a faster rate of convergence. This difference in communication complexity is more obvious when we are using the Rand-1 sketch. In real federated learning settings, there is often constraints on network bandwidth from clients to the server. Thus, trading off between communication complexity and iteration complexity, i.e. selecting the compression mechanism carefully to guarantee a acceptable speed that satisfies the bandwidth constraints, becomes important.

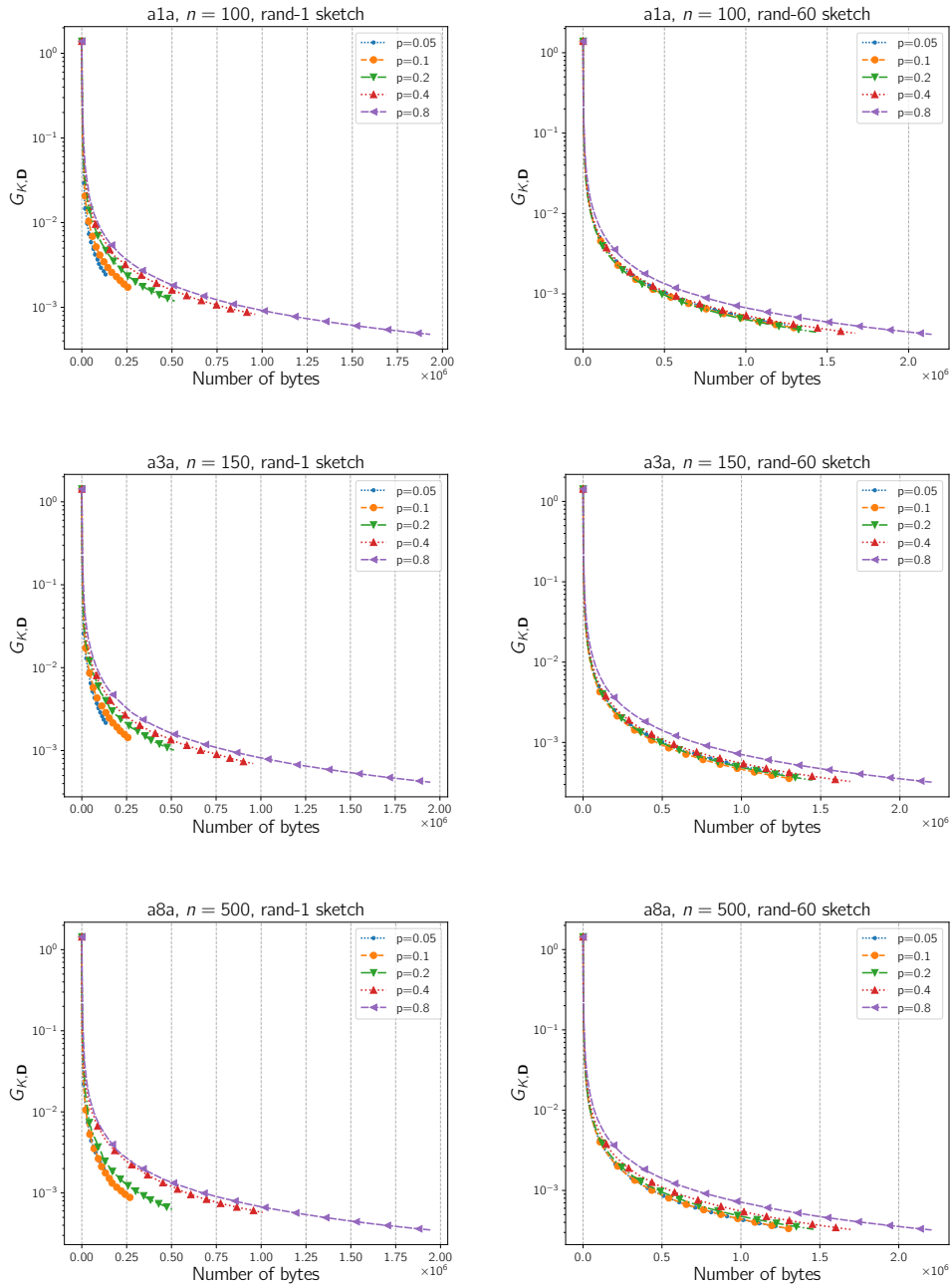


Figure 7: Comparison of **det-MARINA** with stepsize D_{L-1}^* using different probability p . The probability p here is chosen from the set $\{0.05, 0.1, 0.2, 0.4, 0.8\}$. The notation n in the title denote the number of clients considered. The x -axis is now the number of bytes sent from a single node to the server. In each case, **det-MARINA** is run for a fixed number of iterations $K = 5000$.

I.8 COMPARISON OF **DASHA** AND **DET-DASHA**

In this experiment we plan to compare the performance of original **DASHA** with **det-DASHA**. Throughout the experiments, λ is fixed at 0.3. The same Rand- τ sketch is used in the two algo-

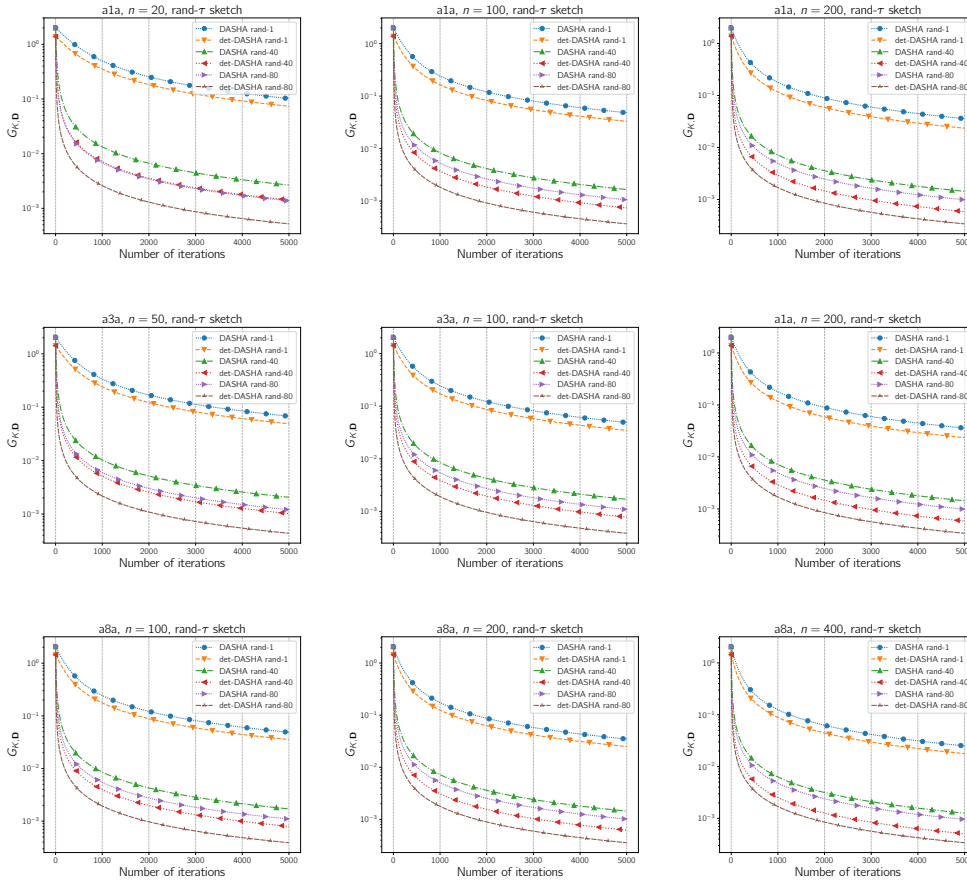


Figure 8: Comparison of **det-DASHA** with matrix stepsize D_{L-1}^{**} and **DASHA** with optimal scalar stepsize γ using different **Rand- τ** sketches. $\lambda = 0.3$ is fixed throughout the experiments. Optimal stepsize is calculated in each case with respect to the sketch used. The x -axis denotes the number of iterations while the notation $G_{K,D}$ for the y -axis denotes the averaged matrix norm of the gradient. The notation n denotes the number of clients in each setting.

rithms. The stepsize condition on **DASHA** when the momentum is set as $a = \frac{1}{2\omega+1}$ is given as

$$\gamma_4 \leq \left(L + \sqrt{\frac{16\omega(2\omega+1)\widehat{L}}{n}} \right)^{-1},$$

according to Theorem 6.1 of Tyurin & Richtárik (2024). Here the L is the smoothness constant of the function f , while \widehat{L} satisfies $\widehat{L}^2 = \frac{1}{n} \sum_{i=1}^n L_i^2$ where L_i is the smoothness constant of local objective f_i . In theory we can pick $\widehat{L} = L$. Similarly, according to Corollary 2, the optimal stepsize matrix D_{L-1}^{**} is given as

$$D_{L-1}^{**} = \frac{2}{1 + \sqrt{1 + 16C_{L-1} \cdot \lambda_{\min}(L)}} \cdot L^{-1}, \quad (68)$$

when the momentum is given as $a = \frac{1}{2\omega_D+1}$. We compare the performance of **DASHA** with ω and **det-DASHA** with D_{L-1}^{**} using the same sketch where the total number of clients are different.

As it can be observed in Figure 8, **det-DASHA** with matrix stepsize D_{L-1}^{**} outperforms **DASHA** with optimal scalar stepsize using the same sketch in every setting we considered. Note that since the same sketch is used in the two algorithm, the number of bits transferred in each iteration is also the same for the two algorithms. This essentially indicates that **det-DASHA** has better iteration

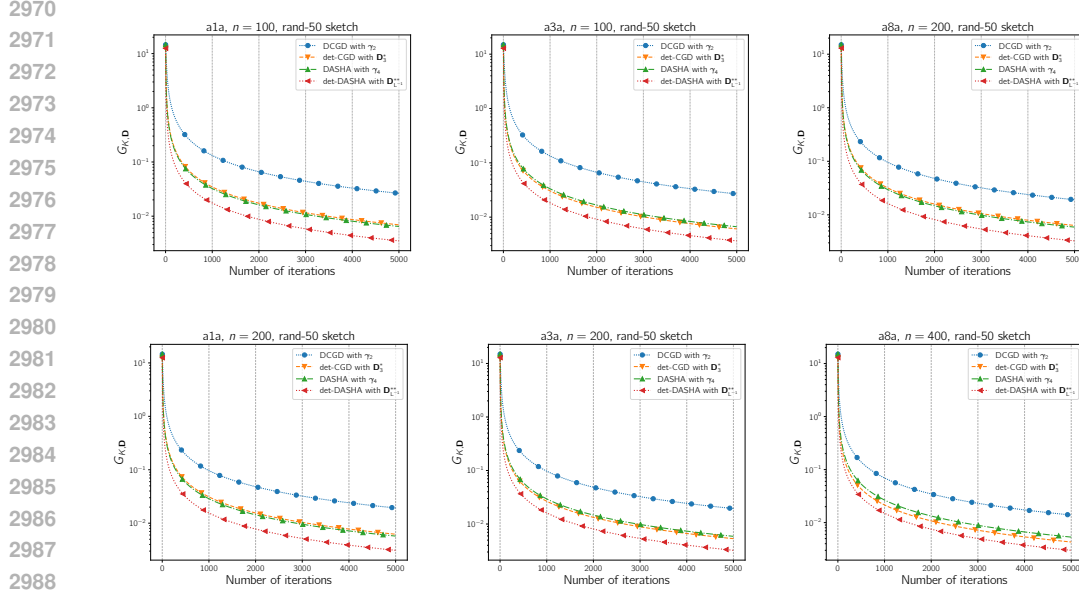


Figure 9: Comparison of **DCGD** with optimal scalar stepsize γ_2 , **det-CGD** with optimal diagonal stepsize D_3^* , **DASHA** with optimal scalar stepsize γ_1 , and **det-DASHA** with optimal stepsize D_{L-1}^* with respect to $W = L^{-1}$. $\lambda = 0.9$ is fixed throughout the experiment. The notation n in the title indicates the number of clients in each case. Rand- τ sketch with $\tau = 50$ are used in all four algorithms.

complexity as well as communication complexity than **DASHA** given that the same sketch is used for the two algorithm.

I.9 COMPARISON OF **DCGD**, **DET-CGD**, **DASHA** AND **DET-DASHA**

In this experiment, we consider the comparison between the two non variance reduced methods **DCGD**, **det-CGD** and the two variance reduced method **DASHA**, **det-DASHA**. The stepsize choices for **DCGD** and **det-CGD** have already been discussed in the previous sections, (for **DCGD** we use γ_2 and for **det-CGD** we use D_3^*), for **DASHA** and **det-DASHA**, we use the stepsize choices of Appendix I.8. Note that ε^2 is set as 0.01, and λ is fixed at 0.9 here. Throughout this experiment, we consider the case where Rand- τ sketch is used in the four algorithms.

It is easy to observe that in each case of Figure 9, **det-DASHA** outperforms the rest of the algorithms. It is expected that **det-DASHA** outperforms **DASHA**, as it is also illustrated by Figure 8, which is a consequence of using matrix stepsize instead of a scalar stepsize. We also see that **det-DASHA** and **DASHA** outperform **det-CGD** and **DCGD** respectively, which demonstrate the advantages of the variance reduction technique. Note that in this case, all four algorithms are using the same sketch, which means that the number of bits transferred in each iteration is the same for the four algorithms, as a result, compared to the other algorithms, **det-DASHA** is better in terms of both iteration complexity and communication complexity.

I.10 COMPARISON OF **DET-DASHA** AND **DET-CGD** WITH DIFFERENT STEPSIZES

In this experiment, we try to compare **det-DASHA** and **det-CGD** with different matrix stepsizes. Throughout this experiment, we will fix $\varepsilon^2 = 0.01$ and $\lambda = 0.9$. The same Rand- τ sketch is used for the two algorithms. For **det-CGD**, we use the stepsize $D_1 = \gamma_{I_d} \cdot I_d$, $D_2 = \gamma_{\text{diag}^{-1}(L)} \cdot \text{diag}^{-1}(L)$ and $D_3 = \gamma_{L^{-1}} \cdot L^{-1}$, for **det-DASHA** we use the stepsize D_{L-1}^* .

It can be observed that in all cases of Figure 10, **det-DASHA** outperforms **det-CGD** with different stepsizes. This further corroborates our theory that **det-DASHA** is variance reduced and thus is better in terms of both iteration complexity, and communication complexity (because in this case

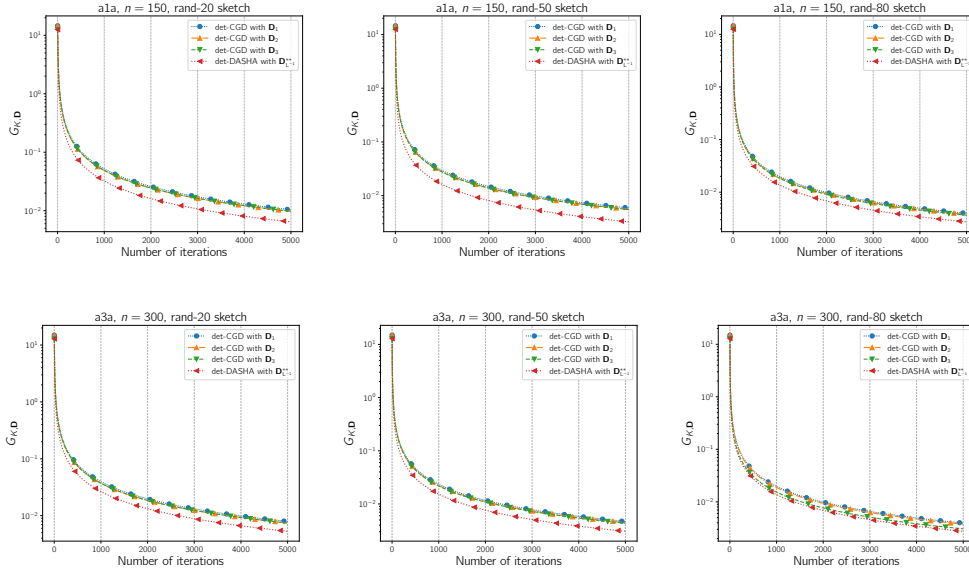


Figure 10: Comparison of **det-DASHA** with stepsize D_{L-1}^{**} and **det-CGD** with three different stepsizes D_1 , D_2 and D_3 . Throughout the experiment, λ is fixed at 0.9, ε^2 is fixed at 0.01, Rand- τ sketch is used for all the algorithms with τ selected from $\{20, 50, 80\}$. The notation n denotes the number of clients in each setting.

the same number of bits are transmitted in each iteration due to the fact that the sketch used is the same).

I.11 COMPARISON OF DIFFERENT STEPSIZES OF **DET-DASHA**

In this experiment, we try to compare **det-DASHA** with different matrix stepsizes. Specifically, we fix matrix W to be three different matrices, I_d , $\text{diag}^{-1}(L)$ and L^{-1} . We denote the optimal stepsizes as $D_{I_d}^{**}$, $D_{\text{diag}^{-1}(L)}^{**}$ and $D_{L^{-1}}^{**}$, respectively. For $D_{L^{-1}}^{**}$, it is already given in (68), for $D_{I_d}^{**}$ and $D_{\text{diag}^{-1}(L)}^{**}$, we use Corollary 2 to compute them. As a result,

$$D_{I_d}^{**} = \frac{2}{1 + \sqrt{1 + 16 \cdot \frac{\omega_{I_d}(4\omega_{I_d} + 1)}{n} \cdot \frac{\lambda_{\min}(L)}{\lambda_{\max}(L)}}} \cdot \frac{I_d}{\lambda_{\max}(L)}, \quad (69)$$

and

$$D_{\text{diag}^{-1}(L)}^{**} = \frac{2}{1 + \sqrt{1 + 16C_{\text{diag}^{-1}(L)} \cdot \lambda_{\min}(L)}} \cdot \text{diag}^{-1}(L). \quad (70)$$

Throughout the experiment, λ is fixed at 0.9, Rand- τ sketch is used for all the algorithms.

We can observe from Figure 11, **det-DASHA** with $D_{L^{-1}}^{**}$ and $D_{\text{diag}^{-1}(L)}^{**}$ both outperform **det-DASHA** with $D_{I_d}^{**}$, which demonstrate the effectiveness of using a matrix stepsize instead of a scalar stepsize. However, depending on the parameters of the problem, it is hard to reach a general conclusion whether $D_{L^{-1}}^{**}$ is better than $D_{\text{diag}^{-1}(L)}^{**}$ or not.

I.12 COMPARISON OF **DET-MARINA** AND **DET-DASHA**

In this section, we aim to provide a comparison of **det-DASHA** and **det-MARINA**. They are similar as they are both variance reduced version of **det-CGD**. However, the variance reduction techniques that are utilized are different. For **det-MARINA**, it is based on **MARINA**, and it requires synchronization from time to time depending on a probability parameter p , while for **det-DASHA** it utilizes

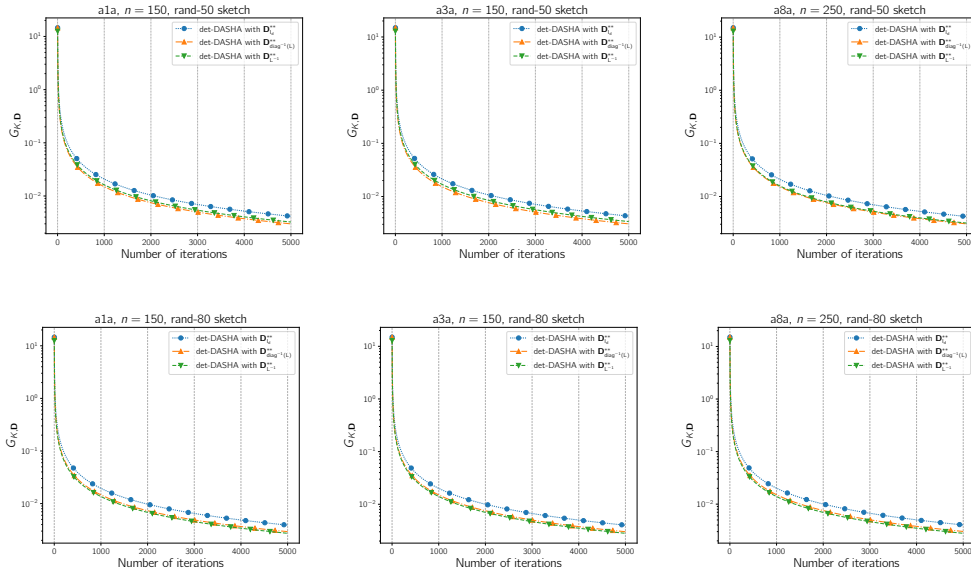


Figure 11: Comparison of **det-DASHA** three different stepsizes D_{L-1}^{**} , $D_{\text{diag}^{-1}(L)}^{**}$ and $D_{I_d}^{**}$. The definition for those matrix stepsize notation are given in (68), (70) and (69) respectively. Throughout the experiment, λ is fixed at 0.9, Rand- τ sketch is used for all the algorithms. The notation n denotes the number of clients in each setting.

the momentum variance reduction technique which was also presented in **DASHA**, it does not need any synchronization at all. Notice that for a fair comparison, we implement the two algorithms so that they use the same sketch. We mainly focus on the communication complexity, i.e. the convergence with respect to the number of bits transferred. Throughout the experiment, $\lambda = 0.9$ is fixed. For **det-DASHA** we pick 3 different kinds of stepsizes $D_{I_d}^{**}$, D_{L-1}^{**} and $D_{\text{diag}^{-1}(L)}^{**}$. For **det-MARINA**, we also pick three different kinds of stepsizes correspondingly $D_{I_d}^*$, D_{L-1}^* and $D_{\text{diag}^{-1}(L)}^*$. We use the same sketch for all of the algorithms we are trying to compare.

It is obvious from Figure 12 that **det-DASHA** always has a better communication complexity comparing to the **det-MARINA** counterpart. Notice that here since each algorithm is run for a fixed number of iterations, so x -axis actually records the total number of bytes transferred for each algorithm. For **det-DASHA**, D_{L-1}^{**} perform similarly to $D_{\text{diag}^{-1}(L)}^{**}$, and both are better than $D_{I_d}^{**}$. This is expected since the same sketch is used, and the number of bytes transferred in each iteration is the same for each variant of **det-DASHA**. The same relation also holds for **det-MARINA**.

I.13 COMPARISON IN TERMS OF FUNCTION VALUES

In this section, we compare **det-MARINA** and **det-DASHA** in terms of function values. The starting points of the two algorithms are set to be the same, and we run the two algorithms for multiple times and we average the function values we obtained in each iteration. For the two algorithms, we use the same sketch, and since we are interested in the performance in terms of communication complexity, we use the number of bytes transferred in the training process as the x -axis. We run each of the algorithm for 20 times, and fix $\lambda = 0.9$. The starting point is fixed throughout the experiment. We pick D_{L-1}^{**} as the stepsize of **det-DASHA**, while D_{L-1}^* as the stepsize of **det-MARINA**.

Observing Figure 13, we can see that the function values continuously decrease as the algorithms progress through more iterations. However, the stability observed here differs from the case of the average (matrix) norm of gradients. Our theoretical framework, as presented in this paper, primarily addresses the average norm of gradients in the non-convex case. Despite this, the experiment reinforces the effectiveness of our algorithms, showcasing consistent decreases in function values.

3132
 3133
 3134
 3135
 3136
 3137
 3138
 3139
 3140
 3141
 3142
 3143
 3144
 3145
 3146
 3147
 3148
 3149
 3150
 3151
 3152
 3153
 3154
 3155
 3156
 3157
 3158
 3159
 3160
 3161
 3162
 3163
 3164
 3165
 3166
 3167
 3168
 3169
 3170
 3171
 3172
 3173
 3174
 3175
 3176
 3177
 3178
 3179
 3180
 3181
 3182
 3183
 3184
 3185

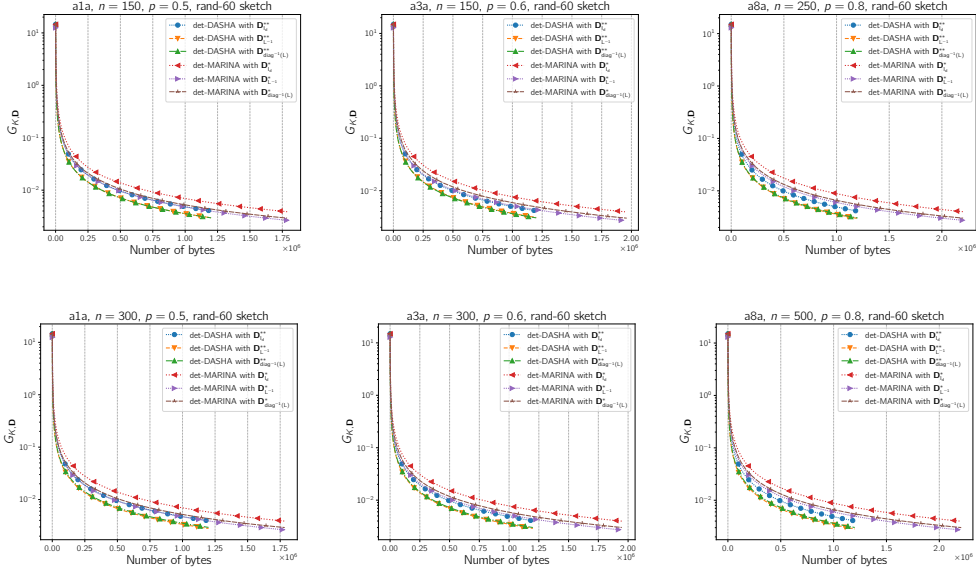


Figure 12: Comparison of **det-DASHA** with three different step sizes $D_{I_d}^{**}$, D_{L-1}^{**} and $D_{\text{diag}^{-1}(L)}^{**}$, and **det-MARINA** with $D_{I_d}^{*}$, D_{L-1}^{*} and $D_{\text{diag}^{-1}(L)}^{*}$ in terms of communication complexity. Throughout the experiment, λ is fixed at 0.9, the same Rand- τ sketch is used for all the algorithms. The notation n denotes the number of clients in each setting. Each algorithm is run for a fixed number of iteration $K = 5000$.

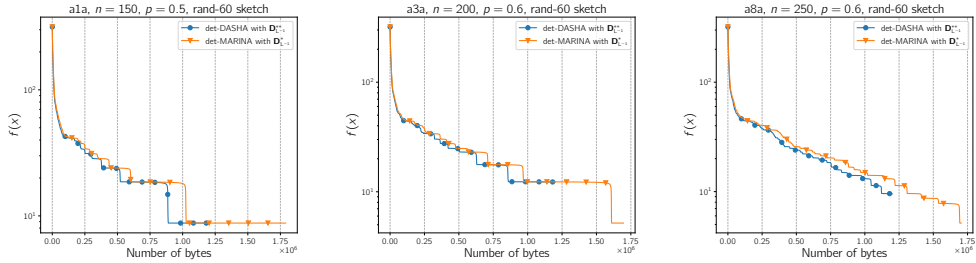


Figure 13: Comparing the performance of **det-DASHA** with D_{L-1}^{**} and **det-MARINA** with D_{L-1}^{*} in terms of the decreasing function values. The function values for each algorithm represent an average of 20 runs using different random seeds. Here, $\lambda = 0.9$ is fixed throughout the experiment, and the starting point for the two algorithms in different runs is the same. The notation n stands for the number of clients, and p represents the probability used in **det-MARINA**. The same Rand- τ sketch is employed for both algorithms.