Structured Preconditioners in Adaptive Optimization: A Unified Analysis

Shuo Xie¹ Tianhao Wang¹ Sashank Reddi² Sanjiv Kumar² Zhiyuan Li¹²

Abstract

We present a novel unified analysis for a broad class of adaptive optimization algorithms with structured (e.g., layerwise, diagonal, and kronecker-factored) preconditioners for both online regret minimization and offline convex optimization. Our analysis not only provides matching rate to several important structured preconditioned algorithms including diagonal AdaGrad, full-matrix AdaGrad, and AdaGrad-Norm, but also gives an improved convergence rate for a one-sided variant of Shampoo over that of original Shampoo. Interestingly, more structured preconditioners (e.g., diagonal Adagrad, AdaGrad-Norm which use less space and compute) are often presented as computationally efficient approximations to full-matrix Adagrad, aiming for improved optimization performance through better approximations. Our unified analysis challenges this prevailing view and reveals, perhaps surprisingly, that more structured preconditioners, despite using less space and computation per step, can outperform their less structured counterparts. To demonstrate this, we show that one-sided Shampoo, which is relatively much cheaper than fullmatrix AdaGrad could outperform it both theoretically and experimentally.

1. Introduction

Adaptive optimization algorithms (Streeter & McMahan, 2010; Duchi et al., 2011; Kingma & Ba, 2014) play a pivotal role in modern machine learning, especially in the expensive training of large foundation models. Within the machine learning community, full-matrix AdaGrad is considered as an ideal adaptive preconditioner for fast convergence in terms of number of steps. The computation of full-matrix AdaGrad preconditioner typically involves inverse square

root of a $d \times d$ matrix where d is the number of parameters. Thus, for large-scale settings, the huge computation and memory cost makes it prohibitively expensive. This has inspired works on designing more efficient adaptive optimizers by using a structured preconditioner, such as coordinate-wise adaptivity employed by AdaGrad (Streeter & McMahan, 2010; Duchi et al., 2011), Kronecker product based preconditioner employed by Shampoo (Gupta et al., 2018; Anil et al., 2020) and layerwise adaptivity employed by LARS (You et al., 2017) and LAMB (You et al., 2019)), which all aim to provide a computational and memory efficient approximation of full-matrix AdaGrad. These algorithms have been shown to be very effective for general deep learning settings. It is often assumed that a better approximation of full-matrix preconditioner usually results in better optimization convergence, as seen with methods like Shampoo (Gupta et al., 2018). In contrast, preconditioners with more structure such as diagonal preconditioner have inferior performance. In this paper, we challenge these prevailing notions by providing theoretical and empirical evidence against them.

Conceptually, one could equate the degree of structure in a preconditioner to the ease of its computation and storage. In this view, full-matrix AdaGrad can be considered as the least structured and most expensive preconditioner while AdaGrad-Norm (Ward et al., 2020), which only maintains a scalar and uses the same preconditioning for every direction, is the most structured and least expensive. The conventional wisdom here is that using a less structured preconditioner, which requires more space and compute per step, reduces the number of steps needed for training. Thus, choosing the right structure balances this trade-off between convergence speed and training step cost.

Among these preconditioning methods, Shampoo (Gupta et al., 2018) has gained notable attention for its Kroneckerfactored preconditioning approach, which promises improved convergence in large-scale optimization tasks (Dahl et al., 2023). Shampoo was originally proposed as a computationally efficient surrogate for full-matrix AdaGrad (Duchi et al., 2011). Despite its popularity, existing analyses of Shampoo (Gupta et al., 2018) are limited and do not provide a full justification for its effectiveness. In particular, we argue that the best-known regret bounds for both Shampoo as well as full-matrix AdaGrad are consistently worse

¹Toyota Technological Institute at Chicago ²Google Research. Correspondence to: Shuo Xie <shuox@ttic.edu>, Zhiyuan Li <zhiyuanli@ttic.edu>.

Proceedings of the 42^{nd} International Conference on Machine Learning, Vancouver, Canada. PMLR 267, 2025. Copyright 2025 by the author(s).

Structured Preconditioners in Adaptive Optimization: A Unified Analysis

Algorithm	Subalgebra ${\cal K}$	$\ m{x}\ _{\mathcal{H}}$	$\left\ \left\ oldsymbol{g}_{1:T} ight\ _{\mathcal{H}}$	Regret Bound = $\ \mathcal{X}\ _{\mathcal{H}} \cdot \ \ g_{1:T}\ \ _{\mathcal{H}}$
AdaGrad-Norm	$c \cdot \boldsymbol{I}_d$ for $c \in \mathbb{R}$	$\frac{\ \boldsymbol{x}\ _2}{\sqrt{d}}$	$\sqrt{d} \sqrt{\sum\limits_{t=1}^T \ oldsymbol{g}_t\ _2^2}$	$\ \mathcal{X}\ _2 \sqrt{\sum_{t=1}^T \ \boldsymbol{g}_t\ _2^2}$ (Streeter & McMahan, 2010)
AdaGrad	Diagonal Matrices	$\left\ oldsymbol{x} ight\ _{\infty}$	$\sum_{i=1}^d \sqrt{\sum_{t=1}^T g_{t,i}^2}$	$\ \mathcal{X}\ _{\infty} \sum_{i=1}^{d} \sqrt{\sum_{t=1}^{T} g_{t,i}^2}$ (Streeter & McMahan, 2010)
Full-matrix AdaGrad	All Matrices	$\left\ oldsymbol{x} ight\ _2$	$ ext{Tr}[(\sum\limits_{t=1}^T oldsymbol{g}_t oldsymbol{g}_t^{ op})^{rac{1}{2}}]$	$\ \mathcal{X}\ _2 \operatorname{Tr}\left[\left(\sum_{t=1}^T \boldsymbol{g}_t \boldsymbol{g}_t^{\top}\right)^{\frac{1}{2}}\right]$ (Duchi et al., 2011)
One-sided Shampoo	$\mathbb{R}^{d_L imes d_L} \otimes I_{d_R}$	$\frac{\left\ \boldsymbol{X}\right\ _{\mathrm{op}}}{\sqrt{d_R}}$	$ ext{Tr}[\left(d_R\sum_{t=1}^T oldsymbol{G}_t oldsymbol{G}_t^{ op} ight)^{rac{1}{2}}]$	$\left\ \mathcal{X}\right\ _{\mathrm{op}} \operatorname{Tr}\left[\left(\sum_{t=1}^{T} \boldsymbol{G}_{t} \boldsymbol{G}_{t}^{\top}\right)^{\frac{1}{2}}\right]$ (Theorem 4.2)

Table 1. Regret bound obtained by Theorem 3.4 for different algorithms. Theorem 4.2 proves the rate specifically for one-sided shampoo and other rows match existing results in literature. \mathcal{H} is chosen as $\mathcal{K} \cap \mathcal{S}^d_+$. In the third column, \boldsymbol{x} denotes the parameter in vector form and \boldsymbol{X} denotes the parameter in matrix form with $\boldsymbol{x} = \overline{\operatorname{vec}}(\boldsymbol{X})$. In the fourth column, \boldsymbol{g}_t denotes the gradient in vector form and \boldsymbol{G}_t denotes the gradient in matrix form with $\boldsymbol{g}_t = \overline{\operatorname{vec}}(\boldsymbol{G}_t)$. We omit $O(\cdot)$ in complexity measures and regret bound for convenience.

than a more memory- and computationally efficient variant of these methods like diagonal AdaGrad and AdaGrad-Norm (Streeter & McMahan, 2010; Duchi et al., 2011; Ward et al., 2020). This demonstrates that more structure on the preconditioner may not necessarily restrict its optimization performance. In fact, our unified analysis uncovers an interesting finding: we show a simpler, more structured variant of Shampoo, called one-sided Shampoo, actually has substantially better regret bound compared to original Shampoo and full-matrix Adagrad (Section 3.3).

1.1. Main contributions

In light of the above discussion, we highlight the main contributions of the paper.

- We present a comprehensive and unified theoretical framework (Theorem 2.1) for adaptive optimization with structured preconditioners, encompassing several popular methods including diagonal AdaGrad, full-matrix AdaGrad, AdaGrad-Norm, and layerwise adaptive methods for both online convex optimization and stochastic smooth convex optimization. In particular, our analysis integrates and extends key insights from existing work (Gupta et al., 2017), which presents a unified and elegant way (Algorithm 1) to derive the aforementioned structured preconditioners, but only allows analysis on a case-by-case basis. To our knowledge, this is the *first* truly unified analysis of a large family of adaptive optimization algorithms.
- To enable a unified analysis, we identify a novel sufficient condition, *well-structured preconditioners* (Definition 3.1), which overcomes a key technical barrier in Gupta et al. (2017); thereby, allowing a unified analysis for several important adaptive algorithms. A more detailed discussion is provided in Section 2 and a summary of our results is presented in Table 1.

- Our unified analysis provides a new regret bound (Theorem 4.2) for a one-sided variant of Shampoo, which is always better than the existing bound for two-sided Shampoo (Gupta et al., 2018) and could be smaller by a multiplicative factor of d, where we assume the matrix-shaped parameter is of size √d-by-√d. This also leads to a novel convergence rate of one-sided shampoo for stochastic convex optimization (Theorem 4.4) via a standard offline-to-online reduction (Levy et al., 2018).
- Conceptually, our findings challenge the conventional wisdom that using a larger set of preconditioners which require more memory and compute leads to better optimization performance in terms of number of steps. In particular, while using a larger set of preconditioners reduces the gradient term $||| \boldsymbol{g}_{1:t} |||_{\mathcal{H}}$ in our regret bound, it increases the other term involving the magnitude of the optimal solution by increasing the norm metric; thus, leading to a worse regret bound (Theorem 2.1 and Table 1). For instance, we demonstrate that one-sided Shampoo can outperform full-matrix AdaGrad both theoretically (Section 4.3) and experimentally (Section 5). This suggest that one-sided Shampoo is not just a computational-efficient surrogate of full-matrix AdaGrad, but could also be fundamentally better, depending on the optimization problem.

1.2. Notations

Let \mathcal{M}^d be the set of all *d*-by-*d* matrices, and $\mathcal{S}^d \subset \mathcal{M}^d$ be the subset of all symmetric matrices. We use \mathcal{S}^d_+ to denote the set of positive semi-definite matrices, and \mathcal{S}^d_{++} to denote the set of positive definite matrices. \mathcal{D}^d is the set of all *d*dimensional diagonal matrices, and $\mathcal{D}^d_+ = \mathcal{D}^d \cap \mathcal{S}^d_+$. We denote by I_d the *d*-by-*d* identity matrix. For matrices A, B, we denote their inner product by $\langle A, B \rangle = \text{Tr}(A^\top B)$.

For any $\boldsymbol{H} \in \mathcal{S}^d_+$ such that $\boldsymbol{H} \neq 0$, we denote $\overline{\boldsymbol{H}} = \boldsymbol{H} / \operatorname{Tr}(\boldsymbol{H})$. For $\boldsymbol{H} \in \mathcal{S}^d_+$, $\|\boldsymbol{x}\|_{\boldsymbol{H}} := \sqrt{\boldsymbol{x}^\top \boldsymbol{H} \boldsymbol{x}}$ is the

(semi-)norm of $x \in \mathbb{R}^d$ with respect to H. For a convex set $\mathcal{H} \subseteq \mathcal{S}^d_+$, we define

$$\|\boldsymbol{x}\|_{\mathcal{H}} := \sup_{\boldsymbol{H} \in \mathcal{H}, \operatorname{Tr}(\boldsymbol{H}) \leq 1} \|\boldsymbol{x}\|_{\boldsymbol{H}}.$$
 (1)

For a convex set $\mathcal{X} \subseteq \mathbb{R}^d$ and any norm $\|\cdot\|$, we define $\|\mathcal{X}\| := \sup_{x \in \mathcal{X}} \|x\|$ and $\|\mathcal{X}\|_{\mathcal{H}} := \sup_{x \in \mathcal{X}} \|x\|_{\mathcal{H}}$. For any $H \succ 0$, the projection of x onto \mathcal{X} with respect to $\|\cdot\|_H$ is defined as $\Pi^H_{\mathcal{X}}(x) := \arg\min_{x' \in \mathcal{X}} \|x - x'\|_H$.

Throughout the paper, we consider the factorization $d = d_L d_R$, and denote the corresponding matrix form of $\boldsymbol{x} \in \mathbb{R}^d$ by $\boldsymbol{X} \in \mathbb{R}^{d_L \times d_R}$. We denote $\boldsymbol{x} = \overline{\operatorname{vec}}(\boldsymbol{X})$ and $\boldsymbol{X} = \overline{\operatorname{vec}}^{-1}(\boldsymbol{x})$ for conversion between the vector form and the matrix form. Then for a function $L(\boldsymbol{x})$ defined on $\boldsymbol{x} \in \mathbb{R}^d$, we extend its definition to matrices by letting $L(\boldsymbol{X})$ denote $L(\overline{\operatorname{vec}}(\boldsymbol{X}))$, and we will use $L(\boldsymbol{x})$ and $L(\boldsymbol{X})$ interchangeably when the context is clear. We define the gradient as $\boldsymbol{g}_t = \nabla L(\boldsymbol{x}_t)$ in vector form and $\boldsymbol{G}_t = \nabla L(\boldsymbol{X}_t)$ in matrix form, so $\boldsymbol{g}_t = \overline{\operatorname{vec}}(\boldsymbol{G}_t)$ and $\boldsymbol{G}_t = \overline{\operatorname{vec}}^{-1}(\boldsymbol{g}_t)$.

For a matrix $\boldsymbol{X} \in \mathbb{R}^{d_L \times d_R}$, we denote its operator norm as $\|\boldsymbol{X}\|_{\text{op}}$ and its Frobenious norm as $\|\boldsymbol{X}\|_{\text{F}}$. For a vector $\boldsymbol{x} \in \mathbb{R}^d$, we denotes its ℓ_{∞} norm by $\|\boldsymbol{x}\|_{\infty} = \max_{i \in [d]} |x_i|$ and ℓ_2 norm by $\|\boldsymbol{x}\|_2$.

2. Background: Unified Adaptive Regularization with Non-Unified Analysis

Seminal work Gupta et al. (2017) presented an adaptive regularization meta-algorithm, AdaReg (Algorithm 1), which can be used to derive various adaptive optimization algorithms known at that time in a unified approach. For example, AdaReg becomes full-matrix AdaGrad (Duchi et al., 2011), diagonal AdaGrad (Duchi et al., 2011), and AdaGrad-Norm (Streeter & McMahan, 2010; Ward et al., 2020) by choosing the set of preconditioners as the set of all PSD matrices, diagonal PSD matrices, and mutipliers of identity matrix respectively (see Table 1). The original AdaReg also allows other choices of potential function Φ , e.g., $\Phi(\cdot) = \log \det(\cdot)$ for Online Newton Step (Hazan et al., 2007), while we are only interested in the case of $\Phi(\cdot) = \eta^2 \operatorname{Tr}(\cdot)$ in this work.

In addition to the unified approach to deriving various adaptive optimization algorithms, Gupta et al. (2017) also attempts to give a unified analysis for the convergence rate or regret of these adaptive algorithms, which can be summarized by the following theorem.

Theorem 2.1 (Gupta et al. (2017)). Let $\{x_t\}_{t=1}^T$ be the iterates of Algorithm 1. Then for any $x^* \in \mathcal{X}$,

$$\sum_{t=1}^{T} L_t(\boldsymbol{x}_t) - \sum_{t=1}^{T} L_t(\boldsymbol{x}^*)$$

$$\leq \frac{1}{2} \left(\langle \boldsymbol{M}_T, \boldsymbol{H}_T^{-1} \rangle + \eta^2 \operatorname{Tr}(\boldsymbol{H}_T) - \eta^2 \operatorname{Tr}(\boldsymbol{H}_0) \right)$$

$$+ \frac{1}{2} \sum_{t=1}^{T} \left(\|\boldsymbol{x}_t - \boldsymbol{x}^*\|_{\boldsymbol{H}_t}^2 - \|\boldsymbol{x}_{t+1} - \boldsymbol{x}^*\|_{\boldsymbol{H}_t}^2 \right).$$
(2)

Algorithm 1 Adaptive Regularization Meta-Algorithm AdaReg (Gupta et al., 2017)

Hyperparam: $\epsilon > 0$, convex set $\mathcal{X} \subseteq \mathbb{R}^d$, learning rate η , preconditioners $\mathcal{H} \subset \mathcal{S}^d_+$ **Input:** initialization x_1 , loss functions $\{L_t\}_{t=1}^T : \mathbb{R}^d \to \mathbb{R}$ $M_0 \leftarrow \epsilon I_d$ **for** t = 1, 2, ..., T **do** $g_t \leftarrow \nabla L_t(x_t)$ $M_t \leftarrow M_{t-1} + g_t g_t^\top$ $H_t \leftarrow \arg \min_{H \in \mathcal{H}} \langle M_t, H^{-1} \rangle + \eta^2 \operatorname{Tr}(H)$ $x_{t+1} \leftarrow \Pi_{\mathcal{X}}^{H_t}(x_t - H_t^{-1}g_t)$ **Return** $x_1, ..., x_T$

The above bound is obtained by first applying a standard bound for online mirror descent to get a bound in the form of $\sum_{t=1}^{T} \|g_t\|_{H_t}$ plus the second term on the RHS of (2). Then the choice of H_t in Algorithm 1 enables the application of FTL-BTL lemma (Kalai & Vempala, 2005) to further bound $\sum_{t=1}^{T} \|g_t\|_{H_t}$ by the first term on the RHS of (2).

To proceed from Theorem 2.1, Gupta et al. (2017) relies on a crucial assumption that $H_{t-1} \preceq H_t$ for each t. Or more generally, for any $M \succ 0$, define

$$P_{\mathcal{H}}(\boldsymbol{M}) := \underset{\boldsymbol{H}\in\mathcal{H}}{\operatorname{arg\,min}} \langle \boldsymbol{M}, \boldsymbol{H}^{-1} \rangle + \eta^2 \operatorname{Tr}(\boldsymbol{H}). \quad (3)$$

Then we hope it holds that

$$P_{\mathcal{H}}(\boldsymbol{M}) \preceq P_{\mathcal{H}}(\boldsymbol{M}') \quad \text{for any} \quad 0 \prec \boldsymbol{M} \preceq \boldsymbol{M}' \quad (4)$$

In other words, we need $P_{\mathcal{H}} : S_{++}^d \to S_{++}^d$ to be *operator* monotone under the semi-definite ordering. With this assumption, a critical step in the derivation of the regret bound in Gupta et al. (2017) is to further rewrite and upper bound the second term on the RHS of Equation (2) by

$$\|\boldsymbol{x}_{1} - \boldsymbol{x}^{*}\|_{\boldsymbol{H}_{1}}^{2} + \sum_{t=2}^{T} \|\boldsymbol{x}_{t} - \boldsymbol{x}^{*}\|_{\boldsymbol{H}_{t} - \boldsymbol{H}_{t-1}}^{2}$$

$$\leq 4\|\boldsymbol{\mathcal{X}}\|_{\boldsymbol{H}_{1}}^{2} + 4\sum_{t=2}^{T} \|\boldsymbol{\mathcal{X}}\|_{\boldsymbol{H}_{t} - \boldsymbol{H}_{t-1}}^{2}.$$
(5)

Note that we need $H_{t-1} \leq H_t$ to ensure that $\|\cdot\|_{H_t-H_{t-1}}^2$ is indeed a (pseudo) norm and that the last inequality holds. Such analysis has been done for a few notable variants of AdaGrad in Gupta et al. (2017), where the condition $H_{t-1} \leq H_t$ is verified in a case-by-case way for specific choice of \mathcal{H} . However, the following question remains unclear for optimizers described by general \mathcal{H} :

Question 1. For a cone $\mathcal{H} \subseteq \mathcal{S}^d_+$, does $P_{\mathcal{H}}(M) \preceq P_{\mathcal{H}}(M')$ hold whenever $0 \prec M \preceq M'$?

Indeed, the answer is *no* for "ill-structured" \mathcal{H} , and we mention two negative examples here for illustration.

Example 2.2. Let $\mathcal{H} = \{ \boldsymbol{A} \otimes \boldsymbol{B} \succeq 0 \mid \boldsymbol{A} \in \mathcal{M}^{d_L}, \boldsymbol{B} \in \mathcal{M}^{d_R} \}$ with $d_L d_R = d$, i.e., the set of preconditioners for

two-sided Shampoo (Algorithm 3). We show that in the special case where $d_L = d_R = 2$, for $\mathbf{M} = \text{diag}(1, \epsilon, \epsilon, \epsilon)$ and $\mathbf{M}' = \text{diag}(1, \epsilon, \epsilon, 1)$, $P_{\mathcal{H}}(\mathbf{M}) \preceq P_{\mathcal{H}}(\mathbf{M}')$ does not hold for sufficiently small $\epsilon > 0$, although we have $\mathbf{M} \preceq \mathbf{M}'$. See Appendix A.4.1 for a detailed proof.

Example 2.3. The second example of ill-structured \mathcal{H} involves tridiagonal PSD matrices, i.e., matrices that only have nonzero elements on the main diagonal, the first diagonal above the main diagonal, and the first diagonal below the main diagonal. Specifically, for \mathcal{H} containing all 3-dimensional PSD matrices, we provide numerical evidence demonstrating that the desired condition Equation (4) breaks for very simple instances. See details in Appendix A.4.2.

These failure modes naturally lead to the second question:

Question 2. Is there a sufficient yet general condition (which covers all existing examples) on H for the inequality to hold?

As one of the main contributions of our work, we give an affirmative answer to Question 2 in Section 3 by proposing a notion of well-structured preconditioners (Definition 3.1) and deriving a unified analysis correspondingly.

3. Unified Analysis for Well-Structured Preconditioners

We establish a unified framework for adaptive optimization with structured preconditioners. In Section 3.1, we propose the notion of well-structured preconditioners and show that they satisfy the desired condition Equation (4). Then in Section 3.2, we present a unified analysis for adaptive optimization with well-structured preconditioners. We discuss several prominent examples in Section 3.3.

3.1. Well-structured preconditioners

For a set of *d*-by-*d* matrices $\mathcal{K} \subseteq \mathcal{M}^d$, we say that \mathcal{K} is a *subalgebra* if it is closed under scalar multiplication, matrix addition, and matrix multiplication. More concretely, we require that for any $\alpha \in \mathbb{R}$ and $A, B \in \mathcal{K}$, it holds that $\alpha A, AB, A + B \in \mathcal{K}$. Based on this, we propose the following core concept of our paper.

Definition 3.1 (Well-structured preconditioner sets). $\mathcal{H} \subseteq S^d_+$ is said to be a *well-structured preconditioner set* if $\mathcal{H} = S^d_+ \cap \mathcal{K}$ for some matrix subalgebra $\mathcal{K} \subseteq \mathcal{M}^d$ with $I_d \in \mathcal{K}$.

As a positive response to Question 2, the following proposition shows that our notion of well-structured preconditioner sets provides a sufficient condition for $P_{\mathcal{H}}(\cdot)$ to be operator monotone. See Appendix A for a proof.

Proposition 3.2. Let \mathcal{H} be a well-structured preconditioner set under Definition 3.1. For any $\mathbf{M} \succ 0$, there exists a unique solution $P_{\mathcal{H}}(\mathbf{M}) \succ 0$ to the optimization problem in (3). Furthermore, for any $M \succ 0$, $P_{\mathcal{H}}(M)$ satisfies the following properties:

(a)
$$\langle \boldsymbol{M}, P_{\mathcal{H}}(\boldsymbol{M})^{-1} \rangle = \eta^2 \operatorname{Tr}(P_{\mathcal{H}}(\boldsymbol{M}))$$

(b) $\overline{P_{\mathcal{H}}(\boldsymbol{M})} = \arg\min_{\boldsymbol{H}\in\mathcal{H},\operatorname{Tr}(\boldsymbol{H})\leq 1} \langle \boldsymbol{M}, \boldsymbol{H}^{-1} \rangle$, where we recall that $\overline{P_{\mathcal{H}}(\boldsymbol{M})} = P_{\mathcal{H}}(\boldsymbol{M}) / \operatorname{Tr}(P_{\mathcal{H}}(\boldsymbol{M}))$.

Moreover, for any $0 \prec \mathbf{M} \preceq \mathbf{M}'$, $P_{\mathcal{H}}(\mathbf{M}') - P_{\mathcal{H}}(\mathbf{M}) \in \mathcal{H}$ holds. In particular, it implies that $P_{\mathcal{H}}(\mathbf{M}) \preceq P_{\mathcal{H}}(\mathbf{M}')$.

The closure under both matrix addition and matrix multiplication is crucial in the proof of Proposition 3.2. Violation of any of the two properties could lead to problems: for the preconditioner set of two-sided Shampoo in Example 2.2, it is not closed under matrix addition; while for Example 2.3 regarding tridiagonal matrices, we note that the set of tridiagonal matrices is not closed under matrix multiplication.

3.2. Unified analysis for adaptive optimization

To proceed with the regret bound, recall from Proposition 3.2 that $\overline{H}_t = H_t / \text{Tr}(H_t)$ is a solution to the following constrained optimization problem

$$\overline{H}_t = \arg\min_{H \in \mathcal{H}, \operatorname{Tr}(H) \le 1} \langle M_t, H^{-1} \rangle.$$
(6)

Also recall that $M_t = \sum_{s=1}^t g_s g_s^{\top}$ (assuming $\epsilon = 0$ for illustration). We interpret the optimal value of this constrained optimization problem as the magnitude of the sequence of gradients $g_{1:t} = (g_1, \ldots, g_t)$ with respect to the best preconditioner in hindsight. Motivated by this, for any sequence of gradients $g_{1:t}$, we define their *adaptive gradient norm* with respect to \mathcal{H} to be

$$\|\|\boldsymbol{g}_{1:t}\|\|_{\mathcal{H}} := \inf_{\boldsymbol{H}\in\mathcal{H}, \operatorname{Tr}(\boldsymbol{H})\leq 1} \sqrt{\left\langle \sum_{s=1}^{t} \boldsymbol{g}_{s} \boldsymbol{g}_{s}^{\top}, \boldsymbol{H}^{-1} \right\rangle}.$$
 (7)

Indeed, this definition of adaptive gradient norm corresponds to the dual norm of the norm $\|\cdot\|_{\mathcal{H}\otimes I_t}$, denoted as $\|\cdot\|_{\mathcal{H}\otimes I_t}^*$. Specifically, we define

$$\|\overline{\operatorname{vec}}(\boldsymbol{g}_{1:t})\|_{\mathcal{H}\otimes \boldsymbol{I}_t}^* = \sup_{\boldsymbol{w}\in\mathbb{R}^{td}:\|\boldsymbol{w}\|_{\mathcal{H}\otimes \boldsymbol{I}_t}\leq 1}\overline{\operatorname{vec}}(\boldsymbol{g}_{1:t})^{\top}\boldsymbol{w}.$$

Then we have $\|\|\boldsymbol{g}_{1:t}\|\|_{\mathcal{H}} = \|\overline{\operatorname{vec}}(\boldsymbol{g}_{1:t})\|_{\mathcal{H}\otimes \boldsymbol{I}_t}^*/\sqrt{t}$ by the following lemma, which is proved in Appendix A.3.

Lemma 3.3. Let \mathcal{H} be a well-structured preconditioner set under Definition 3.1. Then for any $t \ge 1$ and $g_1, \ldots, g_t \in \mathbb{R}^d$, it holds that

$$\inf_{\boldsymbol{H}\in\mathcal{H},\mathrm{Tr}(\boldsymbol{H})\leq 1}\sqrt{\left\langle\sum_{s=1}^{t}\boldsymbol{g}_{s}\boldsymbol{g}_{s}^{\top},\boldsymbol{H}^{-1}\right\rangle}=\frac{1}{\sqrt{t}}\|\overline{\mathrm{vec}}(\boldsymbol{g}_{1:t})\|_{\mathcal{H}\otimes\boldsymbol{I}_{t}}^{*}$$

Given this definition of adaptive gradient norm, combining (6) and the fact that $\text{Tr}(\boldsymbol{H}_t) = \eta^{-2} \langle \boldsymbol{M}_t, \boldsymbol{H}_t^{-1} \rangle$ by Proposition 3.2, we obtain $\text{Tr}(\boldsymbol{H}_t) = \eta^{-1} ||| \boldsymbol{g}_{1:t} |||_{\mathcal{H}}$.

Now recall the upper bound in (5) for the second term in the regret bound (2) from Theorem 2.1. Note that for each t, Proposition 3.2 guarantees that $H_t - H_{t-1} \in \mathcal{H}$, so we can further bound $\|\mathcal{X}\|_{H_t - H_{t-1}}^2 \leq \|\mathcal{X}\|_{\mathcal{H}}^2 \cdot \operatorname{Tr}(H_t - H_{t-1})$. This allows us to telescope the sum in Equation (5) to get

$$\|\mathcal{X}\|_{H_{1}}^{2} + \sum_{t=2}^{T} \|\mathcal{X}\|_{H_{t}-H_{t-1}}^{2} = \frac{\|\mathcal{X}\|_{\mathcal{H}}^{2}}{\eta} \|g_{1:T}\|_{\mathcal{H}}$$
(8)

Observe that the upper bound is factored into two parts: 1) $\|\mathcal{X}\|_{\mathcal{H}}$, the diameter of the domain under norm $\|\cdot\|_{\mathcal{H}}$; and 2) $\|\|g_{1:T}\|\|_{\mathcal{H}}$, the adaptive gradient norm with respect to \mathcal{H} . We pause here for an important remark. Note that $\|\mathcal{X}\|_{H_t-H_{t-1}}^2$ is defined by taking the supremum over $x \in \mathcal{X}$, which precludes telescoping the sum over $t \in [T]$ at first glance. We address this issue by proposing the norm $\|\mathcal{X}\|_{\mathcal{H}}$, which allows us to extract the factor $\operatorname{Tr}(H_t-H_{t-1})$. Again, this unified analysis is possible thanks to Proposition 3.2 for well-structured preconditioner sets, in contrast to the caseby-case analysis done by Gupta et al. (2017). Furthermore, we remark that the above factored bound is crucial for us to identify the correct norm metric for Shampoo, leading to an improved analysis. See Section 4.1 for details.

Finally, combining (8) and the original regret bound in (2) yields the final regret bound for Algorithm 1. This is summarized in the following Theorem 3.4. The complete proof can be found in Appendix D.

Theorem 3.4. Let \mathcal{H} be a well-structured preconditioner set under Definition 3.1. Then for any convex loss functions L_1, \ldots, L_T , the regret of Algorithm 1 compared to any $x^* \in \mathcal{X}$ can be bounded as

$$\sum_{t=1}^{T} L_t(\boldsymbol{x}_t) - \sum_{t=1}^{T} L_t(\boldsymbol{x}^*) \le \left(\frac{D^2}{2\eta} + \eta\right) \left(G + d\sqrt{\epsilon}\right)$$

where $G = ||| g_{1:T} |||_{\mathcal{H}}, D = \max_{t \in [T]} || x_t - x^* ||_{\mathcal{H}}.$

Corollary 3.5. Under the setting of Theorem 2.1, further suppose that \mathcal{X} is a bounded set in \mathbb{R}^d . Then choosing $\eta = \sqrt{2} \|\mathcal{X}\|_{\mathcal{H}}$, the regret bound for Algorithm 1 becomes

$$\sum_{t=1}^{T} L_t(\boldsymbol{x}_t) - \sum_{t=1}^{T} L_t(\boldsymbol{x}^*) \leq 2\sqrt{2} \|\mathcal{X}\|_{\mathcal{H}} \left(G + d\sqrt{\epsilon}\right).$$

Ignoring ϵ , the bound reveals an intrinsic trade-off between $\|\mathcal{X}\|_{\mathcal{H}}$ and $\|\|g_{1:t}\|\|_{\mathcal{H}}$: as \mathcal{H} gets larger, $\|\mathcal{X}\|_{\mathcal{H}}$ increases while $\|\|g_{1:t}\|\|_{\mathcal{H}}$ decreases. The previous common belief that more adaptivity (larger \mathcal{H}) helps optimization could be largely due to the loose upper bound on $\|\mathcal{X}\|_{\mathcal{H}}$, i.e., always measuring the size of domain \mathcal{X} by Frobenius norm instead of potentially much smaller $\|\cdot\|_{\mathcal{H}}$. The fact that Kroneckered-factored subalgebra induces smaller $\|\cdot\|_{\mathcal{H}}$ than the entire matrix subalgebra is the key reason behind our surprising finding that one-sided Shampoo can outperform

full-matrix AdaGrad in terms of the number of steps. See a more detailed analysis in Section 4.3.

Next, we present the convergence rate of Algorithm 1 for stochastic convex smooth loss functions. We use Definition 3.6 to characterize the smoothness of a loss function. It is an extension of Φ -smoothness in (Xie et al., 2025), which only applies to block-diagonal preconditioners (corresponding to blockwise Adam).

Definition 3.6 (\mathcal{H} -smoothness). For a set $\mathcal{H} \subseteq S^d_+$ and any loss function $L : \mathbb{R}^d \to \mathbb{R}$, we define the \mathcal{H} smoothness of L, denoted by $H(L, \mathcal{H})$, as the smallest number $H \ge 0$ such that there exists a matrix $\mathbf{H}^* \in \mathcal{H}$ satisfying $H = \operatorname{Tr}(\mathbf{H}^*)$ and for any $\mathbf{x} \in \mathbb{R}^d$, it holds that $-\mathbf{H}^* \preceq \nabla^2 L(\mathbf{x}) \preceq \mathbf{H}^*$. In the case of convex L, this requirement becomes $\nabla^2 L(\mathbf{x}) \preceq \mathbf{H}^*$. Furthermore, we extend the notation to matrices $\mathbf{A} \in \mathcal{M}^d$ by defining $H(\mathbf{A}, \mathcal{H})$ as $H(\mathbf{x} \mapsto \frac{1}{2}\mathbf{x}^\top \mathbf{A}\mathbf{x}, \mathcal{H})$.

We need the following assumption on the stochastic noise.

Assumption 3.7. For any $t \in [T]$ and any $\boldsymbol{x} \in \mathcal{X}$, $\mathbb{E}[L_t(\boldsymbol{x})] = L(\boldsymbol{x})$ and there exists some $\boldsymbol{\Sigma} \in \mathcal{S}^d_+$ such that $\mathbb{E}[(\nabla L_t(\boldsymbol{x}) - \nabla L(\boldsymbol{x}))(\nabla L_t(\boldsymbol{x}) - \nabla L(\boldsymbol{x}))^\top] \preceq \boldsymbol{\Sigma}$.

Now we are ready to state our main results on the convergence rate of Algorithm 1 for stochastic convex functions.

Theorem 3.8. Let \mathcal{H} be a well-structured preconditioner set under Definition 3.1. Consider any independent stochastic convex loss functions L_1, \ldots, L_T satisfying Assumption 3.7, and let $H(L, \mathcal{H})$ be the \mathcal{H} -smoothness of their expectation L. Suppose the global minimizer of L, denoted by \mathbf{x}^* , is in \mathcal{X} . Then for the iterates $\mathbf{x}_1, \ldots, \mathbf{x}_T$ of Algorithm 1, denoting $\bar{\mathbf{x}}_{1:T} = \frac{1}{T} \sum_{t=1}^{T} \mathbf{x}_t$, it holds that

$$\mathbb{E}\left[L(\bar{\boldsymbol{x}}_{1:T}) - L(\boldsymbol{x}^*)\right] \leq \frac{16}{T} \left\|\mathcal{X}\right\|_{\mathcal{H}}^2 H(L,\mathcal{H}) + \frac{4\sqrt{2}}{\sqrt{T}} \left\|\mathcal{X}\right\|_{\mathcal{H}} \sigma + \frac{4\sqrt{2}d\sqrt{\epsilon}}{T} \left\|\mathcal{X}\right\|_{\mathcal{H}}$$

where $\sigma = \inf_{\boldsymbol{H} \in \mathcal{H}, \operatorname{Tr}(\boldsymbol{H}) \leq 1} \sqrt{\langle \boldsymbol{\Sigma}, \boldsymbol{H}^{-1} \rangle}.$

Our analysis naturally extends to Algorithm 4, which replaces the direct sum of past gradient outer products in Algorithm 1 with an exponential moving average (EMA). This modification is widely used in adaptive optimizers, including Adam (Kingma & Ba, 2014) and AdaSGD (Wang & Wiens, 2020), which is an ema version of AdaGrad. The detailed discussion is in Appendix C.

3.3. Examples of well-structured preconditioner sets

Next, we demonstrate that our Definition 3.1 is general enough to cover existing examples, by discussing several important matrix subalgebra \mathcal{K} . For each associated wellstructured $\mathcal{H} = \mathcal{K} \cap S^4_+$, recall the optimization problem over \mathcal{H} defined in (3). We will show that every minimizer $P_{\mathcal{H}}(\mathbf{M})$ corresponds to the preconditioner used in a specific adaptive optimization algorithm. We list the correspondence relationship below and the results are summarized in Table 1. The detailed derivation and calculation of the norms $\|\mathbf{x}\|_{\mathcal{H}}$ and $\|\|\mathbf{g}_{1:t}\|\|_{\mathcal{H}}$ can be found in Appendix B.

Example 3.9 (AdaGrad-Norm: scalar matrices). For the scalar matrix subalgebra $\mathcal{K} = \{c \cdot I_d \mid c \in \mathbb{R}\}$, we have $\mathcal{H} = \{c \cdot I_d \mid c \geq 0\}$. Then solving (3) for M_t yields

$$\boldsymbol{H}_{t} = \frac{1}{\eta} \sqrt{\mathrm{Tr}(\boldsymbol{M}_{t})/d} \cdot \boldsymbol{I}_{d} = \frac{1}{\eta} \sqrt{\epsilon + \sum_{s=1}^{t} \|\boldsymbol{g}_{s}\|_{2}^{2}/d \cdot \boldsymbol{I}_{d}}.$$

This is the preconditioner used in AdaGrad-Norm.

Example 3.10 (Diagonal AdaGrad: diagonal matrices). For the diagonal matrix subalgebra $\mathcal{K} = \mathcal{D}^d$, we have $\mathcal{H} = \mathcal{D}^d_+$. Correspondingly,

$$\boldsymbol{H}_{t} = \frac{1}{\eta} \operatorname{diag}\left(\sqrt{\epsilon + \sum_{s=1}^{t} \|g_{s,i}\|_{2}^{2} : i \in [d]} \right).$$

This is the preconditioner used in diagonal AdaGrad.

Example 3.11 (Full-matrix AdaGrad: all matrices). For $\mathcal{K} = \mathcal{M}^d$, we have $\mathcal{H} = \mathcal{S}^d_+$. In this case, solving (3) for M_t yields that $H_t = \frac{1}{\eta} M_t^{\frac{1}{2}}$, which corresponds to the update rule of full-matrix AdaGrad.

Example 3.12 (One-sided Shampoo: factored matrices). Let *d* be factored as $d = d_L d_R$. Then for the factored matrix algebra $\mathcal{K} = \mathbb{R}^{d_L \times d_L} \otimes I_{d_R}$, we have $\mathcal{H} = S^{d_L}_+ \otimes I_{d_R}$. Now writing $G_t \in \mathbb{R}^{d_L \times d_R}$ as the matricized version of $g_t \in \mathbb{R}^d$, solving the corresponding problem in (3) leads to

$$\boldsymbol{H}_{t} = \frac{1}{\eta} \left(\boldsymbol{\epsilon} \cdot \boldsymbol{I}_{d_{L}} + \frac{1}{d_{R}} \sum_{s=1}^{t} \boldsymbol{G}_{s} \boldsymbol{G}_{s}^{\top} \right)^{\frac{1}{2}} \otimes \boldsymbol{I}_{d_{R}}$$

which updates x_t the same as one-sided Shampoo¹ displayed in Algorithm 2, where we write the algorithm in the matrix form for convenience. More specifically, note that $g_t = \overline{\text{vec}}(G_t)$ and $H_t = L_t^{\frac{1}{2}} \otimes I_{d_R}$. Moreover, M_t corresponds to L_t by the fact that $\langle M_t, (H_L \otimes I_{d_R})^{-1} \rangle = \langle L_t, H_L^{-1} \rangle$ for any $H_L \in \mathcal{S}_{++}^{d_L}$.

The detailed calculations of the norms $\|\cdot\|_{\mathcal{H}}$ and $\|\cdot\|_{\mathcal{H}}$ for the above examples can be found in Appendix B.

Generate new well-structured preconditioner sets. Beyond the previous examples, it is also possible to generate new well-structured preconditioner sets based on existing ones. We discuss in particular an example of layerwise combination for parameters in a neural network. Specifically, for d parameters of an N-layer neural network, decompose $\mathbb{R}^d = \mathbb{R}^{d_1} \times \cdots \times \mathbb{R}^{d_N}$ where $d = \sum_{n=1}^N d_n$, and each \mathbb{R}^{d_n} corresponds to the d_n parameters in the n-th layer. For each $n \in [N]$, let $\mathcal{K}_n \subseteq \mathcal{M}^{d_n}$ be a matrix subalgebra. Then we define $\mathcal{K} = \bigoplus_{n=1}^N \mathcal{K}_n = \{\bigoplus_{n=1}^N \mathcal{A}_n \mid \mathcal{A}_n \in \mathcal{K}_n, n \in [N]\}$, and it is easy to verify that \mathcal{K} is also a subalgebra. Then the corresponding well-structured preconditioner set $\mathcal{H} = \bigoplus_{n=1}^N (\mathcal{K}_n \cap \mathcal{S}_+^{d_n}) = \bigoplus_{n=1}^N \mathcal{H}_n$ contains preconditioners that apply individual types of transforms to gradient of parameters in different layers². Such \mathcal{H} made up by direct sum of smaller cones also has very compositional property in its induced complexity metrics, namely, $\|\cdot\|_{\mathcal{H}} = \max_{1 \leq n \leq N} \|\cdot\|_{\mathcal{H}_n}$ and $\|\|\cdot\|_{\mathcal{H}} = \sum_{n=1}^N \|\cdot\|_{\mathcal{H}_n}$. This operation provides a useful tool for designing layerwise preconditioning methods (Bernstein & Newhouse, 2024a).

Other possible operations that can generate new matrix subalgebra \mathcal{K}' from the original subalgebra \mathcal{K} include taking Kronecker product with the identity matrix, i.e. $\mathcal{K}' = \{A' = A \otimes I_{d'} \mid A \in \mathcal{K}\}$, and rotation by an orthogonal matrix, i.e. $\mathcal{K}' = \{A' = U^{\top}AU \mid A \in \mathcal{K}\}$ where Uis an orthogonal matrix.

4. Improved Convergence Analysis for One-sided Shampoo

We now turn to one-sided Shampoo (Algorithm 2), a special example of Algorithm 1. In Section 4.1, we present our main results on its regret bound and convergence rate, and then compare to the previous results for two-sided Shampoo in Section 4.2. In Section 4.3, we present a comprehensive comparison between the regret bound of one-sided Shampoo and those of the AdaGrad variants, which suggests why Shampoo can outperform other adaptive algorithms on some real tasks.

4.1. Our results for one-sided Shampoo

We first characterize the norm $\|\cdot\|_{\mathcal{H}}$ for one-sided Shampoo.

Lemma 4.1 ($\|\cdot\|_{\mathcal{H}}$ for one-sided Shampoo). *Recall from Example 3.12 that for one-sided Shampoo (Algorithm 2),* $\mathcal{H} = S_{+}^{d_L} \otimes \mathbf{I}_{d_R}$ where $d_L d_R = d$. Then for any $\mathbf{x} \in \mathbb{R}^d$, it holds that $\|\mathbf{x}\|_{\mathcal{H}} = \frac{1}{\sqrt{d_R}} \|\mathbf{X}\|_{\text{op}}$ with $\mathbf{X} = \overline{\text{vec}}^{-1}(\mathbf{x})$, and thus $\|\mathcal{X}\|_{\mathcal{H}} = \frac{1}{\sqrt{d_R}} \|\mathcal{X}\|_{\text{op}}$.

See Appendix B.4 for the proof.

With this, we can apply Theorem 3.4 to get the regret bound for one-sided Shampoo, as summarized below in Theorem 4.2. See Appendix E.2 for its proof.

Theorem 4.2 (Regret bound for one-sided Shampoo). For convex functions L_1, \ldots, L_T , the regret of one-sided Sham-

¹We add a normalized factor $\frac{1}{d_R}$ compared to the L_t in Algorithm 3 so that it can be exactly derived from Algorithm 1.

²This can be applied to any partition of the parameters, not only for partition based on layers.

Algorithm 2 One-sided Shampoo

Hyperparam: learning rate $\eta > 0$, convex set $\mathcal{X} \subseteq \mathbb{R}^{d_L \times d_R}$, $L_0 = \epsilon I_{d_L}$ for $\epsilon \ge 0$ Input: initialization x_0 , stochastic loss functions $\{L_t\}_{t=1}^T : \mathbb{R}^{d_L \times d_R} \to \mathbb{R}$ for $t = 1, 2, \cdots, T$ do $G_t \leftarrow \nabla L_t(X_{t-1})$ $L_t \leftarrow L_{t-1} + \frac{1}{d_R} G_t G_t^\top$ $X_t \leftarrow \Pi_{\mathcal{X}}^{\frac{1}{2} \otimes I_{d_R}} (X_{t-1} - \eta_t L_t^{-\frac{1}{2}} G_t)$ Return x_T

poo (Algorithm 2) compared to any $X^* \in \mathbb{R}^{d_L \times d_R}$ satisfies

$$\sum_{t=1}^{T} L_t(\boldsymbol{X}_t) - \sum_{t=1}^{T} L_t(\boldsymbol{X}^*) \le \left(\frac{D_{\text{op}}^2}{2d_R\eta} + \eta\right) \left(G + d\sqrt{\epsilon}\right),$$

where $D_{\text{op}} = \max_{t \in [T]} \| \mathbf{X}_t - \mathbf{X}^* \|_{\text{op}}$ and $G = \sqrt{d_R} \operatorname{Tr} \left[\left(\sum_{t=1}^T \mathbf{G}_t \mathbf{G}_t^\top \right)^{\frac{1}{2}} \right]$. When the domain \mathcal{X} is bounded in operator norm, i.e., $\| \mathcal{X} \|_{\text{op}} < \infty$, further choosing $\eta = \sqrt{2/d_R} \| \mathcal{X} \|_{\text{op}}$, it holds

$$\sum_{t=1}^{T} L_t(\boldsymbol{X}_t) - \sum_{t=1}^{T} L_t(\boldsymbol{X}^*)$$
$$\leq 2\sqrt{2} \|\mathcal{X}\|_{\text{op}} \left(\text{Tr} \left[\left(\sum_{t=1}^{T} \boldsymbol{G}_t \boldsymbol{G}_t^\top \right)^{\frac{1}{2}} \right] + \frac{d}{\sqrt{d_R}} \sqrt{\epsilon} \right)$$

Next, before presenting the convergence rate of one-sided Shampoo, we first provide a more interpretable formulation of the \mathcal{H} -smoothness for one-sided Shampoo, which we call *left smoothness*. See Appendix E.1 for a proof of Lemma 4.3.

Lemma 4.3 (Left smoothness for one-sided Shampoo). Let $\mathcal{H} = S_+^{d_L} \otimes \mathbf{I}_{d_R}$ be the well-structured preconditioner set for one-sided Shampoo. Then the \mathcal{H} -smoothness $H(L, \mathcal{H})$ defined in Definition 3.6 is equal to the smallest number $H \ge 0$ such that there exists $\mathbf{H}_{d_L}^* \in \mathbb{R}^{d_L \times d_L}$ satisfying that $H = d_R \operatorname{Tr}(\mathbf{H}_{d_L}^*)$ and that for any $\mathbf{X}, \mathbf{\Delta} \in \mathbb{R}^{d_L \times d_R}$,

$$\left|
abla^2 L(oldsymbol{X}) [oldsymbol{\Delta},oldsymbol{\Delta}]
ight| \leq \left\langle oldsymbol{H}_{d_L}^*,oldsymbol{\Delta}oldsymbol{\Delta}^ op
ight
angle.$$

In this case, the \mathcal{H} -smoothness is denoted by $H_{\text{left}}(L)$.

Then the convergence rate of one-sided Shampoo can be obtained by specializing Theorem 3.8 to $\|\mathcal{X}\|_{\mathcal{H}} = \frac{1}{\sqrt{d_R}} \|\mathcal{X}\|_{\text{op}}$ from Lemma 4.1 and the left smoothness from Lemma 4.3.

Theorem 4.4 (Convergence rate of one-sided Shampoo). Let L_1, \ldots, L_T be stochastic convex loss functions satisfying Assumption 3.7, and let x_1, \ldots, x_T be the corresponding iterates of one-sided Shampoo (Algorithm 2) with learn-

Algorithm 3 Two-sided Shampoo (Gupta et al., 2018)

Hyperparam: learning rate η , convex set $\mathcal{X} \subseteq \mathbb{R}^{d_L \times d_R}$, $L_0 = \epsilon I_{d_L}$, $R_0 = \epsilon I_{d_R}$ for $\epsilon \ge 0$ Input: initialization x_0 , stochastic loss functions $\{L_t\}_{t=1}^T : \mathbb{R}^{d_L \times d_R} \to \mathbb{R}$ for $t = 1, 2, \dots, T$ do $G_t \leftarrow \nabla L_t(X_{t-1})$ $L_t \leftarrow L_{t-1} + G_t G_t^\top$ $R_t \leftarrow R_{t-1} + G_t^\top G_t$ $X_t \leftarrow \Pi_{\mathcal{X}}^{L_t^{\frac{1}{4}} \otimes R_t^{\frac{1}{4}}} (X_{t-1} - \eta_t L_t^{-\frac{1}{4}} G_t R_t^{-\frac{1}{4}})$ Return x_T

ing rate
$$\eta = \sqrt{2} \|\mathcal{X}\|_{\text{op}}$$
. Then for $\bar{x}_{1:T} = \frac{1}{T} \sum_{t=1}^{T} x_t$,
 $\mathbb{E}[L(\bar{x}_{1:T}) - L(x^*)]$
 $\leq \frac{16}{Td_R} \|\mathcal{X}\|_{\text{op}}^2 H_{\text{left}}(L) + \frac{4\sqrt{2}\sigma}{\sqrt{Td_R}} \|\mathcal{X}\|_{\text{op}} + \frac{4\sqrt{2}d\sqrt{\epsilon}}{T} \|\mathcal{X}\|_{\text{op}}$

where $\sigma = \inf_{\mathbf{H} \in \mathcal{H}, \operatorname{Tr}(\mathbf{H}) \leq 1} \sqrt{\langle \boldsymbol{\Sigma}, \mathbf{H}^{-1} \rangle}$, and $H_{\text{left}}(L)$ is the left smoothness of the expected loss L identified in Lemma 4.3.

4.2. Comparison with previous results on Shampoo

We compare our main results for one-sided Shampoo with the original results in Gupta et al. (2018). Here, we restate their original regret bound for easier comparison.

Theorem 4.5 (Regret bound of two-sided Shampoo (Gupta et al., 2018)). For convex functions $\{L_t\}_{t=1}^T$, suppose their gradients $(G_t = \nabla L_t(\mathbf{X}_t))_{t=1}^T$ are matrices of rank at most r. Then the regret of two-sided Shampoo (Algorithm 3³) compared to any $\mathbf{X}^* \in \mathbb{R}^{d_L \times d_R}$ is bounded as

$$\sum_{t=1}^{T} L_t(\boldsymbol{X}_t) - \sum_{t=1}^{T} L_t(\boldsymbol{X}^*) \le \left(\frac{D_{\mathrm{F}}^2}{2\eta} + r\eta\right) \operatorname{Tr}(\boldsymbol{L}_T^{\frac{1}{4}}) \operatorname{Tr}(\boldsymbol{R}_T^{\frac{1}{4}}),$$

where $D_{\mathrm{F}} = \max_{t \in [T]} \| \boldsymbol{X}_t - \boldsymbol{X}^* \|_{\mathrm{F}}, \ \boldsymbol{L}_T = \epsilon \boldsymbol{I}_{d_L} + \sum_{t=1}^T \boldsymbol{G}_t \boldsymbol{G}_t^\top$, and $\boldsymbol{R}_T = \epsilon \boldsymbol{I}_{d_R} + \sum_{t=1}^T \boldsymbol{G}_t^\top \boldsymbol{G}_t$. When $\| \boldsymbol{\mathcal{X}} \|_{\mathrm{F}} < \infty$, we further choose $\eta = \sqrt{2/r} \| \boldsymbol{\mathcal{X}} \|_{\mathrm{F}}$, then

$$\sum_{t=1}^{T} L_t(\boldsymbol{X}_t) - \sum_{t=1}^{T} L_t(\boldsymbol{X}^*) \le \sqrt{2r} \|\mathcal{X}\|_{\mathrm{F}} \operatorname{Tr}(\boldsymbol{L}_T^{\frac{1}{4}}) \operatorname{Tr}(\boldsymbol{R}_T^{\frac{1}{4}}).$$

We now compare our regret bound in Theorem 4.2 and the original regret bound in Theorem 4.5 by Gupta et al. (2018) when $\epsilon = 0$, i.e., $\boldsymbol{L}_T = \sum_{t=1}^T \boldsymbol{G}_t \boldsymbol{G}_t^\top$ and $\boldsymbol{R}_T = \sum_{t=1}^T \boldsymbol{G}_t^\top \boldsymbol{G}_t$. For a matrix $\boldsymbol{M} \in \mathcal{S}_+^d$, it always holds that $\sqrt{\mathrm{Tr}(\boldsymbol{M})} \leq \mathrm{Tr}(\boldsymbol{M}^{\frac{1}{2}})$. Therefore,

$$\begin{aligned} \operatorname{Tr}(\boldsymbol{L}_{T}^{\frac{1}{4}})\operatorname{Tr}(\boldsymbol{R}_{T}^{\frac{1}{4}}) &\geq \operatorname{Tr}(\boldsymbol{L}_{T}^{\frac{1}{4}})\operatorname{Tr}(\boldsymbol{R}_{T})^{\frac{1}{4}} = \operatorname{Tr}(\boldsymbol{L}_{T}^{\frac{1}{4}})\operatorname{Tr}(\boldsymbol{L}_{T})^{\frac{1}{4}} \\ &\geq \operatorname{Tr}(\boldsymbol{L}_{T}^{\frac{1}{4}}) \|\boldsymbol{L}_{T}\|_{\operatorname{op}}^{\frac{1}{4}} \geq \operatorname{Tr}(\boldsymbol{L}_{T}^{\frac{1}{2}}). \end{aligned}$$

³The original two-sided Shampoo analysis in (Gupta et al., 2017) is without per step projection to the bounded domain. We adapt their theorem into the projected version in a standard way.

This implies that regret bound in Theorem 4.5 is always no smaller than the regret bound in Theorem 4.2 because $r \ge 1$ and $\|\mathcal{X}\|_{\mathrm{F}} \ge \|\mathcal{X}\|_{\mathrm{op}}$ as $\|\mathbf{X}\|_{\mathrm{F}} \ge \|\mathbf{X}\|_{\mathrm{op}}$ for any $\mathbf{X} \in \mathcal{X}$.

Moreover, in the worst case, the Frobenius norm can be $\sqrt{\min(d_L, d_R)}$ times larger than the operator norm. As a concrete example, suppose G_t satisfies that $G_t[i, j] = 1$ only for $(j - i) \equiv t \pmod{\min(d_L, d_R)}$ and all other elements are zero. Then each G_t has rank $r = \min(d_L, d_R)$. At step $T = d_L d_R$, $L_T = d_R \cdot \min(d_L, d_R) \cdot I_{d_L}$ and $R_T = d_L \cdot \min(d_L, d_R) \cdot I_{d_R}$, and thus $\operatorname{Tr}(L_T^{\frac{1}{4}}) \operatorname{Tr}(R_T^{\frac{1}{4}}) = d_L^{\frac{5}{4}} d_R^{\frac{5}{4}}$ while $\operatorname{Tr}(L_T^{\frac{1}{2}}) = d_L d_R^{\frac{1}{2}}$. In this case, the regret bound of two-sided Shampoo is $\min(d_L, d_R) d_L^{\frac{1}{4}} d_R^{\frac{3}{4}}$ times larger than the our regret bound for one-sided Shampoo.

Duvvuri et al. (2024) introduced CASPR as an alternative to Shampoo for approximating full-matrix AdaGrad and achieved the same regret bound as Theorem 4.5. Our previous comparison thus applies to both Shampoo and CASPR.

4.3. Comparison with AdaGrad variants

Next we show one-sided Shampoo can achieve the best theoretical upper bound of the suboptimality gap for a specific class of functions, where each loss function has the form

$$L(\boldsymbol{X}) = \langle \boldsymbol{H}, (\boldsymbol{X} - \boldsymbol{X}^*) (\boldsymbol{X} - \boldsymbol{X}^*)^\top \rangle$$
(9)

where $X \in \mathbb{R}^{d_L \times d_R}$, $H \in \mathcal{S}^{d_L}_+$ with $\operatorname{Tr}(H) \leq 1$, and $\|X^*\|_{\operatorname{op}} \leq 1$. For this function class, we compare the largest possible value of convergence rate of different algorithms given by Theorem 3.8, which is summarized in Table 2.

For each algorithm defined by Algorithm 1 with specific \mathcal{H} , we will pick $\mathcal{X} = \{ \boldsymbol{x} \mid \|\boldsymbol{x}\|_{\mathcal{H}} \leq \|\boldsymbol{x}^*\|_{\mathcal{H}} \}$ so that the global minimizer \boldsymbol{x}^* is reachable. To get the convergence rate, it boils down to calculate $\|\mathcal{X}\|_{\mathcal{H}} = \|\boldsymbol{x}^*\|_{\mathcal{H}}$ and $H(L, \mathcal{H})$ associated with each algorithm. We have already derived the explicit form of $\|\mathcal{X}\|_{\mathcal{H}}$ for each algorithm in Section 3.3, as shown in the third column of Table 2. For the \mathcal{H} -smoothness, note that $\nabla^2 L(\boldsymbol{X}) = \boldsymbol{H} \otimes \boldsymbol{I}_{d_R}$ for the loss L in (9), and then it is straightforward to calculate $H(L, \mathcal{H})$ according to Definition 3.6, as shown in the fourth column of Table 2.

Below we present the worst case of the convergence rate on this problem class. We can see that one-sided Shampoo is strictly better the other three adaptive algorithms.

One-sided Shampoo. The worst case of convergence rate for one-sided Shampoo is $\frac{1}{T}$.

AdaGrad-Norm. Since $\max_{\mathrm{Tr}(\boldsymbol{H})\leq 1} \lambda_{\max}(\boldsymbol{H}) = 1$ and $\max_{\|\boldsymbol{X}^*\|_{\mathrm{op}}\leq 1} \|\overline{\mathrm{vec}}(\boldsymbol{X}^*)\|_2 = \max_{\|\boldsymbol{X}^*\|_{\mathrm{op}}\leq 1} \|(\boldsymbol{X}^*)\|_{\mathrm{F}} = \sqrt{\min \{d_L, d_R\}}$, the worst case of convergence rate is $\min (d_L, d_R)/T$.

AdaGrad. For psd matrix $\boldsymbol{H} \in \mathbb{R}^{d_L \times d_L}$ with $\operatorname{Tr}(\boldsymbol{H}) \leq 1$,

it holds that $H \leq I_{d_L}$. Then we have

$$\max_{\operatorname{Tr}(\boldsymbol{H})\leq 1} H(\boldsymbol{H}, \mathcal{D}_{+}^{d_{L}}) = \max_{\operatorname{Tr}(\boldsymbol{H})\leq 1} \min_{\operatorname{diag}(\boldsymbol{D})\succeq \boldsymbol{H}} \operatorname{Tr}(\boldsymbol{D}) \leq d_{L}.$$

On the other hand, if we choose $\boldsymbol{H} = \mathbf{1}_{d_L} \mathbf{1}_{d_L}^\top / d_L$, for any diagonal $\boldsymbol{D} \succeq \boldsymbol{H}$, it holds that $\operatorname{Tr}(\boldsymbol{D}) = \mathbf{1}_{d_L}^\top \boldsymbol{D} \mathbf{1}_{d_L} \geq \mathbf{1}_{d_L}^\top \boldsymbol{H} \mathbf{1}_{d_L} = d_L$. So we prove that $\max_{\operatorname{Tr}(\boldsymbol{H}) \leq 1} H(\boldsymbol{H}, \mathcal{D}_+^{d_L}) = d_L$.

We also know that $\max_{\|\boldsymbol{X}^*\|_{op} \leq 1} \|\overline{\operatorname{vec}}(\boldsymbol{X}^*)\|_{\infty} \leq 1$ since $\|\boldsymbol{A}\|_{\infty} \leq \|\boldsymbol{A}\|_{op}$. If the only nonzero entry of \boldsymbol{X}^* is $\boldsymbol{X}_{1,1}^* = 1$, then $\|\boldsymbol{X}^*\|_{\infty} = \|\boldsymbol{X}^*\|_{op} = 1$. Overall, the worst case of convergence rate is $\frac{d_L d_R}{T}$.

Full-matrix AdaGrad. We have seen in AdaGrad-Norm that $\max_{\|\boldsymbol{X}^*\|_{op} \leq 1} \|\overline{\operatorname{vec}}(\boldsymbol{X}^*)\|_2 = \sqrt{\min\{d_L, d_R\}}$, so the worst case of convergence rate is $\frac{\min(d_L, d_R)d_R}{T}$.

To summarize, we have identified a class of optimization problems, namely loss functions like Equation (9), with Hessian of bounded trace and optimizer of bounded spectral norm, for which one-sided Shampoo has much better worst case convergence rate than any other adaptive algorithms that could be derived from our unified analysis.

5. Experiments

In this section we empirically demonstrate the superior performance of 1-sided shampoo over other variants of AdaReg (Algorithm 1) on a simple but natural setting. Moreover, such superior performance is predicted by our theoretical analysis in Section 4.3, which in turn validates the practical utility of our theory in guiding optimizer selection.

Setup. We consider a linear regression problem $||AX - y||_2^2$ where A is the data matrix and $y = AX^*$ is the label vector generated by ground-truth X^* . Thus, the loss function can be equivalently written as

$$f(\boldsymbol{X}) = \langle \boldsymbol{H}, (\boldsymbol{X} - \boldsymbol{X}^*) (\boldsymbol{X} - \boldsymbol{X}^*)^\top \rangle$$

which is the same function that we studied in Section 4.3 and show that 1-sided shampoo outperforms other adaptive algorithms. We consider $X \in \mathbb{R}^{d \times d}$ with $d = 10^3$. We set the eigenvalues of H by $\sigma_1 = \cdots = \sigma_{10} = 1$ and $\sigma_i = \frac{1}{(i-10)^2}$ for $11 \le i \le 10^3$. Each element of the solution X^* is independently sampled from $\mathcal{N}(0, \frac{1}{d})$. We run AdaGrad-Norm, AdaGrad, one-sided Shampoo and fullmatrix AdaGrad for 100 steps from initialization $X_0 = 0$. We also run the original Shampoo algorithm for comparison. Full-matrix AdaGrad is run in a memory-efficient way and the detail is in Appendix F.1. We will compare the last iterate loss and the average iterate loss separately. The learning rate is tuned over five seeds for last iterate loss and average iterate loss respectfully, selecting the one with the

Structured Preconditioners in Adaptive Optimization: A Unified Analysis

Algorithm	Subalgebra ${\cal K}$	$\left\ oldsymbol{x}^* ight\ _{\mathcal{H}}$	$H(L,\mathcal{H})$	Convergence Rate
AdaGrad-Norm	$\{c\boldsymbol{I}_d \mid c \in \mathbb{R}\}$	$rac{1}{\sqrt{d}} \left\ oldsymbol{x}^* ight\ _2$	$d\cdot\lambda_{\max}(oldsymbol{H})$	$\left\ \mathcal{X} \right\ _2^2 \lambda_{\max}(oldsymbol{H})/T$
AdaGrad	Diagonal matrices, $\mathcal{D}^d(\mathbb{R})$	$\ m{x}^*\ _\infty$	$d_R \cdot H(\boldsymbol{H}, \mathcal{D}^{d_L}_+)$	$\left\ \mathcal{X}\right\ _{\infty}^{2} d_{R} \cdot H(\boldsymbol{H}, \mathcal{D}_{+}^{d_{L}})/T$
Full-Matrix AdaGrad	All matrices, $\mathbb{R}^{d \times d}$	$\left\ oldsymbol{x}^* ight\ _2$	$d_R \operatorname{Tr}(\boldsymbol{H})$	$\left\ \mathcal{X} \right\ _2^2 d_R \operatorname{Tr}(\boldsymbol{H}) / T$
One-Sided Shampoo	$\mathbb{R}^{d_L imes d_L} \otimes oldsymbol{I}_{d_R}$	$rac{1}{\sqrt{d_R}} \left\ oldsymbol{X}^{oldsymbol{*}} ight\ _{\mathrm{op}}$	$d_R\operatorname{Tr}({oldsymbol{H}})$	$\left\ \mathcal{X} \right\ _{\mathrm{op}}^{2} \mathrm{Tr}(\boldsymbol{H}) / T$

Table 2. Convergence rate for the loss function $L(\mathbf{X}) = \langle \mathbf{H}, (\mathbf{X} - \mathbf{X}^*)(\mathbf{X} - \mathbf{X}^*)^\top \rangle$. The results can be obtained by Theorem 3.8 with $\sigma = 0$ and omitting ϵ . For each algorithm we pick the smallest domain which still ensures the minimizer lies in the domain. Here recall that \mathcal{D}^d is the set of all *d*-dimensional diagonal matrices, and $\mathcal{D}^d_+ = \mathcal{D}^d \cap \mathcal{S}^d_+$.

lowest average loss. We use precision float32 and set $\epsilon = 0$ for all the experiments.

We also run the EMA version of the adaptive algorithms, whose results are consistent as shown in Appendix F.2.

Results. The overall results are shown in Section 5. We can see that one-sided Shampoo greatly outperforms other algorithms, corroborating the theoretical analysis in Section 4.3. Moreover, full-matrix AdaGrad is apparently the worst, suggesting more or even full adaptivity does not always help optimization for fixed budget of training steps. One-sided Shampoo is also better than the original Shampoo algorithm.



Figure 1. For both the last iterate loss $f(\mathbf{X}_t)$ and the average iterate training $f(\frac{1}{t}\sum_{s=1}^{t} \mathbf{X}_s)$, one-sided Shampoo performs the best and Full-matrix AdaGrad performs the worst. Here $f(\mathbf{X}) = \langle \mathbf{H}, (\mathbf{X} - \mathbf{X}^*) (\mathbf{X} - \mathbf{X}^*)^{\top} \rangle$.

6. Related Work

Adaptive optimizers and structured preconditioners. The extensive costs for training large-scale deep learning models have motivated the development of efficient optimization algorithms, among which adaptive optimizers have been widely studied because of their ability to exploit the rich geometry of the loss landscape (Pascanu & Bengio, 2013; Martens & Grosse, 2015; Dozat, 2016; Loshchilov & Hutter, 2018; Shazeer & Stern, 2018; Reddi et al., 2019; You et al., 2019; Zhuang et al., 2020; Liu et al., 2023; Yuan et al., 2024). For the sake of memory and computational efficiency, many works involve approximations to full-matrix preconditioners (Ba et al., 2017; George et al., 2018; Martens et al., 2018; Yao et al., 2021; Jahani et al., 2021; Zhang et al., 2022; Duvvuri et al., 2024). Indeed, our results suggest that such compromises might not harm the performance of adaptive optimizers in practice, because more adaptivity is not always helpful, as discussed in Section 4.3 and Section 5.

Understanding Shampoo. There are also recent efforts to understand the Shampoo optimizer from various perspectives. Bernstein & Newhouse (2024b) interpret Shampoo as the steepest descent with respect to the spectral norm of the layerwise matrix-form parameters of the neural network. Recognizing such structures of matrix-form parameters and role of spectral-norm geometry in deep learning has led to development of new optimizer such as Muon (Jordan et al., 2024). In addition, the second-order perspective on Shampoo (Anil et al., 2020) has also led to fruitful results: Morwani et al. (2024) connect the preconditioner in Shampoo to the optimal Kronecker product approximation of the Gauss-Newton component of the Hessian, and Vyas et al. (2024) propose to view Shampoo.

7. Conclusion and Future Works

We present a unified analysis for a broad class of adaptive optimization algorithms with well-structured preconditioners (Definition 3.1) for both online regret minimization and smooth convex optimization. Our analysis not only provides matching rate to several important algorithms including diagonal AdaGrad, full-matrix AdaGrad, and AdaGradNorm, but also gives an improved convergence rate for a one-sided variant of Shampoo over that of the original Shampoo. We reveal a novel trade-off in final convergence rate between domain metric and adaptive gradient norm (Equation (7)) for regret minimization or adaptive smoothness (Definition 3.6) for smooth convex optimization. We hope this insight could be useful towards design of future adaptive optimizers. One important future direction is to identify more subalgebras or other structures that are useful for improving the performance by better adapting to the domain or loss smoothness.

Impact Statement

The goal of this paper is to advance the field of Machine Learning by providing theoretical and experimental insights for adaptive optimization algorithms. There are many potential societal consequences of our work, none of which we feel must be specifically highlighted here.

Acknowledgment

The authors sincerely thank Matt Streeter for his valuable discussions and insightful comments throughout this work.

References

- Anil, R., Gupta, V., Koren, T., Regan, K., and Singer, Y. Scalable second order optimization for deep learning. *arXiv preprint arXiv:2002.09018*, 2020. URL https: //arxiv.org/pdf/2002.09018.
- Ba, J., Grosse, R., and Martens, J. Distributed secondorder optimization using kronecker-factored approximations. In *International Conference on Learning Representations*, 2017. URL https://jimmylba.github. io/papers/nsync.pdf.
- Bernstein, J. and Newhouse, L. Modular duality in deep learning. *arXiv preprint arXiv:2410.21265*, 2024a. URL https://arxiv.org/pdf/2410.21265.
- Bernstein, J. and Newhouse, L. Old optimizer, new norm: An anthology. *arXiv preprint arXiv:2409.20325*, 2024b. URL https://arxiv.org/pdf/2409.20325.
- Dahl, G. E., Schneider, F., Nado, Z., Agarwal, N., Sastry, C. S., Hennig, P., Medapati, S., Eschenhagen, R., Kasimbeg, P., Suo, D., et al. Benchmarking neural network training algorithms. *arXiv preprint arXiv:2306.07179*, 2023. URL https://arxiv.org/pdf/2306.07179.
- Dozat, T. Incorporating nesterov momentum into adam. 2016. URL https://openreview.net/pdf? id=OM0jvwB8jIp57ZJjtNEZ.
- Duchi, J., Hazan, E., and Singer, Y. Adaptive subgradient methods for online learning and stochastic optimization. *Journal of machine learning research*, 2011. URL https://www.jmlr.org/papers/ volume12/duchilla/duchilla.pdf.
- Duvvuri, S. S., Devvrit, F., Anil, R., Hsieh, C.-J., and Dhillon, I. S. Combining axes preconditioners through kronecker approximation for deep learning. In *The Twelfth International Conference on Learning Representations*, 2024. URL https://openreview.net/ pdf?id=8j9hz8DVi8.

- George, T., Laurent, C., Bouthillier, X., Ballas, N., and Vincent, P. Fast approximate natural gradient descent in a kronecker factored eigenbasis. *Advances in Neural Information Processing Systems*, 2018. URL https: //arxiv.org/pdf/1806.03884.
- Gupta, V., Koren, T., and Singer, Y. A unified approach to adaptive regularization in online and stochastic optimization. *arXiv preprint arXiv:1706.06569*, 2017. URL https://arxiv.org/pdf/1706.06569.
- Gupta, V., Koren, T., and Singer, Y. Shampoo: Preconditioned stochastic tensor optimization. In *International Conference on Machine Learning*, 2018. URL https://arxiv.org/pdf/1802.09568.
- Hazan, E., Agarwal, A., and Kale, S. Logarithmic regret algorithms for online convex optimization. *Machine Learning*, 2007. URL https://link.springer.com/ article/10.1007/s10994-007-5016-8.
- Jahani, M., Rusakov, S., Shi, Z., Richtárik, P., Mahoney, M. W., and Takáč, M. Doubly adaptive scaled algorithm for machine learning using second-order information. arXiv preprint arXiv:2109.05198, 2021. URL https://arxiv.org/pdf/2109.05198.
- Jordan, K., Jin, Y., Boza, V., You, J., Cesista, F., Newhouse, L., and Bernstein, J. Muon: An optimizer for hidden layers in neural networks, 2024. URL https://web. archive.org/web/20250122060345/https: //kellerjordan.github.io/posts/muon/.
- Kalai, A. and Vempala, S. Efficient algorithms for online decision problems. Journal of Computer and System Sciences, 2005. URL https://www.sciencedirect.com/science/article/pii/S0022000004001394.
- Kingma, D. P. and Ba, J. Adam: A method for stochastic optimization. arXiv preprint arXiv:1412.6980, 2014. URL https://arxiv.org/pdf/1412.6980.pdf.
- Levy, K. Y., Yurtsever, A., and Cevher, V. Online adaptive methods, universality and acceleration. *Advances in neural information processing systems*, 2018. URL https://arxiv.org/pdf/1809.02864.
- Liu, H., Li, Z., Hall, D. L. W., Liang, P., and Ma, T. Sophia: A scalable stochastic second-order optimizer for language model pre-training. In *The Twelfth International Conference on Learning Representations*, 2023. URL https://arxiv.org/pdf/2305.14342.pdf.
- Loshchilov, I. and Hutter, F. Decoupled weight decay regularization. In *International Conference on Learning Representations*, 2018. URL https://arxiv.org/ pdf/1711.05101.pdf.

- Martens, J. and Grosse, R. Optimizing neural networks with kronecker-factored approximate curvature. In *International conference on machine learning*, 2015. URL http://arxiv.org/pdf/1503.05671.
- Martens, J., Ba, J., and Johnson, M. Kronecker-factored curvature approximations for recurrent neural networks. In International Conference on Learning Representations, 2018. URL https://openreview.net/ pdf?id=HyMTkQZAb.
- Morwani, D., Shapira, I., Vyas, N., Malach, E., Kakade, S., and Janson, L. A new perspective on shampoo's preconditioner. *arXiv preprint arXiv:2406.17748*, 2024. URL https://arxiv.org/pdf/2406.17748.
- Pascanu, R. and Bengio, Y. Revisiting natural gradient for deep networks. arXiv preprint arXiv:1301.3584, 2013. URL https://arxiv.org/pdf/1301.3584.
- Petersen, K. B., Pedersen, M. S., et al. The matrix cookbook. *Technical University of Denmark*, 2008. URL https://ece.uwaterloo.ca/~ece602/ MISC/matrixcookbook.pdf.
- Reddi, S. J., Kale, S., and Kumar, S. On the convergence of adam and beyond. arXiv preprint arXiv:1904.09237, 2019. URL https://arxiv.org/pdf/1904. 09237.
- Shazeer, N. and Stern, M. Adafactor: Adaptive learning rates with sublinear memory cost. In *International Conference on Machine Learning*, 2018. URL https: //arxiv.org/pdf/1804.04235.
- Streeter, M. and McMahan, H. B. Less regret via online conditioning. arXiv preprint arXiv:1002.4862, 2010. URL https://arxiv.org/pdf/1002.4862.
- Vyas, N., Morwani, D., Zhao, R., Shapira, I., Brandfonbrener, D., Janson, L., and Kakade, S. Soap: Improving and stabilizing shampoo using adam. arXiv preprint arXiv:2409.11321, 2024. URL https:// arxiv.org/pdf/2409.11321.
- Wang, J. and Wiens, J. Adasgd: Bridging the gap between sgd and adam. *arXiv preprint arXiv:2006.16541*, 2020. URL https://arxiv.org/pdf/2006.16541.
- Ward, R., Wu, X., and Bottou, L. Adagrad stepsizes: Sharp convergence over nonconvex landscapes. Journal of Machine Learning Research, 2020. URL https://www.jmlr.org/papers/ volume21/18-352/18-352.pdf.

- Xie, S., Mohamadi, M. A., and Li, Z. Adam exploits \$\ell_\infty\$-geometry of loss landscape via coordinatewise adaptivity. In *The Thirteenth International Conference on Learning Representations*, 2025. URL https: //openreview.net/forum?id=PUnD86UEK5.
- Yao, Z., Gholami, A., Shen, S., Mustafa, M., Keutzer, K., and Mahoney, M. Adahessian: An adaptive second order optimizer for machine learning. In *proceedings of the AAAI conference on artificial intelligence*, 2021. URL https://arxiv.org/pdf/2006.00719.
- You, Y., Gitman, I., and Ginsburg, B. Large batch training of convolutional networks. *arXiv preprint arXiv:1708.03888*, 2017. URL https://arxiv. org/pdf/1708.03888.
- You, Y., Li, J., Reddi, S., Hseu, J., Kumar, S., Bhojanapalli, S., Song, X., Demmel, J., Keutzer, K., and Hsieh, C.-J. Large batch optimization for deep learning: Training bert in 76 minutes. *arXiv preprint arXiv:1904.00962*, 2019. URL https://arxiv.org/pdf/1904.00962.
- Yuan, H., Liu, Y., Wu, S., Zhou, X., and Gu, Q. Mars: Unleashing the power of variance reduction for training large models. *arXiv preprint arXiv:2411.10438*, 2024. URL https://arxiv.org/pdf/2411.10438.
- Zhang, L., Shi, S., and Li, B. Eva: Practical second-order optimization with kronecker-vectorized approximation. In *The Eleventh International Conference on Learning Representations*, 2022. URL https://arxiv.org/ pdf/2308.02123.
- Zhuang, J., Tang, T., Ding, Y., Tatikonda, S. C., Dvornek, N., Papademetris, X., and Duncan, J. Adabelief optimizer: Adapting stepsizes by the belief in observed gradients. *Advances in neural information processing systems*, 2020. URL https://arxiv.org/pdf/2010.07468.

A. Proof for Well-Structured Preconditioner Sets

A.1. Definition of ${\cal H}$ and its properties

Recall that \mathcal{K} is a subalgebra of *d*-by-*d* real-valued matrices, and the corresponding well-structured preconditioner set is $\mathcal{H} = \mathcal{K} \cap S^d_+$.

Lemma A.1. For any $A \in \mathcal{K}$ and any polynomial $p, p(A) \in \mathcal{K}$. Furthermore, for any invertible $A \in \mathcal{K}$, its inverse $A^{-1} \in \mathcal{K}$. Also, for any symmetric $A \in \mathcal{K}$, its pseudo inverse $A^{\dagger} \in \mathcal{K}$.

Proof of Lemma A.1. The first statement follows from the fact that \mathcal{K} is a subalgebra, and the second and third statement are consequences of the Cayley-Hamilton theorem.

Lemma A.2. For $\mathcal{H} = \mathcal{K} \cap \mathcal{S}^d_+$ where \mathcal{K} is a subalgebra of d-by-d real-valued matrices, define

$$\mathcal{H}^* = \{ \boldsymbol{A} \in \mathcal{K} : \langle \boldsymbol{H}, \boldsymbol{A} \rangle \ge 0, \forall \boldsymbol{H} \in \mathcal{H} \}.$$
(10)

Then $\mathcal{H}^* = \mathcal{H}$. Consequently, for any $H_1, H_2 \in \mathcal{H}$, if $\langle H_1 - H_2, H \rangle \geq 0$ for all $H \in \mathcal{H}$, then $H_1 \succeq H_2$.

Proof of Lemma A.2. Suppose there exists $A \in \mathcal{H}^*$ such that A has a negative eigenvalue. Let $\lambda_1(A), \ldots, \lambda_d(A)$ be the eigenvalues of A in the decreasing order. Consider the matrix $B = (A - 2 \max(\lambda_1(A), 1)I_d)^2 \succ 0$. The leading eigenspace of B is the same as the eigenspace of A corresponding to its smallest eigenvalue, which is negative. Consequently, for large enough integer n, $B^n/||B^n||_{op}$ is approximately the projection matrix onto the eigenspace of the smallest negative eigenvalue of A. Therefore, for large enough integer n, $\langle B^n/||B^n||_{op}, A \rangle < 0$. However, since $B^n/||B^n||_{op}$ is a polynomial of A, we know that $B^n/||B^n||_{op} \in \mathcal{H}$. This is a contradiction to the definition of \mathcal{H}^* . Hence, we conclude that for any $A \in \mathcal{H}^*$, it holds that $A \succeq 0$, and thus $\mathcal{H}^* \subseteq \mathcal{H}$. Moreover, for any $A \in \mathcal{H}$, it holds that $\langle H, A \rangle \ge 0$ for any $H \in \mathcal{H}$ because both A and H are positive semi-definite. This shows that $\mathcal{H} \subseteq \mathcal{H}^*$, and hence $\mathcal{H}^* = \mathcal{H}$. This completes the proof.

A.2. Properties of $P_{\mathcal{H}}(\cdot)$

For any M > 0, recall the regularized optimization problem in (3). Note that we can assume $\eta = 1$ without loss of generality, because the original problem is equivalent to solve for M/η^2 with regularizer Tr(H) in place of $\eta^2 \text{Tr}(H)$. Therefore, in the rest of this section, we focus on the following optimization problem:

$$P_{\mathcal{H}}(\boldsymbol{M}) := \underset{\boldsymbol{H} \in \mathcal{H}}{\operatorname{arg\,min}} \langle \boldsymbol{M}, \boldsymbol{H}^{-1} \rangle + \operatorname{Tr}(\boldsymbol{H}).$$
(11)

The results proved below applies to the original problem (3) after simple rescaling, and Proposition 3.2 follows from Proposition A.3 below.

For notational convenience, given any $M \succ 0$, we define

$$f_{\boldsymbol{M}}(\boldsymbol{H}) := \langle \boldsymbol{M}, \boldsymbol{H}^{-1} \rangle + \operatorname{Tr}(\boldsymbol{H}).$$
(12)

Proposition A.3. Let \mathcal{H} be a well-structured preconditioner set under Definition 3.1. For any $M \succ 0$, there exists a unique solution $P_{\mathcal{H}}(M) \succ 0$ to the optimization problem in (11). Furthermore, for any $M \succ 0$, $P_{\mathcal{H}}(M)$ satisfies the following properties:

- (a) $\langle \boldsymbol{M}, P_{\mathcal{H}}(\boldsymbol{M})^{-1} \rangle = \operatorname{Tr}(P_{\mathcal{H}}(\boldsymbol{M})).$
- (b) $\overline{P_{\mathcal{H}}(M)} = \arg\min_{H \in \mathcal{H}, \operatorname{Tr}(H) < 1} \langle M, H^{-1} \rangle$ where we recall that $\overline{P_{\mathcal{H}}(M)} = \operatorname{Tr}(P_{\mathcal{H}}(M))^{-1} P_{\mathcal{H}}(M)$.
- (c) For any $\boldsymbol{H} \in \mathcal{H}$, $\langle -P_{\mathcal{H}}(\boldsymbol{M})^{-1}\boldsymbol{M}P_{\mathcal{H}}(\boldsymbol{M})^{-1} + \boldsymbol{I}_{d}, \boldsymbol{H} P_{\mathcal{H}}(\boldsymbol{M}) \rangle = 0.$

Moreover, for any $M_1 \succeq M_2 \succ 0$ *, it holds that* $P_{\mathcal{H}}(M_1) - P_{\mathcal{H}}(M_2) \in \mathcal{H}$ *, and in particular,* $P_{\mathcal{H}}(M_1) \succeq P_{\mathcal{H}}(M_2)$ *.*

Proof of Proposition A.3. We first show that $P_{\mathcal{H}}(\mathbf{M})$ exists and $P_{\mathcal{H}}(\mathbf{M}) \succ 0$. Note that for $\mathcal{H} = \mathcal{K} \cap \mathcal{S}^d_+$, since \mathcal{K} is a linear subspace of \mathcal{M}^d and \mathcal{S}^d_+ is a closed subset of \mathcal{M}^d , we know that \mathcal{H} is also a closed subset of \mathcal{M}^d . Moreover, for any sequence $\{\mathbf{H}_n\}_{n\geq 1}$ such that $\mathbf{H}_n \succ 0$ and either the smallest eigenvalue of \mathbf{H}_n converges to 0 or the largest eigenvalue of \mathbf{H}_n converges to ∞ , the objective value $f_{\mathbf{M}}(\mathbf{H}_n)$ goes to ∞ because $\mathbf{M} \succ 0$. Therefore, there exists $P_{\mathcal{H}}(\mathbf{M}) \succ 0$ that attains the minimum objective value. For the uniqueness of $P_{\mathcal{H}}(\mathbf{M})$, it suffices to note that the objective function $f_{\mathbf{M}}(\mathbf{H})$ is strictly convex in $\mathbf{H} \succ 0$ for $\mathbf{M} \succ 0$.

Proof for property (a). Suppose otherwise that $\langle M, P_{\mathcal{H}}(M)^{-1} \rangle \neq \operatorname{Tr}(P_{\mathcal{H}}(M))$, and consider the following matrix:

$$\boldsymbol{H} = P_{\mathcal{H}}(\boldsymbol{M}) \cdot \sqrt{\langle \boldsymbol{M}, P_{\mathcal{H}}(\boldsymbol{M})^{-1} \rangle / \operatorname{Tr}(P_{\mathcal{H}}(\boldsymbol{M}))} \in \mathcal{H}.$$

For this H, we have $\langle M, H^{-1} \rangle + \operatorname{Tr}(H) = 2\sqrt{\langle M, P_{\mathcal{H}}(M)^{-1} \rangle \cdot \operatorname{Tr}(P_{\mathcal{H}}(M))} < \langle M, P_{\mathcal{H}}(M)^{-1} \rangle + \operatorname{Tr}(P_{\mathcal{H}}(M))$, thus contradicting the optimality of $P_{\mathcal{H}}(M)$. Therefore, it must be true that $\langle M, P_{\mathcal{H}}(M)^{-1} \rangle = \operatorname{Tr}(P_{\mathcal{H}}(M))$.

Proof for property (b) Note that we can rewrite the original optimization problem as follows:

$$egin{aligned} P_{\mathcal{H}}(oldsymbol{M}) &= \operatorname*{arg\,min}_{oldsymbol{H}\in\mathcal{H}} \langle oldsymbol{M},oldsymbol{H}^{-1}
angle + \operatorname{Tr}(oldsymbol{H}) \ &= \operatorname*{arg\,min}_{oldsymbol{H}\in\mathcal{H}} \left\langle oldsymbol{M}, \left(rac{oldsymbol{H}}{\operatorname{Tr}(oldsymbol{H})}
ight)^{-1}
ight
angle \cdot rac{1}{\operatorname{Tr}(oldsymbol{H})} + \operatorname{Tr}(oldsymbol{H}). \end{aligned}$$

Note that solving the above optimization problem is equivalent to first solving $\overline{P_{\mathcal{H}}(M)} = \arg \min_{H \in \mathcal{H}, \operatorname{Tr}(H) \leq 1} \langle M, H^{-1} \rangle$ and then setting $P_{\mathcal{H}}(\underline{M}) = \operatorname{Tr}(P_{\mathcal{H}}(M)) \cdot \overline{P_{\mathcal{H}}(M)}$ where the value of $\operatorname{Tr}(P_{\mathcal{H}}(M))$ ensures the previous property (a). Hence, we see that $\overline{P_{\mathcal{H}}(M)}$ solves the constrained version of the original optimization problem.

Proof for property (c). Since $\nabla_{\boldsymbol{H}} \operatorname{Tr}(\boldsymbol{H}) = \boldsymbol{I}_d$ and $\nabla_{\boldsymbol{H}} \operatorname{Tr}(\boldsymbol{M}^\top \boldsymbol{H}^{-1}) = -\boldsymbol{H}^{-1}\boldsymbol{M}\boldsymbol{H}^{-1}$ (see e.g. Equation (124) in Petersen et al. (2008)), we have

$$\nabla_{\boldsymbol{H}} f_{\boldsymbol{M}}(\boldsymbol{H}) = -\boldsymbol{H}^{-1} \boldsymbol{M} \boldsymbol{H}^{-1} + \boldsymbol{I}_{d}.$$
⁽¹³⁾

Then as \mathcal{H} is a cone, by the optimality of $P_{\mathcal{H}}(M)$, it holds for any $H \in \mathcal{H}$ that

$$0 = \langle \nabla_{\boldsymbol{H}} f_{\boldsymbol{M}}(P_{\mathcal{H}}(\boldsymbol{M})), \boldsymbol{H} - P_{\mathcal{H}}(\boldsymbol{M}) \rangle = \langle -P_{\mathcal{H}}(\boldsymbol{M})^{-1} \boldsymbol{M} P_{\mathcal{H}}(\boldsymbol{M})^{-1} + \boldsymbol{I}_{d}, \boldsymbol{H} - P_{\mathcal{H}}(\boldsymbol{M}) \rangle$$

Proof for the operator monotonicity of $P_{\mathcal{H}}(\cdot)$. By property (c) of $P_{\mathcal{H}}(\cdot)$, we know that for any $M \succ 0$, $\langle M - P_{\mathcal{H}}(M)^2, P_{\mathcal{H}}(M)^{-1}HP_{\mathcal{H}}(M)^{-1} - P_{\mathcal{H}}(M)^{-1} \rangle = 0$. Note that $H \mapsto P_{\mathcal{H}}(M)^{-1}HP_{\mathcal{H}}(M)^{-1}$ is a bijection from \mathcal{H} to \mathcal{H} by Lemma A.1, so we have $\langle M - P_{\mathcal{H}}(M)^2, H - P_{\mathcal{H}}(M)^{-1} \rangle = 0$ for all $H \in \mathcal{H}$. Applying this to both M_1 and M_2 , it follows that for any $H \in \mathcal{H}$

$$\langle \boldsymbol{M}_1 - P_{\mathcal{H}}(\boldsymbol{M}_1)^2, \boldsymbol{H} - P_{\mathcal{H}}(\boldsymbol{M}_1)^{-1} \rangle = \langle \boldsymbol{M}_2 - P_{\mathcal{H}}(\boldsymbol{M}_2)^2, \boldsymbol{H} - P_{\mathcal{H}}(\boldsymbol{M}_2)^{-1} \rangle.$$

Rearranging the above equation, we obtain

$$\begin{split} \langle P_{\mathcal{H}}(\boldsymbol{M}_1)^2 - P_{\mathcal{H}}(\boldsymbol{M}_2)^2, \boldsymbol{H} \rangle &= \langle \boldsymbol{M}_1 - \boldsymbol{M}_2, \boldsymbol{H} \rangle - \langle \boldsymbol{M}_1, P_{\mathcal{H}}(\boldsymbol{M}_1)^{-1} \rangle + \operatorname{Tr}(P_{\mathcal{H}}(\boldsymbol{M}_1)) \\ &+ \langle \boldsymbol{M}_2, P_{\mathcal{H}}(\boldsymbol{M}_2)^{-1} \rangle - \operatorname{Tr}(P_{\mathcal{H}}(\boldsymbol{M}_2)) \\ &= \langle \boldsymbol{M}_1 - \boldsymbol{M}_2, \boldsymbol{H} \rangle \end{split}$$

where the second equality follows from the first property of $P_{\mathcal{H}}(M_1)$ and $P_{\mathcal{H}}(M_2)$ from Proposition A.3. Since $M_1 \succeq M_2$, this implies that $\langle P_{\mathcal{H}}(M_1)^2 - P_{\mathcal{H}}(M_2)^2, H \rangle \ge 0$ for all $H \in \mathcal{H}$. By Lemma A.1, we know that $P_{\mathcal{H}}(M_1)^2 - P_{\mathcal{H}}(M_2)^2 \in \mathcal{H}$, so it further follows from Lemma A.2 that $P_{\mathcal{H}}(M_1)^2 \succeq P_{\mathcal{H}}(M_2)^2$. Since matrix square root is operator monotone, it holds that $P_{\mathcal{H}}(M_1) \succeq P_{\mathcal{H}}(M_2)$. We also know that $P_{\mathcal{H}}(M_1) - P_{\mathcal{H}}(M_2) \in \mathcal{K}$ because \mathcal{K} is a subalgebra. Then we conclude that $P_{\mathcal{H}}(M_1) - P_{\mathcal{H}}(M_2) \in \mathcal{H}$ from the definition of \mathcal{H} . This completes the proof. \Box

We can further extend the definition of $P_{\mathcal{H}}(\cdot)$ to all positive semi-definite matrices. Specifically, for any $M \succeq 0$, we have $M + \epsilon I_d \succ 0$ for any $\epsilon > 0$, so $P_{\mathcal{H}}(M + \epsilon I_d)$ is well-defined. Also, by the operator monotonicity of $P_{\mathcal{H}}(\cdot)$ from

Proposition A.3, $P_{\mathcal{H}}(M + \epsilon I_d) \preceq P_{\mathcal{H}}(M + \epsilon' I_d)$ for $\epsilon' \ge \epsilon > 0$. This implies that $P_{\mathcal{H}}(M + \epsilon I_d)$ has a limit as $\epsilon \to 0$. Therefore, for any $M \succeq 0$, we define

$$P_{\mathcal{H}}(\boldsymbol{M}) = \lim_{\epsilon \searrow 0} P_{\mathcal{H}}(\boldsymbol{M} + \epsilon \boldsymbol{I}_d).$$
(14)

Note that the above equality is also true for $M \succ 0$. To see this, we apply the optimality of every $P_{\mathcal{H}}(M + \epsilon I_d)$ to get

$$\langle \boldsymbol{M} + \epsilon \boldsymbol{I}_d, P_{\mathcal{H}}(\boldsymbol{M})^{-1} \rangle + \operatorname{Tr}(P_{\mathcal{H}}(\boldsymbol{M})) > \langle \boldsymbol{M} + \epsilon \boldsymbol{I}_d, P_{\mathcal{H}}(\boldsymbol{M} + \epsilon \boldsymbol{I}_d)^{-1} \rangle + \operatorname{Tr}(P_{\mathcal{H}}(\boldsymbol{M} + \epsilon \boldsymbol{I}_d)).$$

Also note that $P_{\mathcal{H}}(\boldsymbol{M} + \epsilon \boldsymbol{I}_d) \succ P_{\mathcal{H}}(\boldsymbol{M} - \delta \boldsymbol{I}_d) \succ 0$ for sufficiently small $\delta > 0$ such that $\boldsymbol{M} - \delta \boldsymbol{I}_d \succ 0$. Therefore, letting $\epsilon \rightarrow 0$ on both sides of the previous inequality, we can exchange the order of taking the limit and taking the inverse of $P_{\mathcal{H}}(\boldsymbol{M} + \epsilon \boldsymbol{I}_d)$ to obtain

$$\langle \boldsymbol{M}, P_{\mathcal{H}}(\boldsymbol{M})^{-1} \rangle + \operatorname{Tr}(P_{\mathcal{H}}(\boldsymbol{M})) \geq \left\langle \boldsymbol{M}, \left(\lim_{\epsilon \searrow 0} P_{\mathcal{H}}(\boldsymbol{M} + \epsilon \boldsymbol{I}_d)\right)^{-1} \right\rangle + \operatorname{Tr}\left(\lim_{\epsilon \searrow 0} P_{\mathcal{H}}(\boldsymbol{M} + \epsilon \boldsymbol{I}_d)\right).$$

Then by the optimality of $P_{\mathcal{H}}(M)$ and its uniqueness, we conclude that (14) is also valid for any $M \succ 0$.

Indeed, the definition of $P_{\mathcal{H}}(M)$ in (14) provides a continuous extension of $P_{\mathcal{H}}$ to \mathcal{S}^d_+ , as summarized in the following proposition.

Proposition A.4. Let \mathcal{H} be a well-structured preconditioner set under Definition 3.1. As a function on S_{++}^d , $P_{\mathcal{H}}$ can be continuously extended to be a function on S_{+}^d . Moreover, for any $\mathbf{M} \succeq 0$ such that $\mathbf{M} \neq 0$, $P_{\mathcal{H}}(\mathbf{M}) \succeq 0$ satisfies the following properties:

(a) $\operatorname{span}(M) \subseteq \operatorname{span}(P_{\mathcal{H}}(M))$ and $\langle M, P_{\mathcal{H}}(M)^{\dagger} \rangle + \operatorname{Tr}(P_{\mathcal{H}}(M)) = \inf_{H \in \mathcal{H}} \langle M, H^{-1} \rangle + \operatorname{Tr}(H).$

(b)
$$\langle \boldsymbol{M}, P_{\mathcal{H}}(\boldsymbol{M})^{\dagger} \rangle = \operatorname{Tr}(P_{\mathcal{H}}(\boldsymbol{M})).$$

- (c) $\langle \boldsymbol{M}, \overline{P_{\mathcal{H}}(\boldsymbol{M})}^{\dagger} \rangle = \inf_{\boldsymbol{H} \in \mathcal{H}, \operatorname{Tr}(\boldsymbol{H}) < 1} \langle \boldsymbol{M}, \boldsymbol{H}^{-1} \rangle$ where we recall that $\overline{P_{\mathcal{H}}(\boldsymbol{M})} = \operatorname{Tr}(P_{\mathcal{H}}(\boldsymbol{M}))^{-1} P_{\mathcal{H}}(\boldsymbol{M})$.
- (d) For any $\boldsymbol{H} \in \mathcal{H}$, $\langle P_{\mathcal{H}}(\boldsymbol{M})^{\dagger} \boldsymbol{M} P_{\mathcal{H}}(\boldsymbol{M})^{\dagger} \boldsymbol{\Pi}_{\boldsymbol{M}}, \boldsymbol{H} P_{\mathcal{H}}(\boldsymbol{M}) \rangle = 0$, where $\boldsymbol{\Pi}_{\boldsymbol{M}}$ is the projection matrix onto $\operatorname{span}(P_{\mathcal{H}}(\boldsymbol{M}))$.

Moreover, for any $M_1 \succeq M_2 \succeq 0$ *, it holds that* $P_{\mathcal{H}}(M_1) - P_{\mathcal{H}}(M_2) \in \mathcal{H}$ *, and in particular,* $P_{\mathcal{H}}(M_1) \succeq P_{\mathcal{H}}(M_2)$ *.*

Proof of Proposition A.4. We divide the proof into different parts for different properties of $P_{\mathcal{H}}$.

Proof for continuous extension of $P_{\mathcal{H}}$. We first show that $P_{\mathcal{H}}$ can be continuously extended to S_{+}^d , and we consider the extension of $P_{\mathcal{H}}$ as given in (14). We first show that for any $M \succeq 0$ and any sequence $\{M_n\}_{n=1}^{\infty}$ such that each $M_n \succ 0$ and $\lim_{n\to\infty} M_n = M$, it holds that $\lim_{n\to\infty} P_{\mathcal{H}}(M_n) = P_{\mathcal{H}}(M)$. Note that for any $\delta \in (0, 1)$, there exist $\bar{\epsilon}_n \ge \epsilon_n > 0$ such that $0 \prec (1-\delta)(M + \epsilon_n I_d) \preceq M_n \preceq M + \bar{\epsilon}_n I_d$ for all large enough n and moreover, $\lim_{n\to\infty} \bar{\epsilon}_n = \lim_{n\to\infty} \epsilon_n = 0$. Then by the operator monotonicity of $P_{\mathcal{H}}(\cdot)$, we have $P_{\mathcal{H}}((1-\delta)(M + \epsilon_n I_d)) \preceq P_{\mathcal{H}}(M_n) \preceq P_{\mathcal{H}}(M + \bar{\epsilon}_n I_d)$. Also note that $P_{\mathcal{H}}((1-\delta)(M + \epsilon_n I_d)) = \sqrt{1-\delta}P_{\mathcal{H}}(M + \epsilon_n I_d)$. Letting $n \to \infty$, both $P_{\mathcal{H}}(M + \epsilon_n I_d)$ and $P_{\mathcal{H}}(M + \bar{\epsilon}_n I_d)$ converge to $P_{\mathcal{H}}(M)$, which then implies that $\sqrt{1-\delta}P_{\mathcal{H}}(M) \preceq \lim_{n\to\infty} P_{\mathcal{H}}(M_n) \preceq P_{\mathcal{H}}(M)$. Since δ is arbitrary, we conclude that $\lim_{n\to\infty} P_{\mathcal{H}}(M_n) = P_{\mathcal{H}}(M)$. Next, consider any general sequence $\{M_n\}_{n=1}^{\infty}$ such that $M_n \succeq 0$ and $\lim_{n\to\infty} M_n = M$. For any $\epsilon > 0$, there exists $\delta > 0$ such that for all $M' \succ 0$ with $||M - M'||_F \leq \delta$, it holds that $||P_{\mathcal{H}}(M) - P_{\mathcal{H}}(M_n)| = |Im_{m\to\infty} P_{\mathcal{H}}(M_n + \frac{\delta}{2dm} I_d)$, where for every $m \ge 1$ we have $||M - (M_n + \frac{\delta}{2m\sqrt{d}} I_d)||_F \leq ||M - M_n||_F + ||\frac{\delta}{2m\sqrt{d}} I_d||_F \leq \delta$, so $||P_{\mathcal{H}}(M) - P_{\mathcal{H}}(M_n + \frac{\delta}{2m\sqrt{d}} I_d)||_F \leq \epsilon$ for all n > N. Therefore, it follows that $\lim_{n\to\infty} P_{\mathcal{H}}(M_n) = P_{\mathcal{H}}(M)$. In conclusion, $P_{\mathcal{H}}(\cdot)$ can be extended to be a continuous function on S_{+}^d .

Below, we fix any $M \succeq 0$ such that $M \neq 0$.

Proof for span(M) \subseteq span($P_{\mathcal{H}}(M)$). Let $\Pi_M := P_{\mathcal{H}}(M)P_{\mathcal{H}}(M)^{\dagger}$ be the projection matrix onto span($P_{\mathcal{H}}(M)$), then $\Pi_M \in \mathcal{H}$ by Lemma A.1. It suffices to show that $\langle I_d - \Pi_M, M \rangle = 0$. Using the fact that $\operatorname{Tr}(AB) \leq \operatorname{Tr}(A) \operatorname{Tr}(B)$ for any $A, B \succeq 0$, we can get the following inequality for any $\epsilon > 0$:

$$\langle \boldsymbol{I}_{d} - \boldsymbol{\Pi}_{\boldsymbol{M}}, \boldsymbol{M} \rangle = \operatorname{Tr}((\boldsymbol{I}_{d} - \boldsymbol{\Pi}_{\boldsymbol{M}})\boldsymbol{M}) \leq \operatorname{Tr}((\boldsymbol{I}_{d} - \boldsymbol{\Pi}_{\boldsymbol{M}})P_{\mathcal{H}}(\boldsymbol{M} + \epsilon \boldsymbol{I}_{d})) \cdot \operatorname{Tr}(P_{\mathcal{H}}(\boldsymbol{M} + \epsilon \boldsymbol{I}_{d})^{-1}\boldsymbol{M})$$
$$= \langle \boldsymbol{I}_{d} - \boldsymbol{\Pi}_{\boldsymbol{M}}, P_{\mathcal{H}}(\boldsymbol{M} + \epsilon \boldsymbol{I}_{d}) \rangle \cdot \langle \boldsymbol{M}, \boldsymbol{P}_{\mathcal{H}}(\boldsymbol{M} + \epsilon \boldsymbol{I}_{d})^{-1} \rangle$$
(15)

By Proposition A.3, we know that $\langle \boldsymbol{M}, P_{\mathcal{H}}(\boldsymbol{M} + \epsilon \boldsymbol{I}_d)^{-1} \rangle = \operatorname{Tr}(P_{\mathcal{H}}(\boldsymbol{M} + \epsilon \boldsymbol{I}_d)) - \epsilon \langle \boldsymbol{I}_d, P_{\mathcal{H}}(\boldsymbol{M} + \epsilon \boldsymbol{I}_d)^{-1} \rangle \leq \operatorname{Tr}(P_{\mathcal{H}}(\boldsymbol{M} + \epsilon \boldsymbol{I}_d))$, which is bounded by an absolute constant for all $\epsilon \in (0, 1)$. Therefore, letting $\epsilon \to 0$ on both sided of (15), since $\langle \boldsymbol{I}_d - \boldsymbol{\Pi}_M, P_{\mathcal{H}}(\boldsymbol{M} + \epsilon \boldsymbol{I}_d) \rangle \to \langle \boldsymbol{I}_d - \boldsymbol{\Pi}_M, P_{\mathcal{H}}(\boldsymbol{M}) \rangle = 0$ by the definition of $\boldsymbol{\Pi}_M$, we conclude that $\langle \boldsymbol{I}_d - \boldsymbol{\Pi}_M, \boldsymbol{M} \rangle = 0$. This implies that $\operatorname{span}(\boldsymbol{M}) \subseteq \operatorname{span}(P_{\mathcal{H}}(\boldsymbol{M}))$.

Proof for the optimality of $P_{\mathcal{H}}(M)$. We still use Π_M to denote the projection matrix onto span $(P_{\mathcal{H}}(M))$. For any $\epsilon > 0$, by the optimality of $P_{\mathcal{H}}(M + \epsilon I_d)$, it holds for any $H \in \mathcal{H}$ that

$$\langle \boldsymbol{M} + \epsilon \boldsymbol{I}_d, P_{\mathcal{H}}(\boldsymbol{M} + \epsilon \boldsymbol{I}_d)^{-1} \rangle + \operatorname{Tr}(P_{\mathcal{H}}(\boldsymbol{M} + \epsilon \boldsymbol{I}_d)) \leq \langle \boldsymbol{M} + \epsilon \boldsymbol{I}_d, \boldsymbol{H}^{-1} \rangle + \operatorname{Tr}(\boldsymbol{H}).$$
 (16)

Since span(M) \subseteq span($P_{\mathcal{H}}(M)$), we know that $\langle M, P_{\mathcal{H}}(M + \epsilon I_d)^{-1} \rangle = \langle M, \Pi_M P_{\mathcal{H}}(M + \epsilon I_d)^{-1} \Pi_M \rangle$. Then since $\lim_{\epsilon \to 0} P_{\mathcal{H}}(M + \epsilon I_d) = P_{\mathcal{H}}(M)$, it holds that $\Pi_M P_{\mathcal{H}}(M + \epsilon I_d)^{-1} \Pi_M = (\Pi_M P_{\mathcal{H}}(M + \epsilon I_d) \Pi_M)^{\dagger}$ for sufficiently small $\epsilon > 0$, and also that $P_{\mathcal{H}}(M)^{\dagger} = \lim_{\epsilon \to 0} (\Pi_M P_{\mathcal{H}}(M + \epsilon I_d) \Pi_M)^{\dagger}$. This implies

$$\lim_{\epsilon \searrow 0} \langle \boldsymbol{M}, P_{\mathcal{H}}(\boldsymbol{M} + \epsilon \boldsymbol{I}_d)^{-1} \rangle = \langle \boldsymbol{M}, P_{\mathcal{H}}(\boldsymbol{M})^{\dagger} \rangle.$$
(17)

Therefore, assuming the existence of the limit of $\langle \epsilon I_d, P_{\mathcal{H}}(M + \epsilon I_d)^{-1} \rangle$ as $\epsilon \to 0$, taking the limit of $\epsilon \to 0$ on both sides of (16) yields

$$\langle \boldsymbol{M}, P_{\mathcal{H}}(\boldsymbol{M})^{\dagger} \rangle + \operatorname{Tr}(P_{\mathcal{H}}(\boldsymbol{M})) + \lim_{\epsilon \to 0} \langle \epsilon \boldsymbol{I}_d, P_{\mathcal{H}}(\boldsymbol{M} + \epsilon \boldsymbol{I}_d)^{-1} \rangle \leq \langle \boldsymbol{M}, \boldsymbol{H}^{-1} \rangle + \operatorname{Tr}(\boldsymbol{H}).$$
 (18)

Hence, it suffices to show that $\langle \epsilon I_d, P_{\mathcal{H}}(M + \epsilon I_d)^{-1} \rangle \to 0$ as $\epsilon \to 0$. To prove this, note that for any $\epsilon > 0$, $P_{\mathcal{H}}(M) + \sqrt{\epsilon}(I_d - \Pi_M) \succ 0$, and its inverse is given by $P_{\mathcal{H}}(M)^{\dagger} + \epsilon^{-1/2}(I_d - \Pi_M)$. Also, $P_{\mathcal{H}}(M) + \sqrt{\epsilon}(I_d - \Pi_M) \in \mathcal{H}$ by Lemma A.1. Therefore, by the optimality of $P_{\mathcal{H}}(M + \epsilon I_d)$, we have

$$\begin{split} \langle \boldsymbol{M} + \epsilon \boldsymbol{I}_{d}, P_{\mathcal{H}}(\boldsymbol{M} + \epsilon \boldsymbol{I}_{d})^{-1} \rangle + \operatorname{Tr}(P_{\mathcal{H}}(\boldsymbol{M} + \epsilon \boldsymbol{I}_{d})) \\ &\leq \langle \boldsymbol{M} + \epsilon \boldsymbol{I}_{d}, (P_{\mathcal{H}}(\boldsymbol{M}) + \sqrt{\epsilon}(\boldsymbol{I}_{d} - \boldsymbol{\Pi}_{\boldsymbol{M}}))^{-1} \rangle + \operatorname{Tr}(P_{\mathcal{H}}(\boldsymbol{M}) + \sqrt{\epsilon}(\boldsymbol{I}_{d} - \boldsymbol{\Pi}_{\boldsymbol{M}})) \\ &= \langle \boldsymbol{M} + \epsilon \boldsymbol{I}_{d}, P_{\mathcal{H}}(\boldsymbol{M})^{\dagger} + \epsilon^{-1/2}(\boldsymbol{I}_{d} - \boldsymbol{\Pi}_{\boldsymbol{M}}) \rangle + \operatorname{Tr}(P_{\mathcal{H}}(\boldsymbol{M})) + \sqrt{\epsilon}(d - \operatorname{rank}(P_{\mathcal{H}}(\boldsymbol{M}))) \\ &= \langle \boldsymbol{M}, P_{\mathcal{H}}(\boldsymbol{M})^{\dagger} \rangle + \epsilon \cdot \operatorname{Tr}(P_{\mathcal{H}}(\boldsymbol{M})^{\dagger}) + \operatorname{Tr}(P_{\mathcal{H}}(\boldsymbol{M})) + 2\sqrt{\epsilon}(d - \operatorname{rank}(P_{\mathcal{H}}(\boldsymbol{M}))) \end{split}$$

where we apply $\text{Tr}(I_d - \Pi_M) = d - \text{rank}(P_H(M))$ and the last equality follows from the fact that $\langle M, I_d - \Pi_M \rangle = 0$. Rearranging the above inequality, we obtain

$$\begin{aligned} \epsilon \langle \boldsymbol{I}_d, P_{\mathcal{H}}(\boldsymbol{M} + \epsilon \boldsymbol{I}_d)^{-1} \rangle &\leq \langle \boldsymbol{M}, P_{\mathcal{H}}(\boldsymbol{M})^{\dagger} \rangle - \langle \boldsymbol{M}, P_{\mathcal{H}}(\boldsymbol{M} + \epsilon \boldsymbol{I}_d)^{-1} \rangle + \operatorname{Tr}(P_{\mathcal{H}}(\boldsymbol{M})) - \operatorname{Tr}(P_{\mathcal{H}}(\boldsymbol{M} + \epsilon \boldsymbol{I}_d)) \\ &+ \epsilon \cdot \operatorname{Tr}(P_{\mathcal{H}}(\boldsymbol{M})^{\dagger}) + 2\sqrt{\epsilon}(d - \operatorname{rank}(P_{\mathcal{H}}(\boldsymbol{M}))). \end{aligned}$$

Now applying (17) and noting that $\epsilon \langle I_d, P_{\mathcal{H}}(M + \epsilon I_d)^{-1} \rangle \geq 0$, taking the limit $\epsilon \to 0$, we obtain

$$\lim_{\epsilon \to 0} \langle \epsilon \mathbf{I}_d, P_{\mathcal{H}}(\mathbf{M} + \epsilon \mathbf{I}_d)^{-1} \rangle = 0.$$
⁽¹⁹⁾

Then combining (18) and (19), we conclude that for any $H \in \mathcal{H}$,

$$\langle \boldsymbol{M}, P_{\mathcal{H}}(\boldsymbol{M})^{\dagger} \rangle + \operatorname{Tr}(P_{\mathcal{H}}(\boldsymbol{M})) \leq \langle \boldsymbol{M}, \boldsymbol{H}^{-1} \rangle + \operatorname{Tr}(\boldsymbol{H}).$$

This confirms the optimality of $P_{\mathcal{H}}(\boldsymbol{M})$. Further applying the same argument as in the proof of Proposition A.3 yields the optimality of $\overline{P_{\mathcal{H}}(\boldsymbol{M})} = \text{Tr}(P_{\mathcal{H}}(\boldsymbol{M}))^{-1}P_{\mathcal{H}}(\boldsymbol{M})$.

Proof for property (d) of $P_{\mathcal{H}}(M)$. By the optimality of $P_{\mathcal{H}}(M)$, we have

$$\langle \boldsymbol{M}, P_{\mathcal{H}}(\boldsymbol{M})^{\dagger} \rangle + \operatorname{Tr}(P_{\mathcal{H}}(\boldsymbol{M})) = \inf_{\boldsymbol{H} \in \mathcal{H}} \underbrace{\langle \boldsymbol{M}, \boldsymbol{H}^{-1} \rangle + \operatorname{Tr}(\boldsymbol{H})}_{f_{\boldsymbol{M}}(\boldsymbol{H})}$$

$$= \inf_{\boldsymbol{H} \in \mathcal{H}} \underbrace{\langle \boldsymbol{M}, \boldsymbol{H}^{-1} \rangle + \operatorname{Tr}(\boldsymbol{\Pi}_{\boldsymbol{M}} \boldsymbol{H} \boldsymbol{\Pi}_{\boldsymbol{M}})}_{\tilde{f}_{\boldsymbol{M}}(\boldsymbol{H})}.$$

$$(20)$$

To see why the last equality holds, first note that $\inf_{H \in \mathcal{H}} f_M(H) \ge \inf_{H \in \mathcal{H}} f_M(H)$ because $\operatorname{Tr}(H) \ge \operatorname{Tr}(\Pi_M H \Pi_M)$. Then for the other direction, note that we can always approximate $\tilde{f}_M(H)$ using $f_M(H_\delta)$ where $H_\delta = (1-\delta)\Pi_M H \Pi_M^\top + \delta I_d$, for which $f_M(H_\delta) \le \frac{1}{1-\delta} \langle M, H^{-1} \rangle + (1-\delta) \operatorname{Tr}(\Pi_M H \Pi_M) + \delta d \to \tilde{f}_M(H)$ as $\delta \to 0$. This implies that $\inf_{H \in \mathcal{H}} \tilde{f}_M(H) \ge \inf_{H \in \mathcal{H}} f_M(H)$, and thus the optimal objective values of the two optimization problems are the same. Now define $\widetilde{H}_M = P_{\mathcal{H}}(M) + (I_d - \Pi_M) \in \mathcal{H}$ whose inverse is $\widetilde{H}_M^{-1} = P_{\mathcal{H}}(M)^\dagger + (I_d - \Pi_M)$, then $\widetilde{H}_M \succ 0$ is a solution to the optimization problem in (20) because $\tilde{f}_M(\widetilde{H}_M) = \langle M, P_{\mathcal{H}}(M)^\dagger \rangle + \operatorname{Tr}(P_{\mathcal{H}}(M)) = \inf_{H \in \mathcal{H}} f_M(H)$. The gradient of \tilde{f}_M is given by $\nabla \tilde{f}_M(H) = -H^{-1}MH^{-1} + \Pi_M$. Since \mathcal{H} is a cone, the optimality of \widetilde{H}_M implies that for any $H \in \mathcal{H}$,

$$0 = \langle \nabla f(\widetilde{H}_{M}), H - \widetilde{H}_{M} \rangle = \langle -\widetilde{H}_{M}^{-1}M\widetilde{H}_{M}^{-1} + \Pi_{M}, H - \widetilde{H}_{M} \rangle.$$

By the definition of \widetilde{H}_M , we have $\widetilde{H}_M^{-1}M\widetilde{H}_M^{-1} = P_{\mathcal{H}}(M)^{\dagger}MP_{\mathcal{H}}(M)^{\dagger}$ because $\operatorname{span}(M) \subseteq \operatorname{span}(P_{\mathcal{H}}(M))$. Therefore, it follows that

$$0 = \langle P_{\mathcal{H}}(\boldsymbol{M})^{\dagger} \boldsymbol{M} P_{\mathcal{H}}(\boldsymbol{M})^{\dagger} - \boldsymbol{\Pi}_{\boldsymbol{M}}, \boldsymbol{H} - P_{\mathcal{H}}(\boldsymbol{M}) - (\boldsymbol{I}_{d} - \boldsymbol{\Pi}_{\boldsymbol{M}}) \rangle$$
$$= \langle P_{\mathcal{H}}(\boldsymbol{M})^{\dagger} \boldsymbol{M} P_{\mathcal{H}}(\boldsymbol{M})^{\dagger} - \boldsymbol{\Pi}_{\boldsymbol{M}}, \boldsymbol{H} - P_{\mathcal{H}}(\boldsymbol{M}) \rangle.$$

Proof for the operator monotonicity of the extended $P_{\mathcal{H}}$. For any $M_1 \succeq M_2 \succeq 0$, it follows from Proposition A.3 that for any $\epsilon > 0$, $P_{\mathcal{H}}(M_1 + \epsilon I_d) \succeq P_{\mathcal{H}}(M_2 + \epsilon I_d)$. Taking the limit as $\epsilon \to 0$ on both sides yields $P_{\mathcal{H}}(M_1) \succeq P_{\mathcal{H}}(M_2)$. Similarly, as $P_{\mathcal{H}}(M_1 + \epsilon I_d) - P_{\mathcal{H}}(M_2 + \epsilon I_d) \in \mathcal{H}$ for every $\epsilon > 0$ by Proposition A.3, we have $P_{\mathcal{H}}(M_1) - P_{\mathcal{H}}(M_2) = \lim_{\epsilon \to 0} P_{\mathcal{H}}(M_1 + \epsilon I_d) - P_{\mathcal{H}}(M_2 + \epsilon) \in \mathcal{H}$ because \mathcal{H} is a closed set. This completes the proof. \Box

A.3. Adaptive gradient norm corresponds to the dual norm

Proof of Lemma 3.3. We begin with the case of t = 1. Fix any $\boldsymbol{g} \in \mathbb{R}^d$. First for any $\boldsymbol{w} \in \mathbb{R}^d$ with $\|\boldsymbol{w}\|_{\mathcal{H}} \leq 1$, by Cauchy-Schwarz inequality, it holds for all $\boldsymbol{H} \in \mathcal{H}$ with $\operatorname{Tr}(\boldsymbol{H}) \leq 1$ that

$$oldsymbol{g}^{ op}oldsymbol{w} \leq \sqrt{oldsymbol{g}^{ op}oldsymbol{H}^{-1}oldsymbol{g}} \cdot oldsymbol{\|w\|}_{\mathcal{H}} \leq \sqrt{\langleoldsymbol{g}oldsymbol{g}^{ op},oldsymbol{H}^{-1}
angle} \cdot oldsymbol{\|w\|}_{\mathcal{H}} \leq \sqrt{\langleoldsymbol{g}oldsymbol{g}^{ op},oldsymbol{H}^{-1}
angle}$$

where the second inequality follows from the definition of $\|w\|_{\mathcal{H}}$. Now, further taking infimum over $H \in \mathcal{H}$ with $\operatorname{Tr}(H) \leq 1$ and then taking supremum over $w \in \mathbb{R}^d$ with $\|w\|_{\mathcal{H}} \leq 1$, we obtain that

$$\|\boldsymbol{g}\|_{\mathcal{H}}^* \leq \sqrt{\inf_{\boldsymbol{H}\in\mathcal{H},\mathrm{Tr}(\boldsymbol{H})\leq 1} \langle \boldsymbol{g}\boldsymbol{g}^{\top},\boldsymbol{H}^{-1}\rangle}.$$
(21)

Next, we need to show that the above inequality holds also for the other direction.

By Proposition A.4, we define $\boldsymbol{H}_* = \operatorname{Tr}(P_{\mathcal{H}}(\boldsymbol{g}\boldsymbol{g}^{\top}))^{-1}P_{\mathcal{H}}(\boldsymbol{g}\boldsymbol{g}^{\top}) = \inf_{\boldsymbol{H} \in \mathcal{H}, \operatorname{Tr}(\boldsymbol{H}) \leq 1} \langle \boldsymbol{M}, \boldsymbol{H}^{-1} \rangle$. We choose correspondingly $\boldsymbol{w}_* = \frac{\boldsymbol{H}_*^{\dagger}\boldsymbol{g}}{\|\boldsymbol{H}_*^{\dagger}\boldsymbol{g}\|_{\mathcal{H}}}$, which satisfies that $\|\boldsymbol{w}_*\|_{\mathcal{H}} = 1$. Then

$$\|\boldsymbol{g}\|_{\mathcal{H}}^* \ge \boldsymbol{g}^\top \boldsymbol{w}_* = \frac{\boldsymbol{g}^\top \boldsymbol{H}_*^\dagger \boldsymbol{g}}{\|\boldsymbol{H}_*^\dagger \boldsymbol{g}\|_{\mathcal{H}}}.$$
(22)

Therefore, it suffices to show that $\boldsymbol{g}^{\top}\boldsymbol{w}_{*}\geq\sqrt{\boldsymbol{g}^{\top}\boldsymbol{H}_{*}^{\dagger}\boldsymbol{g}}$, which is equivalent to

$$oldsymbol{g}^ opoldsymbol{H}^\dagger_*oldsymbol{g} \geq egin{array}{c} oldsymbol{H}^\dagger_*oldsymbol{g} & \|oldsymbol{H}^\dagger_*oldsymbol{g}\|_{\mathcal{H}}^2 = \sup_{oldsymbol{H}\in\mathcal{H}, \mathrm{Tr}(oldsymbol{H}) \leq 1}oldsymbol{g}^ opoldsymbol{H}^\dagger_*oldsymbol{H}oldsymbol{H}^\dagger_*oldsymbol{g}.$$

By the property (d) of $P_{\mathcal{H}}(\boldsymbol{g}\boldsymbol{g}^{\top})$ from Proposition A.4, we have $\langle P_{\mathcal{H}}(\boldsymbol{g}\boldsymbol{g}^{\top})^{\dagger}\boldsymbol{g}\boldsymbol{g}^{\top}P_{\mathcal{H}}(\boldsymbol{g}\boldsymbol{g}^{\top})^{\dagger} - \boldsymbol{\Pi}_{\boldsymbol{g}\boldsymbol{g}^{\top}}, \boldsymbol{H} - P_{\mathcal{H}}(\boldsymbol{g}\boldsymbol{g}^{\top}) \rangle = 0$ for any $\boldsymbol{H} \in \mathcal{H}$. Rearranging this equality, we obtain

$$\boldsymbol{g}^{\top} P_{\mathcal{H}} (\boldsymbol{g} \boldsymbol{g}^{\top})^{\dagger} \boldsymbol{H} P_{\mathcal{H}} (\boldsymbol{g} \boldsymbol{g}^{\top})^{\dagger} \boldsymbol{g} = \boldsymbol{g}^{\top} P_{\mathcal{H}} (\boldsymbol{g} \boldsymbol{g}^{\top})^{\dagger} \boldsymbol{g} + \langle \boldsymbol{\Pi}_{\boldsymbol{g} \boldsymbol{g}^{\top}}, \boldsymbol{H} - P_{\mathcal{H}} (\boldsymbol{g} \boldsymbol{g}^{\top}) \rangle$$

= $\boldsymbol{g}^{\top} P_{\mathcal{H}} (\boldsymbol{g} \boldsymbol{g}^{\top})^{\dagger} \boldsymbol{g} + \langle \boldsymbol{\Pi}_{\boldsymbol{g} \boldsymbol{g}^{\top}}, \boldsymbol{H} \rangle - \operatorname{Tr} (P_{\mathcal{H}} (\boldsymbol{g} \boldsymbol{g}^{\top}))$

where the second equality is because $\Pi_{gg^{\top}}$ is exactly the projection matrix onto span $(P_{\mathcal{H}}(gg^{\top}))$. Applying the above equality with $\operatorname{Tr}(P_{\mathcal{H}}(gg^{\top})) \cdot \overline{H}$ in place of H and plugging in the definition of H_* , we further have

$$\boldsymbol{g}^{\top} \boldsymbol{H}_{*}^{\dagger} \overline{\boldsymbol{H}} \boldsymbol{H}_{*}^{\dagger} \boldsymbol{g} = \boldsymbol{g}^{\top} \boldsymbol{H}_{*}^{\dagger} \boldsymbol{g} + \operatorname{Tr}(P_{\mathcal{H}}(\boldsymbol{g}\boldsymbol{g}^{\top})) \langle \boldsymbol{\Pi}_{\boldsymbol{g}\boldsymbol{g}^{\top}}, \overline{\boldsymbol{H}} \rangle - \operatorname{Tr}(P_{\mathcal{H}}(\boldsymbol{g}\boldsymbol{g}^{\top}))$$

$$\leq \boldsymbol{g}^{\top} \boldsymbol{H}_{*}^{\dagger} \boldsymbol{g} + \operatorname{Tr}(P_{\mathcal{H}}(\boldsymbol{g}\boldsymbol{g}^{\top}))^{2} \operatorname{Tr}(\overline{\boldsymbol{H}}) - \operatorname{Tr}(P_{\mathcal{H}}(\boldsymbol{g}\boldsymbol{g}^{\top}))^{2}$$

$$= \boldsymbol{g}^{\top} \boldsymbol{H}_{*}^{\dagger} \boldsymbol{g}.$$

This implies that for any $H \in \mathcal{H}$ with $\operatorname{Tr}(H) \leq 1$, we have $g^{\top}H_*^{\dagger}HH_*^{\dagger}g \leq g^{\top}H_*^{\dagger}g$. Therefore, it follows from (22) that

$$\|\boldsymbol{g}\|_{\mathcal{H}}^{*} \geq \sqrt{\boldsymbol{g}^{\top} \boldsymbol{H}_{*}^{\dagger} \boldsymbol{g}} \sqrt{\frac{\boldsymbol{g}^{\top} \boldsymbol{H}_{*}^{\dagger} \boldsymbol{g}}{\inf_{\boldsymbol{H} \in \mathcal{H}, \operatorname{Tr}(\boldsymbol{H}) \leq 1} \boldsymbol{g}^{\top} \boldsymbol{H}_{*}^{\dagger} \boldsymbol{H} \boldsymbol{H}_{*}^{\dagger} \boldsymbol{g}}} \\ \geq \sqrt{\boldsymbol{g}^{\top} \boldsymbol{H}_{*}^{\dagger} \boldsymbol{g}} = \sqrt{\frac{\inf_{\boldsymbol{H} \in \mathcal{H}, \operatorname{Tr}(\boldsymbol{H}) \leq 1} \langle \boldsymbol{g} \boldsymbol{g}^{\top}, \boldsymbol{H}^{-1} \rangle}$$
(23)

where the equality follows from the definition of H_* .

Finally, combining (21) and (23), we conclude that

$$\|\boldsymbol{g}\|_{\mathcal{H}}^* = \inf_{\boldsymbol{H}\in\mathcal{H}, \operatorname{Tr}(\boldsymbol{H})\leq 1} \sqrt{\langle \boldsymbol{g}\boldsymbol{g}^{\top}, \boldsymbol{H}^{-1} \rangle}.$$
(24)

For the general case of $\boldsymbol{g}_{1:t} = (\boldsymbol{g}_1, \dots, \boldsymbol{g}_t) \in \mathbb{R}^{d \times t}$, note that for any $\boldsymbol{H} \in \mathcal{H}$

$$\left\langle \sum_{s=1}^{t} \boldsymbol{g}_{s} \boldsymbol{g}_{s}^{\top}, \boldsymbol{H}^{-1} \right\rangle = \langle \overline{\operatorname{vec}}(\boldsymbol{g}_{1:t}) \overline{\operatorname{vec}}(\boldsymbol{g}_{1:t})^{\top}, \boldsymbol{H}^{-1} \otimes \boldsymbol{I}_{t} \rangle = \langle \overline{\operatorname{vec}}(\boldsymbol{g}_{1:t}) \overline{\operatorname{vec}}(\boldsymbol{g}_{1:t})^{\top}, (\boldsymbol{H} \otimes \boldsymbol{I}_{t})^{-1} \rangle.$$

Then applying (24) with $\overline{\text{vec}}(g_{1:t})$ in place of g and $\mathcal{H} \otimes I_t$ in place of \mathcal{H} , we obtain

$$\begin{split} \|\overline{\operatorname{vec}}(\boldsymbol{g}_{1:t})\|_{\mathcal{H}\otimes\boldsymbol{I}_{t}}^{*} &= \inf_{\boldsymbol{H}\otimes\boldsymbol{I}_{t}\in\mathcal{H}\otimes\boldsymbol{I}_{t},\operatorname{Tr}(\boldsymbol{H}\otimes\boldsymbol{I}_{t})\leq1} \sqrt{\langle\overline{\operatorname{vec}}(\boldsymbol{g}_{1:t})\overline{\operatorname{vec}}(\boldsymbol{g}_{1:t})^{\top},(\boldsymbol{H}\otimes\boldsymbol{I}_{t})^{-1}\rangle} \\ &= \inf_{\boldsymbol{H}\in\mathcal{H},\operatorname{Tr}(\boldsymbol{H})\leq1} \sqrt{t\langle\overline{\operatorname{vec}}(\boldsymbol{g}_{1:t})\overline{\operatorname{vec}}(\boldsymbol{g}_{1:t})^{\top},(\boldsymbol{H}\otimes\boldsymbol{I}_{t})^{-1}\rangle} \\ &= \sqrt{t}\cdot\inf_{\boldsymbol{H}\in\mathcal{H},\operatorname{Tr}(\boldsymbol{H})\leq1} \sqrt{\left\langle\sum_{s=1}^{t}\boldsymbol{g}_{s}\boldsymbol{g}_{s}^{\top},\boldsymbol{H}^{-1}\right\rangle} \end{split}$$

where the second equality is because $Tr(H \otimes I_t) = Tr(H) Tr(I_t) = t \cdot Tr(H)$. This completes the proof.

A.4. Examples of ill-structured preconditioner sets

A.4.1. PRECONDITIONER SET OF TWO-SIDED SHAMPOO

Consider $\mathcal{H} = \{ \boldsymbol{H} \in \mathcal{S}^d_+ : \boldsymbol{H} = \boldsymbol{U} \otimes \boldsymbol{V} \text{ for some } \boldsymbol{U} \in \mathbb{R}^{d_L \times d_L}, \boldsymbol{V} \in \mathbb{R}^{d_R \times d_R} \}$. In particular, we consider the special case of $d_L = d_R = 2$, and the following two matrices

$$G_1(\epsilon) = \operatorname{diag}(1, \epsilon, \epsilon, \epsilon), \qquad G_2(\epsilon) = \operatorname{diag}(1, \epsilon, \epsilon, 1)$$

where $\epsilon > 0$ is a small constant to be determined later. Clearly $G_1(\epsilon) \preceq G_2(\epsilon)$ when $\epsilon \leq 1$, but we will show that $H_{G_1(\epsilon)} \preceq H_{G_2(\epsilon)}$ does not hold for sufficiently small $\epsilon > 0$. For any $H = U \otimes V \in \mathcal{H}$, it holds that both U and V are PSD, and we explicitly parametrize them as

$$\boldsymbol{U} = \begin{pmatrix} u_1 & u_2 \\ u_2 & u_3 \end{pmatrix}, \qquad \boldsymbol{V} = \begin{pmatrix} v_1 & v_2 \\ v_2 & v_3 \end{pmatrix}$$

where $u_1, u_3 \ge 0, u_1u_3 \ge u_2^2$, and similarly, $v_1, v_3 \ge 0, v_1v_3 \ge v_3^2$. Correspondingly,

$$oldsymbol{H} = oldsymbol{U} \otimes oldsymbol{V} = egin{pmatrix} u_1 v_1 & u_1 v_2 & u_2 v_1 & u_2 v_2 \ u_1 v_2 & u_1 v_3 & u_2 v_2 & u_2 v_3 \ u_2 v_1 & u_2 v_2 & u_3 v_1 & u_3 v_2 \ u_2 v_2 & u_3 v_2 & u_3 v_2 & u_3 v_3 \end{pmatrix}.$$

We first analyze $H_{G_2(\epsilon)}$. For convenience, we consider

$$\begin{aligned} \boldsymbol{H}_{\boldsymbol{G}_{2}(\epsilon)}^{-1} &= \operatorname*{arg\,min}_{\boldsymbol{H}\in\mathcal{H}} \langle \boldsymbol{G}_{2}(\epsilon), \boldsymbol{H} \rangle + \operatorname{Tr}(\boldsymbol{H}^{-1}) \\ &= \operatorname*{arg\,min}_{\boldsymbol{U},\boldsymbol{V}\in\mathcal{S}_{+}^{2}} \underbrace{\langle \boldsymbol{G}_{2}(\epsilon), \boldsymbol{U}\otimes\boldsymbol{V} \rangle + \operatorname{Tr}((\boldsymbol{U}\otimes\boldsymbol{V})^{-1})}_{=:f_{\boldsymbol{G}_{2}(\epsilon)}(\boldsymbol{U},\boldsymbol{V})} \end{aligned}$$

Recall the properties of the Kronecker product that $(U \otimes V)^{-1} = U^{-1} \otimes V^{-1}$ and that $\text{Tr}(U^{-1} \otimes V^{-1}) = \text{Tr}(U^{-1}) \cdot \text{Tr}(V^{-1})$. Then further using the definition of $G_2(\epsilon)$, we have

$$f_{\mathbf{G}_{2}(\epsilon)}(\mathbf{U},\mathbf{V}) = u_{1}v_{1} + u_{3}v_{3} + \operatorname{Tr}(\mathbf{U}^{-1}) \cdot \operatorname{Tr}(\mathbf{V}^{-1}) + \epsilon(u_{1}v_{3} + u_{3}v_{1})$$

$$= u_{1}v_{1} + u_{3}v_{3} + \frac{u_{1} + u_{3}}{u_{1}u_{3} - u_{2}^{2}} \cdot \frac{v_{1} + v_{3}}{v_{1}v_{3} - v_{2}^{2}} + \epsilon(u_{1}v_{3} + u_{3}v_{1})$$

where we apply the explicit expression of U^{-1} and V^{-1} . First observe that to minimize $f_{G_2(\epsilon)}(U, V)$, we must have $u_2 = v_2 = 0$ because otherwise $\frac{u_1+u_3}{u_1u_3-u_2^2} > \frac{u_1+u_3}{u_1u_3}$ and similarly $\frac{v_1+v_3}{v_1v_3-v_2^2} > \frac{v_1+v_3}{v_1v_3}$. Therefore, it suffices to further minimize the following:

$$f_{G_2(\epsilon)}(\boldsymbol{U}, \boldsymbol{V}) = \underbrace{u_1 v_1 + u_3 v_3 + \frac{u_1 + u_3}{u_1 u_3} \cdot \frac{v_1 + v_3}{v_1 v_3}}_{\tilde{f}_{G_2}(\boldsymbol{U}, \boldsymbol{V})} + \epsilon(u_1 v_3 + u_3 v_1).$$

For the first term $\tilde{f}_{G_2}(\boldsymbol{U}, \boldsymbol{V})$, we can apply the AM-GM inequality to get $\tilde{f}_{G_2}(\boldsymbol{U}, \boldsymbol{V}) \ge (4\sqrt{u_1u_3v_1v_3})^{1/3}$, where the equality is achieved when $u_1 = u_3$, $v_1 = v_3$, and $u_1v_1 = \sqrt{2}$. Moreover, for sufficiently small $\epsilon > 0$, the contribution of the second term $\epsilon(u_1v_3 + u_3v_1)$ is negligible, and thus we conclude that $\lim_{\epsilon \to 0} \boldsymbol{H}_{G_2(\epsilon)} = \frac{\sqrt{2}}{2} \boldsymbol{I}_4$.

Similarly, for $G_1(\epsilon)$, we have

$$\begin{split} f_{\boldsymbol{G}_{1}(\epsilon)}(\boldsymbol{U},\boldsymbol{V}) &= u_{1}v_{1} + \frac{u_{1} + u_{3}}{u_{1}u_{3} - u_{2}^{2}} \cdot \frac{v_{1} + v_{3}}{v_{1}v_{3} - v_{2}^{2}} + \epsilon(u_{1}v_{3} + u_{3}v_{1} + u_{3}v_{3}) \\ &\geq u_{1}v_{1} + \frac{u_{1} + u_{3}}{u_{1}u_{3}} \cdot \frac{v_{1} + v_{3}}{v_{1}v_{3}} + \epsilon(u_{1}v_{3} + u_{3}v_{1} + u_{3}v_{3}) \\ &= \underbrace{u_{1}v_{1} + \frac{1}{u_{1}v_{1}}}_{\tilde{f}_{\boldsymbol{G}_{1}}(\boldsymbol{U},\boldsymbol{V})} + \frac{1}{u_{1}v_{3}} + \frac{1}{u_{3}v_{1}} + \frac{1}{u_{3}v_{3}} + \epsilon(u_{1}v_{3} + u_{3}v_{1} + u_{3}v_{3}) \end{split}$$

where the equality holds when $u_2 = v_2 = 0$. The first term $\hat{f}_{G_1}(U, V)$ attains the minimum value when $u_1v_1 = 1$, and the remainder can be made small by choosing large u_3, v_3 when ϵ is sufficiently small. Therefore, we conclude that $\lim_{\epsilon \to 0} H_{G_1(\epsilon)} = \operatorname{diag}(1, 0, 0, 0)$.

Now comparing the limits of $H_{G_1}(\epsilon)$ and $H_{G_2}(\epsilon)$ as $\epsilon \to 0$, we can see that for sufficiently small $\epsilon > 0$, it holds that $H_{G_1(\epsilon)}[1,1] > H_{G_2(\epsilon)}[1,1]$. Hence, for sufficiently small $\epsilon > 0$, $H_{G_1}(\epsilon) \preceq H_{G_2(\epsilon)}$ does not hold.

A.4.2. TRIDIAGONAL MATRICES

Consider $\mathcal{H} = \{ \mathbf{H} \in S^d_+ : \mathbf{H} \text{ is tridiagonal} \}$. Here, for \mathbf{H} to be tridiagonal, it has nonzero elements on the main diagonal, the first diagonal above the main diagonal, and the first diagonal below the main diagonal. We consider the specific example of 3-by-3 tridiagonal matrices, and examine $P_{\mathcal{H}}(\cdot)$ for the following two matrices

$$\boldsymbol{M} = \begin{pmatrix} 2 & 1 & 1 \\ 1 & 2 & 1 \\ 1 & 1 & 2 \end{pmatrix}, \quad \boldsymbol{M}' = \begin{pmatrix} 10000 & 1 & 1 \\ 1 & 2 & 1 \\ 1 & 1 & 2 \end{pmatrix}.$$

It is clear that $0 \prec M \preceq M'$. We numerically solve for $P_{\mathcal{H}}(M)$ and $P_{\mathcal{H}}(M')$ to get

$$P_{\mathcal{H}}(\boldsymbol{M}) \approx \begin{pmatrix} 1.382548 & 0.297594 & 0\\ 0.297594 & 1.318491 & 0.297594\\ 0 & 0.297594 & 1.382548 \end{pmatrix}, \qquad P_{\mathcal{H}}(\boldsymbol{M}') \approx \begin{pmatrix} 100.000004 & 0.007229 & 0\\ 0.007229 & 1.365999 & 0.366002\\ 0 & 0.366002 & 1.366032 \end{pmatrix}.$$

Note that $P_{\mathcal{H}}(M) \preceq P_{\mathcal{H}}(M')$ does not hold because the last diagonal entry of $P_{\mathcal{H}}(M)$ is larger than that of $P_{\mathcal{H}}(M')$.

B. Calculations for Examples of Well-Structured Preconditioner Sets

As mentioned in Section 3.3, we provide calculations for how to derive each specific algorithm from Algorithm 1 with specific choice of \mathcal{H} and explain each entry in Table 1. Recall from Algorithm 1 that

$$oldsymbol{M}_t = \epsilon oldsymbol{I}_d + \sum_{s=1}^t oldsymbol{g}_s oldsymbol{g}_s^ op$$

B.1. AdaGrad-Norm

For AdaGrad-Norm, we have $\mathcal{H} = \{ c \cdot \mathbf{I}_d \mid c \geq 0 \}.$

Calculation for H_t . Then for any $H = c \cdot I_d \in \mathcal{H}$,

$$\langle \boldsymbol{M}_t, \boldsymbol{H}^{-1} \rangle + \eta^2 \operatorname{Tr}(\boldsymbol{H}) = \frac{1}{c} \operatorname{Tr}(\boldsymbol{M}_t) + \eta^2 c d \ge 2\eta \sqrt{d \cdot \operatorname{Tr}(\boldsymbol{M}_t)}$$

where the equality is achieved by choosing

$$c = \frac{1}{\eta} \sqrt{\frac{1}{d} \operatorname{Tr}(\boldsymbol{M}_t)} = \frac{1}{\eta} \sqrt{\epsilon + \frac{1}{d} \sum_{s=1}^t \|\boldsymbol{g}_s\|_2^2}.$$

This corresponds exactly to the update rule of AdaGrad-Norm, which adjusts the global learning rate based on the accumulated ℓ_2 norm of past gradients.

Calculation for $\|\cdot\|_{\mathcal{H}}$. For the associated $\|\cdot\|_{\mathcal{H}}$, we have

$$\|\boldsymbol{x}\|_{\mathcal{H}} = \sup_{0 \le c \le 1/d} \sqrt{c \cdot \boldsymbol{x}^{\top} \boldsymbol{I}_d \boldsymbol{x}} = \frac{\|\boldsymbol{x}\|_2}{\sqrt{d}}.$$

Calculation for $||| g_{1:t} |||_{\mathcal{H}}$. We can calculate the adaptive gradient norm similarly:

$$\|\|\boldsymbol{g}_{1:t}\|\|_{\mathcal{H}} = \inf_{0 \le c \le 1/d} \sqrt{\left\langle \sum_{s=1}^{t} \boldsymbol{g}_{s} \boldsymbol{g}_{s}^{\top}, (c\boldsymbol{I}_{d})^{-1} \right\rangle} = \sqrt{d} \sqrt{\operatorname{Tr}\left(\sum_{s=1}^{t} \boldsymbol{g}_{s} \boldsymbol{g}_{s}^{\top}\right)} = \sqrt{d} \sqrt{\sum_{s=1}^{t} \|\boldsymbol{g}_{s}\|_{2}^{2}}$$

B.2. Diagonal AdaGrad

For diagonal AdaGrad, $\mathcal{H} = \mathcal{D}^d_+ = \{ \operatorname{diag}(c_1, \ldots, c_d) \mid c_1, \ldots, c_d \geq 0 \}.$

Calculation for H_t . For any $H = \text{diag}(c_1, \ldots, c_d) \in \mathcal{H}$, we have

$$\langle \boldsymbol{M}_t, \boldsymbol{H}^{-1} \rangle + \eta^2 \operatorname{Tr}(\boldsymbol{H}) = \sum_{i=1}^d \left(\frac{\boldsymbol{M}_{t,ii}}{c_i} + \eta^2 c_i \right) \ge \sum_{i=1}^d 2\eta \sqrt{\boldsymbol{M}_{t,ii}}$$

where the equality is achieved by choosing

$$c_i = \frac{1}{\eta} \sqrt{M_{t,ii}} = \frac{1}{\eta} \sqrt{\epsilon + \sum_{s=1}^t g_{s,i}^2}, \quad \text{for } i = 1, \dots, d.$$

This corresponds to the update rule diagonal AdaGrad, which computes the historical sum of squared gradients for each individual coordinate.

Calculation for $\|\cdot\|_{\mathcal{H}}$. For the associated $\|\cdot\|_{\mathcal{H}}$, we have

$$\|\boldsymbol{x}\|_{\mathcal{H}} = \sup_{\sum_{i=1}^{d} c_i \leq 1} \sqrt{\boldsymbol{x}^{\top} \operatorname{diag}(c_1, \dots, c_d) \boldsymbol{x}} = \sup_{\sum_{i=1}^{d} c_i \leq 1} \sqrt{\sum_{i=1}^{d} c_i x_i^2} = \max_{i \in [d]} |x_i| = \|\boldsymbol{x}\|_{\infty}$$

Calculation for $|||g_{1:t}|||_{\mathcal{H}}$. We can calculate the adaptive gradient norm similarly:

$$\|\|\boldsymbol{g}_{1:t}\|\|_{\mathcal{H}} = \inf_{\sum_{i=1}^{d} c_i \le 1} \sqrt{\left\langle \sum_{s=1}^{t} \boldsymbol{g}_s \boldsymbol{g}_s^{\top}, \operatorname{diag}(c_1 \dots, c_d)^{-1} \right\rangle} \\ = \inf_{\sum_{i=1}^{d} c_i \le 1} \sqrt{\sum_{i=1}^{d} \frac{1}{c_i} \sum_{s=1}^{t} g_{s,i}^2} \\ = \sum_{i=1}^{d} \sqrt{\sum_{s=1}^{t} g_{s,i}^2}$$

where the last equality is because of the Cauchy inequality: for $c_1, \ldots, c_d \ge 0$ such that $\sum_{i=1}^d c_i \le 1$,

$$\sqrt{\sum_{i=1}^{d} \frac{1}{c_i} \sum_{s=1}^{t} g_{s,i}^2} \ge \sqrt{\sum_{i=1}^{d} \frac{1}{c_i} \sum_{s=1}^{t} g_{s,i}^2} \sqrt{\sum_{i=1}^{d} c_i} \ge \sum_{i=1}^{d} \sqrt{\sum_{s=1}^{t} g_{s,i}^2}.$$

B.3. Full-matrix AdaGrad

For full-matrix AdaGrad, $\mathcal{H} = \mathcal{S}_{+}^{d}$.

Calculation for H_t . Note that $\langle M_t, H^{-1} \rangle + \eta^2 \operatorname{Tr}(H)$ is a convex function of $H \succ 0$, so we can get H_t by first calculating the gradient of the objective function and then setting it to zero. Specifically, for any $H \succ 0$, we have

$$\nabla_{\boldsymbol{H}} \left(\left\langle \boldsymbol{M}_t, \boldsymbol{H}^{-1} \right\rangle + \eta^2 \operatorname{Tr}(\boldsymbol{H}) \right) = -\boldsymbol{H}^{-1} \boldsymbol{M}_t \boldsymbol{H}^{-1} + \eta^2 \boldsymbol{I}_d.$$

Setting it to zero yields

$$oldsymbol{H}_t = rac{1}{\eta} oldsymbol{M}_t^{rac{1}{2}} = rac{1}{\eta} igg(\epsilon oldsymbol{I}_d + \sum_{s=1}^t oldsymbol{g}_s oldsymbol{g}_s^{ op} igg)^{rac{1}{2}}.$$

This corresponds to the update rule of full-matrix AdaGrad.

Calculation for $\|\cdot\|_{\mathcal{H}}$. For the associated $\|\cdot\|_{\mathcal{H}}$, we have

$$\begin{split} \|\boldsymbol{x}\|_{\mathcal{H}} &= \sup_{\boldsymbol{H} \succeq 0, \mathrm{Tr}(\boldsymbol{H}) \leq 1} \sqrt{\boldsymbol{x}^{\top} \boldsymbol{H} \boldsymbol{x}} \\ &= \sup_{\boldsymbol{H} \succeq 0, \mathrm{Tr}(\boldsymbol{H}) \leq 1} \sqrt{\langle \boldsymbol{x} \boldsymbol{x}^{\top}, \boldsymbol{H} \rangle} = \sqrt{\lambda_1(\boldsymbol{x} \boldsymbol{x}^{\top})} = \|\boldsymbol{x}\|_2 \end{split}$$

Calculation for $|||g_{1:t}|||_{\mathcal{H}}$. We can adapt the calculation for H_t to the constrained optimization problem in the definition of $|||g_{1:t}|||_{\mathcal{H}}$ to get

$$\begin{split} \|\|\boldsymbol{g}_{1:t}\|\|_{\mathcal{H}} &= \inf_{\boldsymbol{H}\in\mathcal{H}, \operatorname{Tr}(\boldsymbol{H})\leq 1} \sqrt{\left\langle \sum_{s=1}^{t} \boldsymbol{g}_{s}\boldsymbol{g}_{s}^{\top}, \boldsymbol{H}^{-1} \right\rangle} \\ &= \sqrt{\left\langle \left\langle \sum_{s=1}^{t} \boldsymbol{g}_{s}\boldsymbol{g}_{s}^{\top}, \left(\frac{1}{\operatorname{Tr}[(\sum_{s=1}^{t} \boldsymbol{g}_{s}\boldsymbol{g}_{s}^{\top})^{\frac{1}{2}}] \left(\sum_{s=1}^{t} \boldsymbol{g}_{s}\boldsymbol{g}_{s}^{\top} \right)^{\frac{1}{2}} \right)^{-1} \right\rangle} \\ &= \operatorname{Tr}\left[\left(\sum_{s=1}^{t} \boldsymbol{g}_{s}\boldsymbol{g}_{s}^{\top} \right)^{\frac{1}{2}} \right]. \end{split}$$

B.4. One-sided Shampoo

For one-sided Shampoo, $\mathcal{H} = \mathcal{S}^{d_L}_+ \otimes I_{d_R}$.

Calculation for H_t . For any $H = H_L \otimes I_{d_R} \in \mathcal{H}, H^{-1} = H_L^{-1} \otimes I_{d_R}$, so we have $\operatorname{Tr}(H) = d_R \operatorname{Tr}(H_L)$ and $\langle M_t, H^{-1} \rangle = \langle M_t, (H_L \otimes I_{d_R})^{-1} \rangle$

$$= \langle \epsilon \boldsymbol{I}_{d}, \boldsymbol{H}_{L}^{-1} \otimes \boldsymbol{I}_{d_{R}} \rangle + \sum_{s=1}^{t} \langle \overline{\operatorname{vec}}(\boldsymbol{G}_{s}) \overline{\operatorname{vec}}(\boldsymbol{G}_{s})^{\top}, \boldsymbol{H}_{L}^{-1} \otimes \boldsymbol{I}_{d_{R}} \rangle$$
$$= \epsilon d_{R} \operatorname{Tr}(\boldsymbol{H}_{L}^{-1}) + \sum_{s=1}^{t} \operatorname{Tr}(\boldsymbol{G}_{s}^{\top} \boldsymbol{H}_{L}^{-1} \boldsymbol{G}_{s})$$
$$= \operatorname{Tr}\left[\left(\epsilon d_{R} \boldsymbol{I}_{d_{L}} + \sum_{s=1}^{t} \boldsymbol{G}_{s} \boldsymbol{G}_{s}^{\top} \right) \boldsymbol{H}_{L}^{-1} \right]$$

where we use the fact that $\langle \overline{\operatorname{vec}}(\boldsymbol{X}) \overline{\operatorname{vec}}(\boldsymbol{X})^{\top}, \boldsymbol{H}_L \otimes \boldsymbol{I}_{d_R} \rangle = \langle \boldsymbol{X} \boldsymbol{X}^{\top}, \boldsymbol{H}_L \rangle$ for any $\boldsymbol{X} \in \mathbb{R}^{d_L \times d_R}$. Again, since the objective function is convex in $\boldsymbol{H}_L \succ 0$, we can derive \boldsymbol{H}_L by first calculating the gradient and then setting it to zero. Taking derivative with respect to \boldsymbol{H}_L , we obtain

$$\nabla_{\boldsymbol{H}_{L}}\left(\langle \boldsymbol{M}_{t}, \boldsymbol{H}^{-1}\rangle + \eta^{2}\operatorname{Tr}(\boldsymbol{H})\right) = -\boldsymbol{H}_{L}^{-1}\left(\epsilon d_{R}\boldsymbol{I}_{d_{L}} + \sum_{s=1}^{t}\boldsymbol{G}_{s}\boldsymbol{G}_{s}^{\top}\right)\boldsymbol{H}_{L}^{-1} + \eta^{2}d_{R}\boldsymbol{I}_{d_{L}}.$$

Setting it to 0, we obtain

$$\begin{split} \arg\min_{\boldsymbol{H}_{L}\in\mathcal{S}_{+}^{d_{L}}} \langle \boldsymbol{M}_{t}, (\boldsymbol{H}_{L}\otimes\boldsymbol{I}_{d_{R}})^{-1} \rangle + \eta^{2}d_{R}\operatorname{Tr}(\boldsymbol{H}_{L}\otimes\boldsymbol{I}_{d_{R}}) &= \frac{1}{\eta\sqrt{d_{R}}} \bigg(\epsilon d_{R}\boldsymbol{I}_{d_{L}} + \sum_{s=1}^{t}\boldsymbol{G}_{s}\boldsymbol{G}_{s}^{\top} \bigg)^{\frac{1}{2}} \\ &= \frac{1}{\eta} \bigg(\epsilon \boldsymbol{I}_{d_{L}} + \frac{1}{d_{R}}\sum_{s=1}^{t}\boldsymbol{G}_{s}\boldsymbol{G}_{s}^{\top} \bigg)^{\frac{1}{2}}. \end{split}$$

This is exactly the L_t in Algorithm 2. Therefore, the preconditioner H_t in one-sided Shampoo is given by

$$\boldsymbol{H}_{t} = \frac{1}{\eta} \left(\epsilon \boldsymbol{I}_{d_{L}} + \frac{1}{d_{R}} \sum_{s=1}^{t} \boldsymbol{G}_{s} \boldsymbol{G}_{s}^{\mathsf{T}} \right)^{\frac{1}{2}} \otimes \boldsymbol{I}_{d_{R}}.$$

As a result, Algorithm 1 with $\mathcal{H} = S^{d_L}_+ \otimes I_{d_R}$ recovers Algorithm 2.

Calculation for $\|\cdot\|_{\mathcal{H}}$. For the associated $\|\cdot\|_{\mathcal{H}}$, we have

$$\begin{split} \|\boldsymbol{x}\|_{\mathcal{H}} &= \sup_{\boldsymbol{H} = \boldsymbol{H}_L \otimes \boldsymbol{I}_{d_R}, \boldsymbol{H}_L \succeq 0, \operatorname{Tr}(\boldsymbol{H}) \leq 1} \sqrt{\boldsymbol{x}^\top \boldsymbol{H} \boldsymbol{x}} \\ &= \sup_{\boldsymbol{H}_L \succeq 0, \operatorname{Tr}(\boldsymbol{H}_L) \leq 1/d_R} \sqrt{\langle \operatorname{vec} \left(\boldsymbol{X} \right) \overline{\operatorname{vec}} \left(\boldsymbol{X} \right)^\top, \boldsymbol{H}_L \otimes \boldsymbol{I}_{d_R} \rangle} \\ &= \sup_{\boldsymbol{H}_L \succeq 0, \operatorname{Tr}(\boldsymbol{H}_L) \leq 1/d_R} \sqrt{\langle \boldsymbol{X} \boldsymbol{X}^\top, \boldsymbol{H}_L \rangle} \\ &= \frac{1}{\sqrt{d_R}} \|\boldsymbol{X}\|_{\operatorname{op}} \end{split}$$

where the last equality is achieved at $H_L = uu^{\top}/d_R$ for u being the leading eigenvector of XX^{\top} . The derivation above provides a proof for Lemma 4.1.

Calculation for $|||g_{1:t}|||_{\mathcal{H}}$. Again, for the adaptive gradient norm, we can adapt the calculation for H_t to the constrained optimization problem in the definition of $|||g_{1:t}|||_{\mathcal{H}}$ to get

$$\begin{split} \|\|\boldsymbol{g}_{1:t}\|\|_{\mathcal{H}} &= \inf_{\boldsymbol{H}\in\mathcal{H},\operatorname{Tr}(\boldsymbol{H})\leq 1} \sqrt{\left\langle \sum_{s=1}^{t} \boldsymbol{g}_{s}\boldsymbol{g}_{s}^{\top}, \boldsymbol{H}^{-1} \right\rangle} \\ &= \inf_{\boldsymbol{H}_{L}\in\mathcal{S}_{+}^{d_{L}},\operatorname{Tr}(\boldsymbol{H}_{L})\leq \frac{1}{d_{R}}} \sqrt{\left\langle \sum_{s=1}^{t} \boldsymbol{G}_{s}\boldsymbol{G}_{s}^{\top}, \boldsymbol{H}_{L}^{-1} \right\rangle} \\ &= \sqrt{\left\langle \left\langle \sum_{s=1}^{t} \boldsymbol{G}_{s}\boldsymbol{G}_{s}^{\top}, \left[\frac{1}{d_{R}\operatorname{Tr}[(\sum_{s=1}^{t} \boldsymbol{G}_{s}\boldsymbol{G}_{s}^{\top})^{\frac{1}{2}}] \left(\sum_{s=1}^{t} \boldsymbol{G}_{s}\boldsymbol{G}_{s}^{\top} \right)^{\frac{1}{2}} \right]^{-1} \right\rangle} \\ &= \sqrt{d_{R}}\operatorname{Tr}\left[\left(\sum_{s=1}^{t} \boldsymbol{G}_{s}\boldsymbol{G}_{s}^{\top} \right)^{\frac{1}{2}} \right] = \operatorname{Tr}\left[\left(d_{R}\sum_{s=1}^{t} \boldsymbol{G}_{s}\boldsymbol{G}_{s}^{\top} \right)^{\frac{1}{2}} \right]. \end{split}$$

C. Analysis for EMA adaptive optimization

As mentioned in Section 3, our analysis can be generalized to Algorithm 4, which is an EMA variant of Algorithm 1. To simplify the theoretical analysis, we instead focus on Algorithm 5 in which the current gradient is not scaled by $1 - \beta_2$ when computing M_t . It serves as a bridge between Algorithms 1 and 4. When $\beta_2 = 1$, Algorithm 5 exactly recovers Algorithm 1. When the learning rate $\tilde{\eta}$ is rescaled by $\sqrt{1 - \beta_2}$ and ϵ is rescaled by $1/(1 - \beta_2)$, Algorithm 5 is equivalent to Algorithm 4. Theorem C.1 extends Theorem 3.4 to get online regret bound for Algorithm 5. The proof is in Appendix D.3.

Theorem C.1. Let \mathcal{H} be a well-structured preconditioner set under Definition 3.1. Then for any convex loss functions $\tilde{L}_1, \ldots, \tilde{L}_T$, the regret of Algorithm 5 compared to any $\mathbf{x}^* \in \mathcal{X}$ can be bounded as

$$\sum_{t=1}^{T} \beta_2^{\frac{T-t}{2}} \left[\tilde{L}_t(\boldsymbol{x}_t) - \tilde{L}_t(\boldsymbol{x}^*) \right] \le \left(\frac{D^2}{2\tilde{\eta}} + \tilde{\eta} \right) \inf_{\boldsymbol{H} \in \mathcal{H}, \operatorname{Tr}(\boldsymbol{H}) \le 1} \sqrt{\left\langle \tilde{\boldsymbol{M}}_T, \boldsymbol{H}^{-1} \right\rangle}$$

where $D = \max_{t \in [T]} \| \boldsymbol{x}_t - \boldsymbol{x}^* \|_{\mathcal{H}}$.

Theorem C.2 provides the convergence rate for Algorithm 5, which generalizes Theorem 3.8. The proof is in Appendix D.3. **Theorem C.2.** Let \mathcal{H} be a well-structured preconditioner set under Definition 3.1. Consider any independent stochastic convex loss functions $\tilde{L}_1, \ldots, \tilde{L}_T$ satisfying Assumption 3.7, and let $H(\tilde{L}, \mathcal{H})$ be the \mathcal{H} -smoothness of their expectation \tilde{L} . Suppose the global minimizer of L, denoted by \mathbf{x}^* , is in \mathcal{X} . Then for the iterates $\mathbf{x}_1, \ldots, \mathbf{x}_T$ of Algorithm 5, denoting $\bar{\mathbf{x}}_{1:T} = (\sum_{t=1}^T \beta_2^{\frac{T-t}{2}})^{-1} \sum_{t=1}^T \beta_2^{\frac{T-t}{2}} \mathbf{x}_t$, it holds that

$$\mathbb{E}\tilde{L}(\bar{\boldsymbol{x}}_{1:T}) - \tilde{L}_{t}(\boldsymbol{x}^{*}) \leq \frac{16}{\sum_{t=1}^{T} \beta_{2}^{\frac{T-t}{2}}} \|\mathcal{X}\|_{\mathcal{H}}^{2} H(\tilde{L},\mathcal{H}) + \frac{4\sqrt{2}}{\sqrt{\sum_{t=1}^{T} \beta_{2}^{T-t}}} \|\mathcal{X}\|_{\mathcal{H}} \sigma + \frac{4\sqrt{2}}{\sum_{t=1}^{T} \beta_{2}^{\frac{-t}{2}}} \|\mathcal{X}\|_{\mathcal{H}} d\sqrt{\epsilon}$$

Algorithm 4 EMA Adaptive Regularization Meta- Algorithm	Algorithm 5 Weighted Adaptive Regularization Meta- Algorithm
Hyperparam: $\epsilon > 0$, convex set $\mathcal{X} \subseteq \mathbb{R}^d$, learning rate η , preconditioners $\mathcal{H} \subset \mathcal{S}^d_+, \beta_2 \in (0, 1)$ Input: initialization x_1 , loss functions $\{L_t\}_{t=1}^T : \mathbb{R}^d \to \mathbb{R}$ $M_0 \leftarrow \epsilon L_d$	Hyperparam: $\epsilon > 0$, convex set $\mathcal{X} \subseteq \mathbb{R}^d$, learning rate $\tilde{\eta}$, preconditioners $\mathcal{H} \subset \mathcal{S}^d_+, \beta_2 \in (0, 1)$ Input: initialization \boldsymbol{x}_1 , loss functions $\{\tilde{L}_t\}_{t=1}^T : \mathbb{R}^d \to \mathbb{R}$ $\tilde{M}_t \leftarrow \epsilon L$.
for $t = 1, 2, \dots, T$ do $g_t \leftarrow \nabla L_t(x_t)$	$\mathbf{for} \ t = 1, 2, \dots, T \ \mathbf{do}$ $\tilde{\mathbf{a}}_t \leftarrow \nabla \tilde{L}_t(\tilde{\mathbf{x}}_t)$
$ \begin{split} \mathbf{M}_t &\leftarrow \beta_2 \mathbf{M}_{t-1} + (1-\beta_2) \mathbf{g}_t \mathbf{g}_t^\top \\ \mathbf{H}_t &\leftarrow \arg\min_{\mathbf{H} \in \mathcal{H}} \left\langle \mathbf{M}_t, \mathbf{H}^{-1} \right\rangle + \eta^2 \operatorname{Tr}(\mathbf{H}) \\ \mathbf{H}_t &\leftarrow \mathbf{M}_t \left(\mathbf{H}_t^{-1} \right) \end{split} $	$egin{aligned} & \tilde{m{M}}_t \leftarrow eta_2 m{M}_{t-1} + ilde{m{g}}_t m{ ilde{g}}_t^{ op} \ & ilde{m{H}}_t \leftarrow rg\min_{m{ ilde{H}} \in \mathcal{H}} \left\langle m{M}_t, m{m{H}}^{-1} ight angle + ilde{\eta}^2 \operatorname{Tr}(m{ ilde{H}}) \end{aligned}$
$egin{aligned} oldsymbol{x}_{t+1} \leftarrow \Pi_{\mathcal{X}^{+}}^{\mathcal{X}^{+}} oldsymbol{\left(x_{t} - H_{t}^{-1} oldsymbol{g}_{t} ight)} \ extbf{Return } oldsymbol{x}_{1}, \dots, oldsymbol{x}_{T} \end{aligned}$	$ ilde{m{x}}_{t+1} \leftarrow \Pi_{\mathcal{X}}^{ ilde{m{H}}_t} \left(ilde{m{x}}_t - ilde{m{H}}_t^{-1} ilde{m{g}}_t ight)^{-1}$ Return $ ilde{m{x}}_1, \dots, ilde{m{x}}_T$

where $\sigma = \inf_{\boldsymbol{H} \in \mathcal{H}, \operatorname{Tr}(\boldsymbol{H}) \leq 1} \sqrt{\langle \boldsymbol{\Sigma}, \boldsymbol{H}^{-1} \rangle}.$

When $\beta_2 = 1$, Theorem C.2 recovers the result in Theorem 3.8 and provides a $O(T^{-\frac{1}{2}})$ convergence rate. When $\beta_2 < 1$, the optimality gap is upper bounded by $O((1 - \beta_2) \|\mathcal{X}\|_{\mathcal{H}}^2 H(\tilde{L}, \mathcal{H}) + \sqrt{1 - \beta_2} \|\mathcal{X}\|_{\mathcal{H}} \sigma)$ when $T = \Theta\left(\frac{1}{1 - \beta_2}\right)$.

D. Proof for the Unified Analysis

D.1. Regret bound

We first present the proof for the main result on the regret bound for Algorithm 1 with a well-structured preconditioner set \mathcal{H} .

Theorem 3.4. Let \mathcal{H} be a well-structured preconditioner set under Definition 3.1. Then for any convex loss functions L_1, \ldots, L_T , the regret of Algorithm 1 compared to any $x^* \in \mathcal{X}$ can be bounded as

$$\sum_{t=1}^{T} L_t(\boldsymbol{x}_t) - \sum_{t=1}^{T} L_t(\boldsymbol{x}^*) \le \left(\frac{D^2}{2\eta} + \eta\right) \left(G + d\sqrt{\epsilon}\right)$$

where $G = ||| g_{1:T} |||_{\mathcal{H}}, D = \max_{t \in [T]} || x_t - x^* ||_{\mathcal{H}}.$

Proof of Theorem 3.4. First we will analyze the property of each H_t . Recall from Proposition 3.2 that H_t satisfies $\langle M_t, H_t^{-1} \rangle = \eta^2 \operatorname{Tr}(H_t)$ and that $\overline{H}_t := \operatorname{Tr}(H_t)^{-1}H_t = \arg \min_{H \in \mathcal{H}, \operatorname{Tr}(H) \leq 1} \langle M_t, H^{-1} \rangle$. Therefore,

$$\langle \boldsymbol{M}_t, \boldsymbol{H}_t^{-1} \rangle = \eta^2 \operatorname{Tr}(\boldsymbol{H}_t) = \eta \sqrt{\langle \boldsymbol{M}_t, \boldsymbol{H}_t^{-1} \rangle \operatorname{Tr}(\boldsymbol{H}_t)} = \eta \sqrt{\langle \boldsymbol{M}_t, \overline{\boldsymbol{H}}_t^{-1} \rangle}.$$
 (25)

Now recall the regret bound from Theorem 2.1:

$$\sum_{t=1}^{T} L_{t}(\boldsymbol{x}_{t}) - \sum_{t=1}^{T} L_{t}(\boldsymbol{x}^{*}) \leq \frac{1}{2} \left(\langle \boldsymbol{M}_{T}, \boldsymbol{H}_{T}^{-1} \rangle + \eta^{2} \operatorname{Tr}(\boldsymbol{H}_{T}) - \eta^{2} \operatorname{Tr}(\boldsymbol{H}_{0}) \right) + \frac{1}{2} \sum_{t=1}^{T} \left(\|\boldsymbol{x}_{t} - \boldsymbol{x}^{*}\|_{\boldsymbol{H}_{t}}^{2} - \|\boldsymbol{x}_{t+1} - \boldsymbol{x}^{*}\|_{\boldsymbol{H}_{t}}^{2} \right) = \frac{1}{2} \left(2\eta \sqrt{\langle \boldsymbol{M}_{T}, \overline{\boldsymbol{H}}_{T}^{-1} \rangle} - \eta^{2} \operatorname{Tr}(\boldsymbol{H}_{0}) \right) + \frac{1}{2} \sum_{t=1}^{T} \left(\|\boldsymbol{x}_{t} - \boldsymbol{x}^{*}\|_{\boldsymbol{H}_{t}}^{2} - \|\boldsymbol{x}_{t+1} - \boldsymbol{x}^{*}\|_{\boldsymbol{H}_{t}}^{2} \right)$$
(26)

where the equality follows from the facts in (25). Next, for the second term on the right-hand side of (26), we rearrange the summation to obtain

$$\begin{split} \sum_{t=1}^{T} \left(\| \boldsymbol{x}_t - \boldsymbol{x}^* \|_{\boldsymbol{H}_t}^2 - \| \boldsymbol{x}_{t+1} - \boldsymbol{x}^* \|_{\boldsymbol{H}_t}^2 \right) &\leq \| \boldsymbol{x}_1 - \boldsymbol{x}^* \|_{\boldsymbol{H}_1}^2 + \sum_{t=2}^{T} \left(\| \boldsymbol{x}_t - \boldsymbol{x}^* \|_{\boldsymbol{H}_t}^2 - \| \boldsymbol{x}_t - \boldsymbol{x}^* \|_{\boldsymbol{H}_{t-1}}^2 \right) \\ &= \| \boldsymbol{x}_1 - \boldsymbol{x}^* \|_{\boldsymbol{H}_1}^2 + \sum_{t=2}^{T} \| \boldsymbol{x}_t - \boldsymbol{x}^* \|_{\boldsymbol{H}_t - \boldsymbol{H}_{t-1}}^2 \,. \end{split}$$

Notice that $M_t - M_{t-1} = g_t g_t^\top \succeq 0$, and thus $H_t - H_{t-1} \in \mathcal{H}$ by Proposition 3.2. This implies

$$\|\boldsymbol{x}_{t} - \boldsymbol{x}^{*}\|_{\boldsymbol{H}_{t} - \boldsymbol{H}_{t-1}}^{2} = (\boldsymbol{x}_{t} - \boldsymbol{x}^{*})^{\top} (\boldsymbol{H}_{t} - \boldsymbol{H}_{t-1}) (\boldsymbol{x}_{t} - \boldsymbol{x}^{*}) \leq \operatorname{Tr}(\boldsymbol{H}_{t} - \boldsymbol{H}_{t-1}) \|\boldsymbol{x}_{t} - \boldsymbol{x}^{*}\|_{\mathcal{H}}^{2}$$

where the inequality follows from the definition of $\|\cdot\|_{\mathcal{H}}$ in (1). It then follows that

$$\begin{split} \sum_{t=1}^{T} \left(\| \boldsymbol{x}_{t} - \boldsymbol{x}^{*} \|_{\boldsymbol{H}_{t}}^{2} - \| \boldsymbol{x}_{t+1} - \boldsymbol{x}^{*} \|_{\boldsymbol{H}_{t}}^{2} \right) &\leq \operatorname{Tr}(\boldsymbol{H}_{1}) \| \boldsymbol{x}_{1} - \boldsymbol{x}^{*} \|_{\mathcal{H}}^{2} + \sum_{t=2}^{T} \operatorname{Tr}(\boldsymbol{H}_{t} - \boldsymbol{H}_{t-1}) \| \boldsymbol{x}_{t} - \boldsymbol{x}^{*} \|_{\mathcal{H}}^{2} \\ &\leq \operatorname{Tr}(\boldsymbol{H}_{T}) \max_{1 \leq t \leq T} \| \boldsymbol{x}_{t} - \boldsymbol{x}^{*} \|_{\mathcal{H}}^{2} \\ &= \frac{1}{\eta} \sqrt{\langle \boldsymbol{M}_{T}, \overline{\boldsymbol{H}}_{T}^{-1} \rangle} \max_{1 \leq t \leq T} \| \boldsymbol{x}_{t} - \boldsymbol{x}^{*} \|_{\mathcal{H}}^{2} \end{split}$$

where the equality again follows from (25). Plugging this back into (26), we obtain

$$\sum_{t=1}^{T} L_t(\boldsymbol{x}_t) - \sum_{t=1}^{T} L_t(\boldsymbol{x}^*) \le \eta \sqrt{\langle \boldsymbol{M}_T, \overline{\boldsymbol{H}}_T^{-1} \rangle} + \frac{1}{2\eta} \sqrt{\langle \boldsymbol{M}_T, \overline{\boldsymbol{H}}_T^{-1} \rangle} \max_{1 \le t \le T} \|\boldsymbol{x}_t - \boldsymbol{x}^*\|_{\mathcal{H}}^2$$
$$= \left(\frac{\max_{1 \le t \le T} \|\boldsymbol{x}_t - \boldsymbol{x}^*\|_{\mathcal{H}}^2}{2\eta} + \eta\right) \inf_{\boldsymbol{H} \in \mathcal{H}, \operatorname{Tr}(\boldsymbol{H}) \le 1} \sqrt{\langle \boldsymbol{M}_T, \boldsymbol{H}^{-1} \rangle}$$
$$\le \left(\frac{\max_{1 \le t \le T} \|\boldsymbol{x}_t - \boldsymbol{x}^*\|_{\mathcal{H}}^2}{2\eta} + \eta\right) \sqrt{\langle \boldsymbol{M}_T - \boldsymbol{M}_0, \boldsymbol{H}^{-1} \rangle + \langle \boldsymbol{M}_0, \boldsymbol{H}^{-1} \rangle}$$

for any $\boldsymbol{H} \in \mathcal{H}$ with $\operatorname{Tr}(\boldsymbol{H}) = 1$. In particular, we choose $\boldsymbol{H} = \alpha \boldsymbol{H}_T^* + \frac{1-\alpha}{d} \boldsymbol{I}_d$ for some $\alpha \in (0,1)$ where $\boldsymbol{H}_T^* = \arg \min_{\boldsymbol{H} \in \mathcal{H}, \operatorname{Tr}(\boldsymbol{H}) \leq 1} \sqrt{\langle \boldsymbol{M}_T - \boldsymbol{M}_0, \boldsymbol{H}^{-1} \rangle}$, so $\langle \boldsymbol{M}_T - \boldsymbol{M}_0, \boldsymbol{H}^{-1} \rangle = \| \boldsymbol{g}_{1:T} \|_{\mathcal{H}}^2$ according to the definition of the adaptive gradient norm in (7). Since $\boldsymbol{H}^{-1} \leq \frac{1}{\alpha} (\boldsymbol{H}_T^*)^{-1}$ and $\boldsymbol{H}^{-1} \leq \frac{d}{1-\alpha} \boldsymbol{I}_d$, we further have

$$\sum_{t=1}^{T} L_t(\boldsymbol{x}_t) - \sum_{t=1}^{T} L_t(\boldsymbol{x}^*) \le \left(\frac{\max_{1 \le t \le T} \|\boldsymbol{x}_t - \boldsymbol{x}^*\|_{\mathcal{H}}^2}{2\eta} + \eta\right) \sqrt{\frac{1}{\alpha} \langle \boldsymbol{M}_T - \boldsymbol{M}_0, (\boldsymbol{H}_T^*)^{-1} \rangle} + \frac{d}{1 - \alpha} \langle \boldsymbol{M}_0, \boldsymbol{I}_d \rangle$$
$$= \left(\frac{\max_{1 \le t \le T} \|\boldsymbol{x}_t - \boldsymbol{x}^*\|_{\mathcal{H}}^2}{2\eta} + \eta\right) \sqrt{\frac{1}{\alpha} \|\boldsymbol{g}_{1:T}\|_{\mathcal{H}}^2 + \frac{d^2\epsilon}{1 - \alpha}}.$$

Finally choosing $\alpha = \frac{\|\boldsymbol{g}_{1:T}\|_{\mathcal{H}}}{\|\|\boldsymbol{g}_{1:T}\|\|_{\mathcal{H}} + d\sqrt{\epsilon}}$, we obtain

$$\sum_{t=1}^{T} L_t(\boldsymbol{x}_t) - \sum_{t=1}^{T} L_t(\boldsymbol{x}^*) = \left(\frac{\max_{1 \le t \le T} \|\boldsymbol{x}_t - \boldsymbol{x}^*\|_{\mathcal{H}}^2}{2\eta} + \eta\right) \left(\|\|\boldsymbol{g}_{1:T}\|\|_{\mathcal{H}} + d\sqrt{\epsilon}\right)$$
proof.

This completes the proof.

1

The proof for Corollary 3.5 is straightforward.

Corollary D.1. Under the setting of Theorem 2.1, further suppose that \mathcal{X} is a bounded set in \mathbb{R}^d . Then choosing $\eta = \sqrt{2} \|\mathcal{X}\|_{\mathcal{H}}$, the regret bound for Algorithm 1 becomes

$$\sum_{t=1}^{T} L_t(\boldsymbol{x}_t) - \sum_{t=1}^{T} L_t(\boldsymbol{x}^*) \le 2\sqrt{2} \|\mathcal{X}\|_{\mathcal{H}} \left(G + d\sqrt{\epsilon}\right)$$

Proof of Corollary 3.5. Note that $D = \max_{t \in [T]} \| \boldsymbol{x}_t - \boldsymbol{x}^* \|_{\mathcal{H}} \leq \max_{t \in [T]} \| \boldsymbol{x}_t \|_{\mathcal{H}} + \| \boldsymbol{x}^* \|_{\mathcal{H}} \leq 2 \| \mathcal{X} \|_{\mathcal{H}}$ because $\boldsymbol{x}_1, \ldots, \boldsymbol{x}_T$ and \boldsymbol{x}^* are all in \mathcal{X} . The proof is completed by setting $\eta = \sqrt{2} \| \mathcal{X} \|_{\mathcal{H}}$ to minimize $2 \frac{\| \mathcal{X} \|_{\mathcal{H}}^2}{\eta} + \eta$.

D.2. Convergence rate

Next, we present the proof for the convergence rate of Algorithm 1 with a well-structured preconditioner set \mathcal{H} .

Theorem 3.8. Let \mathcal{H} be a well-structured preconditioner set under Definition 3.1. Consider any independent stochastic convex loss functions L_1, \ldots, L_T satisfying Assumption 3.7, and let $H(L, \mathcal{H})$ be the \mathcal{H} -smoothness of their expectation L. Suppose the global minimizer of L, denoted by \mathbf{x}^* , is in \mathcal{X} . Then for the iterates $\mathbf{x}_1, \ldots, \mathbf{x}_T$ of Algorithm 1, denoting $\bar{\mathbf{x}}_{1:T} = \frac{1}{T} \sum_{t=1}^T \mathbf{x}_t$, it holds that

$$\begin{split} & \mathbb{E}\left[L(\bar{\boldsymbol{x}}_{1:T}) - L(\boldsymbol{x}^*)\right] \\ \leq & \frac{16}{T} \left\|\mathcal{X}\right\|_{\mathcal{H}}^2 H(L,\mathcal{H}) + \frac{4\sqrt{2}}{\sqrt{T}} \left\|\mathcal{X}\right\|_{\mathcal{H}} \sigma + \frac{4\sqrt{2}d\sqrt{\epsilon}}{T} \left\|\mathcal{X}\right\|_{\mathcal{H}} \end{split}$$

where $\sigma = \inf_{\boldsymbol{H} \in \mathcal{H}, \operatorname{Tr}(\boldsymbol{H}) \leq 1} \sqrt{\langle \boldsymbol{\Sigma}, \boldsymbol{H}^{-1} \rangle}.$

Proof of Theorem 3.8. Let $H^* \in \mathcal{H}$ be the matrix given in Definition 3.6 for the loss function L, such that $\nabla^2 L(\mathbf{x}) \preceq H^*$. Therefore, for any $\mathbf{x}, \mathbf{x}' \in \mathcal{X}$,

$$L(\mathbf{x}') \leq L(\mathbf{x}) + \nabla L(\mathbf{x})^{\top} (\mathbf{x}' - \mathbf{x}) + \frac{1}{2} (\mathbf{x}' - \mathbf{x})^{\top} \mathbf{H}^* (\mathbf{x}' - \mathbf{x})$$

= $L(\mathbf{x}) + \frac{1}{2} (\mathbf{x}' - \mathbf{x} + (\mathbf{H}^*)^{-1} \nabla L(\mathbf{x}))^{\top} \mathbf{H}^* (\mathbf{x}' - \mathbf{x} + (\mathbf{H}^*)^{-1} \nabla L(\mathbf{x})) - \frac{1}{2} \nabla L(\mathbf{x})^{\top} (\mathbf{H}^*)^{-1} \nabla L(\mathbf{x}).$

Then we have that

$$\begin{split} L(\boldsymbol{x}^*) &= \min_{\boldsymbol{x}'} L(\boldsymbol{x}') \\ &\leq \min_{\boldsymbol{x}'} L(\boldsymbol{x}) + \frac{1}{2} \left(\boldsymbol{x}' - \boldsymbol{x} + (\boldsymbol{H}^*)^{-1} \nabla L(\boldsymbol{x}) \right)^\top \boldsymbol{H}^* \left(\boldsymbol{x}' - \boldsymbol{x} + (\boldsymbol{H}^*)^{-1} \nabla L(\boldsymbol{x}) \right) - \frac{1}{2} \nabla L(\boldsymbol{x})^\top (\boldsymbol{H}^*)^{-1} \nabla L(\boldsymbol{x}) \\ &= L(\boldsymbol{x}) - \frac{1}{2} \nabla L(\boldsymbol{x})^\top (\boldsymbol{H}^*)^{-1} \nabla L(\boldsymbol{x}). \end{split}$$

Applying the above inequality to x_1, \ldots, x_T and denoting $\bar{g}_t = \nabla L(x_t)$, we then have

$$\mathbb{E}\left[\sum_{t=1}^{T} (L(\boldsymbol{x}_{t}) - L(\boldsymbol{x}^{*}))\right] \geq \frac{1}{2} \mathbb{E}\left[\sum_{t=1}^{T} \bar{\boldsymbol{g}}_{t}^{\top} (\boldsymbol{H}^{*})^{-1} \bar{\boldsymbol{g}}_{t}\right] \\
= \frac{1}{2} \left\langle \mathbb{E}\left[\sum_{t=1}^{T} \bar{\boldsymbol{g}}_{t} \bar{\boldsymbol{g}}_{t}^{\top}\right], (\boldsymbol{H}^{*})^{-1} \right\rangle \\
= \frac{1}{2 \operatorname{Tr}(\boldsymbol{H}^{*})} \left\langle \mathbb{E}\left[\sum_{t=1}^{T} \bar{\boldsymbol{g}}_{t} \bar{\boldsymbol{g}}_{t}^{\top}\right], (\boldsymbol{H}^{*}/\operatorname{Tr}(\boldsymbol{H}^{*}))^{-1} \right\rangle \\
\geq \frac{1}{2H(L,\mathcal{H})} \inf_{\operatorname{Tr}(\boldsymbol{H}) \leq 1, \boldsymbol{H} \in \mathcal{H}} \left\langle \mathbb{E}\left[\sum_{t=1}^{T} \bar{\boldsymbol{g}}_{t} \bar{\boldsymbol{g}}_{t}^{\top}\right], \boldsymbol{H}^{-1} \right\rangle.$$
(27)

For any $\delta > 0$, we choose $H_g \in \mathcal{H}$ such that $\langle \mathbb{E}[\sum_{t=1}^T \bar{g}_t \bar{g}_t^\top], H_g^{-1} \rangle \leq \delta + \inf_{H \in \mathcal{H}, \operatorname{Tr}(H) \leq 1} \langle \mathbb{E}[\sum_{t=1}^T \bar{g}_t \bar{g}_t^\top], H^{-1} \rangle$ and $\operatorname{Tr}(H_g) \leq 1$. Similarly, we choose $H_{\Sigma} \in \mathcal{H}$ such that $\langle \Sigma, H_{\Sigma}^{-1} \rangle \leq \delta + \inf_{H \in \mathcal{H}, \operatorname{Tr}(H) \leq 1} \langle \Sigma, H^{-1} \rangle$ and $\operatorname{Tr}(H_{\Sigma}) \leq 1$. Correspondingly, we define $H' = \alpha H_g + (1 - \alpha) H_{\Sigma}$, which satisfies that $H' \in \mathcal{H}$ and $\operatorname{Tr}(H') \leq 1$. Then by Jensen's inequality, we have

$$\mathbb{E}\inf_{\boldsymbol{H}\in\mathcal{H},\mathrm{Tr}(\boldsymbol{H})\leq 1}\sqrt{\langle \boldsymbol{M}_T-\boldsymbol{M}_0,\boldsymbol{H}^{-1}\rangle}\leq \mathbb{E}\sqrt{\langle \boldsymbol{M}_T-\boldsymbol{M}_0,(\boldsymbol{H}')^{-1}\rangle}\leq \sqrt{\mathbb{E}\langle \boldsymbol{M}_T-\boldsymbol{M}_0,(\boldsymbol{H}')^{-1}\rangle}.$$

Plugging in $M_T - M_0 = \sum_{t=1}^T g_t g_t^{\top}$, since $\mathbb{E}[\sum_{t=1}^t g_t g_t^{\top}] \leq \mathbb{E}[\sum_{t=1}^T \bar{g}_t \bar{g}_t^{\top}] + T\Sigma$, we further have

$$\mathbb{E} \inf_{\boldsymbol{H}\in\mathcal{H},\mathrm{Tr}(\boldsymbol{H})\leq 1} \sqrt{\langle \boldsymbol{M}_{T}-\boldsymbol{M}_{0},\boldsymbol{H}^{-1}\rangle} \leq \sqrt{\left\langle \mathbb{E}\left[\sum_{t=1}^{T} \bar{\boldsymbol{g}}_{t} \bar{\boldsymbol{g}}_{t}^{\mathsf{T}}\right], (\boldsymbol{H}')^{-1}\right\rangle} + T\left\langle \boldsymbol{\Sigma}, (\boldsymbol{H}')^{-1}\right\rangle}$$
$$\leq \sqrt{\frac{1}{\alpha} \left\langle \mathbb{E}\left[\sum_{t=1}^{T} \bar{\boldsymbol{g}}_{t} \bar{\boldsymbol{g}}_{t}^{\mathsf{T}}\right], \boldsymbol{H}_{\boldsymbol{g}}^{-1}\right\rangle} + \frac{1}{1-\alpha} T\left\langle \boldsymbol{\Sigma}, \boldsymbol{H}_{\boldsymbol{\Sigma}}^{-1}\right\rangle}$$
$$= \sqrt{\left\langle \mathbb{E}\left[\sum_{t=1}^{T} \bar{\boldsymbol{g}}_{t} \bar{\boldsymbol{g}}_{t}^{\mathsf{T}}\right], \boldsymbol{H}_{\boldsymbol{g}}^{-1}\right\rangle} + \sqrt{T\left\langle \boldsymbol{\Sigma}, \boldsymbol{H}_{\boldsymbol{\Sigma}}^{-1}\right\rangle}$$

where in the last step we choose α to be

$$\alpha = \frac{\sqrt{\langle \mathbb{E}[\sum_{t=1}^{T} \bar{\boldsymbol{g}}_{t} \bar{\boldsymbol{g}}_{t}^{\top}], \boldsymbol{H}_{\boldsymbol{g}}^{-1} \rangle}}{\sqrt{\langle \mathbb{E}[\sum_{t=1}^{T} \bar{\boldsymbol{g}}_{t} \bar{\boldsymbol{g}}_{t}^{\top}], \boldsymbol{H}_{\boldsymbol{g}}^{-1} \rangle} + \sqrt{T \langle \boldsymbol{\Sigma}, \boldsymbol{H}_{\boldsymbol{\Sigma}}^{-1} \rangle}}$$

Recall that $\sigma = \inf_{\boldsymbol{H} \in \mathcal{H}, \operatorname{Tr}(\boldsymbol{H}) \leq 1} \sqrt{\langle \boldsymbol{\Sigma}, \boldsymbol{H}^{-1} \rangle}$. Then it follows from (27) and the definitions of $\boldsymbol{H}_{\boldsymbol{g}}, \boldsymbol{H}_{\boldsymbol{\Sigma}}$ that

$$\mathbb{E}\Big[\inf_{\boldsymbol{H}\in\mathcal{H},\operatorname{Tr}(\boldsymbol{H})\leq 1}\sqrt{\langle \boldsymbol{M}_{T}-\boldsymbol{M}_{0},\boldsymbol{H}^{-1}\rangle}\Big]\leq \sqrt{2H(L,\mathcal{H})\cdot\mathbb{E}\Big[\sum_{t=1}^{T}(L_{t}(\boldsymbol{x}_{t})-L_{t}(\boldsymbol{x}^{*}))\Big]+\delta+\sqrt{T\sigma^{2}+T\delta}}$$

Since $\delta > 0$ is arbitrary, this implies that

$$\mathbb{E}\Big[\inf_{\boldsymbol{H}\in\mathcal{H},\operatorname{Tr}(\boldsymbol{H})\leq 1}\sqrt{\langle \boldsymbol{M}_{T}-\boldsymbol{M}_{0},\boldsymbol{H}^{-1}\rangle}\Big]\leq \sqrt{2H(L,\mathcal{H})\cdot\mathbb{E}\Big[\sum_{t=1}^{T}(L_{t}(\boldsymbol{x}_{t})-L_{t}(\boldsymbol{x}^{*}))\Big]+\sqrt{T\sigma^{2}}}$$

Now recall the regret bound from Corollary 3.5, and it follows from the above inequality that

$$\mathbb{E}\left[\sum_{t=1}^{T} (L_t(\boldsymbol{x}_t) - L_t(\boldsymbol{x}^*))\right] \le 2\sqrt{2} \|\mathcal{X}\|_{\mathcal{H}} \cdot \mathbb{E}\left[\inf_{\boldsymbol{H}\in\mathcal{H}, \operatorname{Tr}(\boldsymbol{H})\le 1} \sqrt{\langle \boldsymbol{M}_t - \boldsymbol{M}_0, \boldsymbol{H}^{-1} \rangle} + d\sqrt{\epsilon}\right] \\ \le 2\sqrt{2} \|\mathcal{X}\|_{\mathcal{H}} \left(\sqrt{2H(L, \mathcal{H}) \cdot \mathbb{E}\left[\sum_{t=1}^{T} (L_t(\boldsymbol{x}_t) - L_t(\boldsymbol{x}^*))\right]} + \sqrt{T\sigma^2} + d\sqrt{\epsilon}\right).$$

Solving the above inequality yields

$$\mathbb{E}\left[\frac{1}{T}\sum_{t=1}^{T}(L_t(\boldsymbol{x}_t) - L_t(\boldsymbol{x}^*))\right] \leq \frac{16}{T} \|\mathcal{X}\|_{\mathcal{H}}^2 H(L,\mathcal{H}) + \frac{4\sqrt{2}}{\sqrt{T}} \|\mathcal{X}\|_{\mathcal{H}} \sigma + \frac{4\sqrt{2}}{T} \|\mathcal{X}\|_{\mathcal{H}} d\sqrt{\epsilon}.$$

This completes the proof.

D.3. Analysis for EMA style optimizers

Here we present the proof for the theorems in Appendix C.

Theorem C.1. Let \mathcal{H} be a well-structured preconditioner set under Definition 3.1. Then for any convex loss functions $\tilde{L}_1, \ldots, \tilde{L}_T$, the regret of Algorithm 5 compared to any $\mathbf{x}^* \in \mathcal{X}$ can be bounded as

$$\sum_{t=1}^{T} \beta_2^{\frac{T-t}{2}} \left[\tilde{L}_t(\boldsymbol{x}_t) - \tilde{L}_t(\boldsymbol{x}^*) \right] \leq \left(\frac{D^2}{2\tilde{\eta}} + \tilde{\eta} \right) \inf_{\boldsymbol{H} \in \mathcal{H}, \operatorname{Tr}(\boldsymbol{H}) \leq 1} \sqrt{\left\langle \tilde{\boldsymbol{M}}_T, \boldsymbol{H}^{-1} \right\rangle}$$

where $D = \max_{t \in [T]} \|\boldsymbol{x}_t - \boldsymbol{x}^*\|_{\mathcal{H}}$.

Proof of Theorem C.1. We can choose $L_t = \beta_2^{-\frac{t}{2}} \tilde{L}_t$, $\epsilon = \tilde{\epsilon}$ and $\eta = \tilde{\eta}$. Then we claim Algorithm 5 is equivalent to Algorithm 1 with hyperparameter ϵ, η and loss functions $\{L_t\}_{t=1}^T$, which will be shown by induction.

Assume $\tilde{\boldsymbol{x}}_s = \boldsymbol{x}_s$ for $s \leq t$. Then we know $\boldsymbol{g}_s = \nabla L_s(\boldsymbol{x}_s) = \beta_2^{-\frac{s}{2}} \nabla \tilde{L}_s(\tilde{\boldsymbol{x}}_s)$ for $s \leq t$.

We consider the update in the step t of Algorithm 5.

$$\begin{split} \tilde{\boldsymbol{M}}_{t} &= \sum_{i=1}^{t} \beta_{2}^{t-i} \tilde{\boldsymbol{g}}_{i} \tilde{\boldsymbol{g}}_{i}^{\top} + \beta_{2}^{t} \tilde{\boldsymbol{M}}_{0} = \beta_{2}^{t} \left[\sum_{i=1}^{t} \boldsymbol{g}_{t} \boldsymbol{g}_{t}^{\top} + \epsilon \boldsymbol{I}_{d} \right] = \beta_{2}^{t} \boldsymbol{M}_{t}, \\ \tilde{\boldsymbol{H}}_{t} &= \operatorname*{arg\,min}_{\tilde{\boldsymbol{H}} \in \mathcal{H}} \left\langle \tilde{\boldsymbol{M}}_{t}, \tilde{\boldsymbol{H}}^{-1} \right\rangle + \tilde{\eta}^{2} \operatorname{Tr}(\tilde{\boldsymbol{H}}) = \operatorname*{arg\,min}_{\boldsymbol{H} \in \mathcal{H}} \left\langle \beta_{2}^{t} \boldsymbol{M}_{t}, \boldsymbol{H}^{-1} \right\rangle + \eta^{2} \operatorname{Tr}(\boldsymbol{H}) = \beta_{2}^{\frac{t}{2}} \boldsymbol{H}_{t}, \\ \\ \tilde{\boldsymbol{H}}_{t}^{-1} \tilde{\boldsymbol{g}}_{t} &= \beta_{2}^{-\frac{t}{2}} \boldsymbol{H}_{t}^{-1} \beta_{2}^{\frac{t}{2}} \boldsymbol{g}_{t} = \boldsymbol{H}_{t}^{-1} \boldsymbol{g}_{t}, \\ \\ \tilde{\boldsymbol{x}}_{t+1} &= \Pi_{\mathcal{X}}^{\tilde{\boldsymbol{H}}_{t}} \left(\tilde{\boldsymbol{x}}_{t} - \tilde{\boldsymbol{H}}_{t}^{-1} \tilde{\boldsymbol{g}}_{t} \right) = \Pi_{\mathcal{X}}^{\tilde{\boldsymbol{H}}_{t}} \left(\boldsymbol{x}_{t} - \boldsymbol{H}_{t}^{-1} \boldsymbol{g}_{t} \right) = \Pi_{\mathcal{X}}^{H_{t}} \left(\boldsymbol{x}_{t} - \boldsymbol{H}_{t}^{-1} \boldsymbol{g}_{t} \right) = \boldsymbol{x}_{t+1}. \end{split}$$

Therefore, we can obtain the regret bound for Algorithm 4 with Theorem 3.4.

$$\sum_{t=1}^{T} \beta_2^{-\frac{t}{2}} \left[\tilde{L}_t(\boldsymbol{x}_t) - \tilde{L}_t(\boldsymbol{x}^*) \right] = \sum_{t=1}^{T} \left[L_t(\boldsymbol{x}_t) - L_t(\boldsymbol{x}^*) \right]$$
$$\leq \left(\frac{\max_{1 \le t \le T} \|\boldsymbol{x}_t - \boldsymbol{x}^*\|_{\mathcal{H}}^2}{2\eta} + \eta \right) \inf_{\boldsymbol{H} \in \mathcal{H}, \operatorname{Tr}(\boldsymbol{H}) \le 1} \sqrt{\langle \boldsymbol{M}_T, \boldsymbol{H}^{-1} \rangle}$$
$$= \beta_2^{-\frac{T}{2}} \left(\frac{\max_{1 \le t \le T} \|\boldsymbol{x}_t - \boldsymbol{x}^*\|_{\mathcal{H}}^2}{2\tilde{\eta}} + \tilde{\eta} \right) \inf_{\boldsymbol{H} \in \mathcal{H}, \operatorname{Tr}(\boldsymbol{H}) \le 1} \sqrt{\langle \tilde{\boldsymbol{M}}_T, \boldsymbol{H}^{-1} \rangle}$$

and

$$\sum_{t=1}^{T} \beta_{2}^{-\frac{t-1}{2}} \frac{\beta_{2}^{-\frac{1}{2}} - 1}{\beta_{2}^{-\frac{T}{2}} - 1} \left[\tilde{L}_{t}(\boldsymbol{x}_{t}) - \tilde{L}_{t}(\boldsymbol{x}^{*}) \right] \leq \frac{1 - \sqrt{\beta_{2}}}{1 - \beta_{2}^{\frac{T}{2}}} \left(\frac{\max_{1 \leq t \leq T} \|\boldsymbol{x}_{t} - \boldsymbol{x}^{*}\|_{\mathcal{H}}^{2}}{2\tilde{\eta}} + \tilde{\eta} \right) \inf_{\boldsymbol{H} \in \mathcal{H}, \operatorname{Tr}(\boldsymbol{H}) \leq 1} \sqrt{\left\langle \tilde{\boldsymbol{M}}_{T}, \boldsymbol{H}^{-1} \right\rangle}$$

Theorem C.2. Let \mathcal{H} be a well-structured preconditioner set under Definition 3.1. Consider any independent stochastic convex loss functions $\tilde{L}_1, \ldots, \tilde{L}_T$ satisfying Assumption 3.7, and let $H(\tilde{L}, \mathcal{H})$ be the \mathcal{H} -smoothness of their expectation \tilde{L} . Suppose the global minimizer of L, denoted by \mathbf{x}^* , is in \mathcal{X} . Then for the iterates $\mathbf{x}_1, \ldots, \mathbf{x}_T$ of Algorithm 5, denoting $\bar{\mathbf{x}}_{1:T} = (\sum_{t=1}^T \beta_2^{\frac{T-t}{2}})^{-1} \sum_{t=1}^T \beta_2^{\frac{T-t}{2}} \mathbf{x}_t$, it holds that

$$\mathbb{E}\tilde{L}(\bar{\boldsymbol{x}}_{1:T}) - \tilde{L}_{t}(\boldsymbol{x}^{*}) \leq \frac{16}{\sum_{t=1}^{T} \beta_{2}^{\frac{T-t}{2}}} \left\| \mathcal{X} \right\|_{\mathcal{H}}^{2} H(\tilde{L},\mathcal{H}) + \frac{4\sqrt{2}}{\sqrt{\sum_{t=1}^{T} \beta_{2}^{T-t}}} \left\| \mathcal{X} \right\|_{\mathcal{H}} \sigma + \frac{4\sqrt{2}}{\sum_{t=1}^{T} \beta_{2}^{\frac{-t}{2}}} \left\| \mathcal{X} \right\|_{\mathcal{H}} d\sqrt{\epsilon}$$

where $\sigma = \inf_{\boldsymbol{H} \in \mathcal{H}, \operatorname{Tr}(\boldsymbol{H}) \leq 1} \sqrt{\langle \boldsymbol{\Sigma}, \boldsymbol{H}^{-1} \rangle}.$

Proof of Theorem C.2. Similar to the proof of Theorem 3.8, we know

$$\tilde{L}(\tilde{\boldsymbol{x}}_t) - \tilde{L}(\tilde{\boldsymbol{x}}^*) \ge \left\langle \mathbb{E} \tilde{\boldsymbol{g}}_t \mathbb{E} \tilde{\boldsymbol{g}}_t^\top, (\tilde{\boldsymbol{H}}^*)^{-1} \right\rangle$$

and

$$\sum_{t=1}^{T} \beta_2^{\frac{T-t}{2}} \tilde{L}(\tilde{\boldsymbol{x}}_t) - \tilde{L}(\tilde{\boldsymbol{x}}^*) \geq \left\langle \sum_{t=1}^{T} \beta_2^{\frac{T-t}{2}} \mathbb{E} \tilde{\boldsymbol{g}}_t \mathbb{E} \tilde{\boldsymbol{g}}_t^{\top}, (\tilde{\boldsymbol{H}}^*)^{-1} \right\rangle$$
$$\geq \frac{1}{2H(\tilde{L}, \mathcal{H})} \inf_{\operatorname{Tr}(\boldsymbol{H}) \leq 1, \boldsymbol{H} \in \mathcal{H}} \left\langle \sum_{t=1}^{T} \beta_2^{\frac{T-t}{2}} \mathbb{E} \tilde{\boldsymbol{g}}_t \mathbb{E} \tilde{\boldsymbol{g}}_t^{\top}, \boldsymbol{H}^{-1} \right\rangle$$

For any $\delta>0,$ we choose $oldsymbol{H_g}\in\mathcal{H}$ such that $\mathrm{Tr}(oldsymbol{H_g})\leq 1$ and

$$\left\langle \sum_{t=1}^{T} \beta_2^{\frac{T-t}{2}} \mathbb{E} \tilde{\boldsymbol{g}}_t \mathbb{E} \tilde{\boldsymbol{g}}_t^{\top}, \boldsymbol{H}_{\boldsymbol{g}}^{-1} \right\rangle \leq \delta + \inf_{\boldsymbol{H} \in \mathcal{H}, \operatorname{Tr}(\boldsymbol{H}) \leq 1} \left\langle \sum_{t=1}^{T} \beta_2^{\frac{T-t}{2}} \mathbb{E} \tilde{\boldsymbol{g}}_t \mathbb{E} \tilde{\boldsymbol{g}}_t^{\top}, \boldsymbol{H}^{-1} \right\rangle.$$

We also choose $H_{\Sigma} \in \mathcal{H}$ such that $\langle \Sigma, H_{\Sigma}^{-1} \rangle \leq \delta + \inf_{H \in \mathcal{H}, \operatorname{Tr}(H) \leq 1} \langle \Sigma, H^{-1} \rangle$ and $\operatorname{Tr}(H_{\Sigma}) \leq 1$. We define $H' = \alpha H_g + (1 - \alpha) H_{\Sigma}$. Then $H' \in \mathcal{H}$ and $\operatorname{Tr}(H') \leq 1$. We have that

$$\mathbb{E} \inf_{\boldsymbol{H} \in \mathcal{H}, \operatorname{Tr}(\boldsymbol{H}) \leq 1} \sqrt{\left\langle \tilde{\boldsymbol{M}}_{T} - \beta_{2}^{T} \tilde{\boldsymbol{M}}_{0}, \boldsymbol{H}^{-1} \right\rangle} \leq \mathbb{E} \sqrt{\left\langle \tilde{\boldsymbol{M}}_{T} - \beta_{2}^{T} \tilde{\boldsymbol{M}}_{0}, (\boldsymbol{H}')^{-1} \right\rangle} \\ \leq \sqrt{\mathbb{E} \left\langle \tilde{\boldsymbol{M}}_{T} - \beta_{2}^{T} \tilde{\boldsymbol{M}}_{0}, (\boldsymbol{H}')^{-1} \right\rangle} \\ \leq \sqrt{\left\langle \sum_{t=1}^{T} \beta_{2}^{T-t} \mathbb{E} \tilde{\boldsymbol{g}}_{t} \mathbb{E} \tilde{\boldsymbol{g}}_{t}^{\top}, (\boldsymbol{H}')^{-1} \right\rangle} + \sum_{t=1}^{T} \beta_{2}^{T-t} \left\langle \boldsymbol{\Sigma}, (\boldsymbol{H}')^{-1} \right\rangle}.$$

Now plugging in the definition of H', we obtain

$$\mathbb{E}\inf_{\boldsymbol{H}\in\mathcal{H},\mathrm{Tr}(\boldsymbol{H})\leq 1}\sqrt{\left\langle\tilde{\boldsymbol{M}}_{T}-\boldsymbol{\beta}_{2}^{T}\tilde{\boldsymbol{M}}_{0},\boldsymbol{H}^{-1}\right\rangle}\leq\sqrt{\left\langle\sum_{t=1}^{T}\boldsymbol{\beta}_{2}^{T-t}\mathbb{E}\tilde{\boldsymbol{g}}_{t}\mathbb{E}\tilde{\boldsymbol{g}}_{t}^{\top},(\boldsymbol{\alpha}\boldsymbol{H}_{\boldsymbol{g}})^{-1}\right\rangle}+\sum_{t=1}^{T}\boldsymbol{\beta}_{2}^{T-t}\left\langle\boldsymbol{\Sigma},((1-\boldsymbol{\alpha})\boldsymbol{H}_{\boldsymbol{\Sigma}})^{-1}\right\rangle}$$
$$=\sqrt{\frac{1}{\alpha}\left\langle\sum_{t=1}^{T}\boldsymbol{\beta}_{2}^{\frac{T-t}{2}}\mathbb{E}\tilde{\boldsymbol{g}}_{t}\mathbb{E}\tilde{\boldsymbol{g}}_{t}^{\top},(\boldsymbol{H}_{\boldsymbol{g}})^{-1}\right\rangle}+\frac{1}{1-\boldsymbol{\alpha}}\sum_{t=1}^{T}\boldsymbol{\beta}_{2}^{T-t}\left\langle\boldsymbol{\Sigma},(\boldsymbol{H}_{\boldsymbol{\Sigma}})^{-1}\right\rangle}$$
$$\leq\sqrt{\left\langle\sum_{t=1}^{T}\boldsymbol{\beta}_{2}^{\frac{T-t}{2}}\mathbb{E}\tilde{\boldsymbol{g}}_{t}\mathbb{E}\tilde{\boldsymbol{g}}_{t}^{\top},(\boldsymbol{H}_{\boldsymbol{g}})^{-1}\right\rangle}+\sqrt{\sum_{t=1}^{T}\boldsymbol{\beta}_{2}^{T-t}\left\langle\boldsymbol{\Sigma},(\boldsymbol{H}_{\boldsymbol{\Sigma}})^{-1}\right\rangle}.$$

where in the last step we choose α to be

$$\alpha = \frac{\sqrt{\left\langle \sum_{t=1}^{T} \beta_2^{\frac{T-t}{2}} \mathbb{E} \tilde{\boldsymbol{g}}_t \mathbb{E} \tilde{\boldsymbol{g}}_t^{\top}, (\boldsymbol{H}_{\boldsymbol{g}})^{-1} \right\rangle}}{\sqrt{\left\langle \sum_{t=1}^{T} \beta_2^{\frac{T-t}{2}} \mathbb{E} \tilde{\boldsymbol{g}}_t \mathbb{E} \tilde{\boldsymbol{g}}_t^{\top}, (\boldsymbol{H}_{\boldsymbol{g}})^{-1} \right\rangle} + \sqrt{\sum_{t=1}^{T} \beta_2^{T-t} \left\langle \boldsymbol{\Sigma}, (\boldsymbol{H}_{\boldsymbol{\Sigma}})^{-1} \right\rangle}}$$

Then we have that

$$\mathbb{E}\inf_{\boldsymbol{H}\in\mathcal{H},\mathrm{Tr}(\boldsymbol{H})\leq 1}\sqrt{\left\langle \tilde{\boldsymbol{M}}_{T}-\boldsymbol{\beta}_{2}^{T}\tilde{\boldsymbol{M}}_{0},\boldsymbol{H}^{-1}\right\rangle}\leq\sqrt{2H(\tilde{\boldsymbol{L}},\mathcal{H})\left(\mathbb{E}\sum_{t=1}^{T}\boldsymbol{\beta}_{2}^{\frac{T-t}{2}}[L_{t}(\boldsymbol{x}_{t})-L_{t}(\boldsymbol{x}^{*})]\right)+\delta}+\sqrt{\sum_{t=1}^{T}\boldsymbol{\beta}_{2}^{T-t}}\sqrt{\sigma^{2}+\delta}$$

When taking δ to 0, we have that

$$\mathbb{E}\sum_{t=1}^{T}\beta_{2}^{\frac{T-t}{2}}\left[\tilde{L}_{t}(\boldsymbol{x}_{t})-\tilde{L}_{t}(\boldsymbol{x}^{*})\right] \leq 2\sqrt{2}D\left[\sqrt{2H(\tilde{L},\mathcal{H})\left(\mathbb{E}\sum_{t=1}^{T}\beta_{2}^{\frac{T-t}{2}}[L_{t}(\boldsymbol{x}_{t})-L_{t}(\boldsymbol{x}^{*})]\right)}+\sqrt{\sum_{t=1}^{T}\beta_{2}^{T-t}}\sigma+\beta_{2}^{\frac{T}{2}}d\sqrt{\epsilon}\right]$$

and

$$\mathbb{E}\sum_{t=1}^{T}\beta_{2}^{\frac{T-t}{2}}\left[\tilde{L}_{t}(\boldsymbol{x}_{t})-\tilde{L}_{t}(\boldsymbol{x}^{*})\right] \leq 16D^{2}H(\tilde{L},\mathcal{H})+4\sqrt{2\sum_{t=1}^{T}\beta_{2}^{T-t}}D\sigma+4\sqrt{2}D\beta_{2}^{\frac{T}{2}}d\sqrt{\epsilon}$$

If we choose $\bar{x}_{1:T} = \frac{1}{\sum_{t=1}^{T} \beta_2^{\frac{T-t}{2}}} \sum_{t=1}^{T} \beta_2^{\frac{T-t}{2}} x_t$, then from the convexity of \tilde{L} we know that

$$\begin{split} \mathbb{E}[\tilde{L}(\bar{\boldsymbol{x}}_{1:T}) - \tilde{L}(\boldsymbol{x}^*)] &\leq \mathbb{E}\frac{1}{\sum_{t=1}^{T} \beta_2^{\frac{T-t}{2}}} \sum_{t=1}^{T} \beta_2^{\frac{T-t}{2}} \left[\tilde{L}_t(\boldsymbol{x}_t) - \tilde{L}_t(\boldsymbol{x}^*) \right] \\ &\leq \frac{1}{\sum_{t=1}^{T} \beta_2^{\frac{T-t}{2}}} \left(16D^2 H(\tilde{L}, \mathcal{H}) + 4\sqrt{2} \sum_{t=1}^{T} \beta_2^{T-t} D\sigma + 4\sqrt{2} D\beta_2^{\frac{T}{2}} d\sqrt{\epsilon} \right) \\ &\leq \frac{16}{\sum_{t=1}^{T} \beta_2^{\frac{T-t}{2}}} D^2 H(\tilde{L}, \mathcal{H}) + \frac{4\sqrt{2}}{\sqrt{\sum_{t=1}^{T} \beta_2^{T-t}}} D\sigma + \frac{4\sqrt{2}}{\sum_{t=1}^{T} \beta_2^{\frac{-t}{2}}} Dd\sqrt{\epsilon} \end{split}$$
etes the proof.

This completes the proof.

E. Proof for One-Sided Shampoo

E.1. Proof for left smoothness

Lemma E.1 (Left smoothness for one-sided Shampoo). Let $\mathcal{H} = S^{d_L}_+ \otimes I_{d_R}$ be the well-structured preconditioner set for one-sided Shampoo. Then the \mathcal{H} -smoothness $H(L, \mathcal{H})$ defined in Definition 3.6 is equal to the smallest number $H \ge 0$ such that there exists $\mathbf{H}^*_{d_L} \in \mathbb{R}^{d_L \times d_L}$ satisfying that $H = d_R \operatorname{Tr}(\mathbf{H}^*_{d_L})$ and that for any $\mathbf{X}, \mathbf{\Delta} \in \mathbb{R}^{d_L \times d_R}$,

$$\left| \nabla^2 L(\boldsymbol{X})[\boldsymbol{\Delta}, \boldsymbol{\Delta}] \right| \leq \left\langle \boldsymbol{H}_{d_L}^*, \boldsymbol{\Delta} \boldsymbol{\Delta}^\top \right\rangle$$

In this case, the H-smoothness is denoted by $H_{left}(L)$.

Proof of Lemma 4.3. First, for any $X, \Delta \in \mathbb{R}^{d_L \times d_R}$ and $x = \overline{\operatorname{vec}}(X)$, we have

$$\overline{\operatorname{vec}}(\boldsymbol{\Delta})^{\top} \nabla^2 L(\boldsymbol{x}) \, \overline{\operatorname{vec}}(\boldsymbol{\Delta}) = \nabla^2 L(\boldsymbol{X}) [\boldsymbol{\Delta}, \boldsymbol{\Delta}]$$

Therefore, given the Kronecker product form of all $H \in \mathcal{H}$, to find $H \in \mathcal{H}$ with the smallest trace such that $-H \preceq \nabla^2 L(\mathbf{x}) \preceq H$, it is equivalent to find $H = H_{d_L} \otimes I_{d_R} \in \mathcal{H}$ with the smallest trace such that for any $\Delta \in \mathbb{R}^{d_L \times d_R}$,

$$egin{aligned} |
abla^2 L(oldsymbol{X})[oldsymbol{\Delta},oldsymbol{\Delta}]| &\leq \overline{ ext{vec}}(oldsymbol{\Delta})^{ op}(oldsymbol{H}_{d_L}\otimesoldsymbol{I}_{d_R})\,\overline{ ext{vec}}(oldsymbol{\Delta}) \ &= \operatorname{Tr}(oldsymbol{\Delta}^ opoldsymbol{H}_{d_L}oldsymbol{\Delta}) \ &= \langle oldsymbol{H}_{d_L},oldsymbol{\Delta}oldsymbol{\Delta}^ op
angle. \end{aligned}$$

Further note that $\operatorname{Tr}(\boldsymbol{H}) = d_R \operatorname{Tr}(\boldsymbol{H}_{d_L})$, and thus we conclude that it is equivalent to find $\boldsymbol{H}_{d_L} \in \mathcal{S}^{d_L}_+$ with the smallest trace such that the above inequality holds for all $\boldsymbol{\Delta} \in \mathbb{R}^{d_L \times d_R}$. This completes the proof.

E.2. Proof for regret bound

Theorem 4.2 (Regret bound for one-sided Shampoo). For convex functions L_1, \ldots, L_T , the regret of one-sided Shampoo (Algorithm 2) compared to any $\mathbf{X}^* \in \mathbb{R}^{d_L \times d_R}$ satisfies

$$\sum_{t=1}^{T} L_t(\boldsymbol{X}_t) - \sum_{t=1}^{T} L_t(\boldsymbol{X}^*) \le \left(\frac{D_{\text{op}}^2}{2d_R\eta} + \eta\right) \left(G + d\sqrt{\epsilon}\right),$$

where $D_{\text{op}} = \max_{t \in [T]} \| \boldsymbol{X}_t - \boldsymbol{X}^* \|_{\text{op}}$ and $G = \sqrt{d_R} \operatorname{Tr} \left[\left(\sum_{t=1}^T \boldsymbol{G}_t \boldsymbol{G}_t^\top \right)^{\frac{1}{2}} \right]$. When the domain \mathcal{X} is bounded in operator norm, i.e., $\| \mathcal{X} \|_{\text{op}} < \infty$, further choosing $\eta = \sqrt{2/d_R} \| \mathcal{X} \|_{\text{op}}$, it holds

$$\sum_{t=1}^{T} L_t(\boldsymbol{X}_t) - \sum_{t=1}^{T} L_t(\boldsymbol{X}^*)$$

$$\leq 2\sqrt{2} \|\mathcal{X}\|_{\text{op}} \left(\operatorname{Tr} \left[\left(\sum_{t=1}^{T} \boldsymbol{G}_t \boldsymbol{G}_t^{\mathsf{T}} \right)^{\frac{1}{2}} \right] + \frac{d}{\sqrt{d_R}} \sqrt{\epsilon} \right).$$

Proof of Theorem 4.2. We will apply Theorem 3.4 to one-sided Shampoo. According to the analysis in Appendix B.4, Algorithm 1 with $\mathcal{H} = \left(\mathbb{R}^{d_L \times d_L} \otimes I_{d_R}\right) \cap \mathcal{S}^d_+$ recovers one-sided Shampoo. We can plug in $\|\boldsymbol{x}\|_{\mathcal{H}} = \frac{\|\boldsymbol{X}\|_{\text{op}}}{\sqrt{d_R}}$ and $\|\|\boldsymbol{g}_{1:T}\|\|_{\mathcal{H}} = \sqrt{d_R} \operatorname{Tr} \left[\left(\sum_{t=1}^T \boldsymbol{G}_t \boldsymbol{G}_t^\top \right)^{\frac{1}{2}} \right]$ into Theorem 3.4 and get that $\sum_{t=1}^T L_t(\boldsymbol{X}_t) - \sum_{t=1}^T L_t(\boldsymbol{X}^*) \le \sqrt{2} \left(\frac{D_{\text{op}}^2}{2d_R \eta} + \eta \right) \left(G + \min\left(d\sqrt{\epsilon}, \frac{d^2\epsilon}{2G} \right) \right)$ with $D_{\text{op}} = \max_{t \in [T]} \|\boldsymbol{X}_t - \boldsymbol{X}^*\|_{\text{op}}$ and $G = \sqrt{d_R} \operatorname{Tr} \left[\left(\sum_{t=1}^T \boldsymbol{G}_t \boldsymbol{G}_t^\top \right)^{\frac{1}{2}} \right]$.

F. Additional Results for Experiments

F.1. Efficient implementation of full-matrix AdaGrad.

Directly applying full-matrix AdaGrad to this 10⁶-dimensional problem is impractical. Instead, we consider the eigendecomposition of H to be $U^{\top} \Sigma U$ and define the transformation $\mathcal{T}(X) = UX$. We further define d orthogonal matrices $V_1, \ldots, V_d \in \mathbb{R}^{d \times d}$ such that the first row of V_i is in the same direction of $\mathcal{T}(X^* - X_0)_{i,:}$ and define $V = \text{diag}(V_1, \ldots, V_d)$. We can know that $V \overline{\text{vec}}(\mathcal{T}(X^* - X_0) = v \otimes e_1$ where $v_i = \|\mathcal{T}(X^* - X_0)_{i,:}\|_2$ for $i \in [d]$.

Then it holds that

$$\begin{split} f(\boldsymbol{X}) &= \left\langle \boldsymbol{H}, (\boldsymbol{X} - \boldsymbol{X}^*) (\boldsymbol{X} - \boldsymbol{X}^*)^\top \right\rangle = \left\langle \boldsymbol{\Sigma}, (\boldsymbol{U}\boldsymbol{X} - \boldsymbol{U}\boldsymbol{X}^*) (\boldsymbol{U}\boldsymbol{X} - \boldsymbol{U}\boldsymbol{X}^*)^\top \right\rangle \\ &= \left\langle \boldsymbol{\Sigma}, (\mathcal{T}(\boldsymbol{X}) - \mathcal{T}(\boldsymbol{X}^*)) (\mathcal{T}(\boldsymbol{X}) - \mathcal{T}(\boldsymbol{X}^*))^\top \right\rangle \\ &= \left\langle \boldsymbol{\Sigma} \otimes \boldsymbol{I}_d, \overline{\operatorname{vec}} (\mathcal{T}(\boldsymbol{X} - \boldsymbol{X}^*)^\top) \overline{\operatorname{vec}} (\mathcal{T}(\boldsymbol{X} - \boldsymbol{X}^*)^\top)^\top \right\rangle \\ &= \left\langle \boldsymbol{\Sigma} \otimes \boldsymbol{I}_d, \boldsymbol{V} \overline{\operatorname{vec}} (\mathcal{T}(\boldsymbol{X} - \boldsymbol{X}^*)^\top) \overline{\operatorname{vec}} (\mathcal{T}(\boldsymbol{X} - \boldsymbol{X}^*)^\top)^\top \boldsymbol{V}^\top \right\rangle \end{split}$$

Further denoting $\tilde{\boldsymbol{x}} = \boldsymbol{V} \overline{\operatorname{vec}}(\mathcal{T}(\boldsymbol{X})^{\top})$ and $\tilde{\boldsymbol{y}} = \boldsymbol{V} \overline{\operatorname{vec}}(\mathcal{T}(\boldsymbol{X}_0)^{\top})$, then we obtain

$$\begin{split} f(\boldsymbol{X}) &= \sum_{i=1}^{d} \sigma_{i} \left\| \tilde{\boldsymbol{x}}_{(i-1)d+1:id} - (\boldsymbol{v} \otimes \boldsymbol{e}_{1})_{(i-1)d+1:id} \right\|_{2}^{2} \\ &= \sum_{i=1}^{d} \sigma_{i} \bigg[(\tilde{\boldsymbol{x}}_{(i-1)d+1} - \tilde{\boldsymbol{y}}_{(i-1)d+1} - \boldsymbol{v}_{i})^{2} + \sum_{j=2}^{d} (\tilde{\boldsymbol{x}}_{(i-1)d+j} - \tilde{\boldsymbol{y}}_{(i-1)d+j})^{2} \bigg] \end{split}$$

Running full-matrix AdaGrad on $f(\mathbf{X})$ starting from \mathbf{X}_0 can be implemented equivalently by using $\tilde{\mathbf{x}}$ as variable starting from $\tilde{\mathbf{y}}$. Only $\tilde{\mathbf{x}}_{(i-1)d+1}$ will receive non-zero gradient so full-matrix AdaGrad actually only cares these d coordinates, which reduces the original problem to a problem only with d variables.

F.2. Results for EMA algorithms

As mentioned in Section 5, we compare AdaSGD, Adam, one-sided EMA Shampoo and full-matrix AdaSGD, which are EMA version of AdaGrad-Norm, diagonal AdaGrad, one-sided Shampoo and full-matrix Adagrad. The results are plotted in Figure 2. We set $\beta_2 = 0.95$ and disable first-order momentum, i.e., $\beta_1 = 0$ in Adam.

We tried 60 learning rates between 1×10^{-4} and 1×10^{2} . The relationship between loss and learning rate is shown in Figure 3.



Figure 2. We plot the last iterate training loss $f(\mathbf{X}_t) = \langle \mathbf{H}, (\mathbf{X}_t - \mathbf{X}^*) (\mathbf{X}_t - \mathbf{X}^*)^\top \rangle$ and the average iterate training loss $f(\frac{1-\beta_2}{1-\beta_2}\sum_{s=1}^t \beta_2^{t-s} \mathbf{X}_s)$ over steps for optimizers obtained from Algorithm 5.



Figure 3. We plot the last iterate loss and average iterate loss versus learning rate. For each learning rate the plotted value is the average of last iterate loss and average iterate loss across five random seeds.