

Faster Adaptive Optimization via Expected Gradient Outer Product Reparameterization

PENDING

Abstract. Adaptive optimization algorithms—such as Adagrad, Adam, and their variants—have found widespread use in machine learning, signal processing and many other settings. Several methods in this family are not rotationally equivariant, meaning that simple reparameterizations (i.e. change of basis) can drastically affect their convergence. However, their sensitivity to the choice of parameterization has not been systematically studied; it is not clear how to identify a “favorable” change of basis in which these methods perform best. In this paper we propose a reparameterization method and demonstrate both theoretically and empirically its potential to improve their convergence behavior. Our method is an orthonormal transformation based on the *expected gradient outer product* (EGOP) matrix, which can be approximated using either full-batch or stochastic gradient oracles. Our theoretical results show that in the neighborhoods of local minima, the sensitivity of adaptive algorithms to choice-of-basis is influenced by spectral decay in the EGOP matrix. We illustrate the potential impact of EGOP reparameterization by presenting empirical evidence and theoretical arguments that common machine learning tasks with “natural” data exhibit EGOP spectral decay.

Key words. optimization, adaptive gradient methods, Adam, Adagrad, expected gradient outer product, reparameterization, low-rank, preconditioning, neural networks

MSC codes. 68T07 (Computer science: artificial NN’s and deep learning), 90C99 (Mathematical programming)

1. Introduction. Adaptive optimization algorithms are popular methods in modern machine learning [8]. Optimizers in this family include the seminal Adagrad and Adam algorithms as well as variants such as AdamW, Adadelata, and Adamax [22, 26, 32, 41]. These algorithms are termed “adaptive” because they maintain and update different learning rates for each coordinate in parameter space.¹ Despite their popularity, fully understanding the impact these adaptive learning rates have on convergence remains an area of ongoing research [21, 25, 27]. Notably, coordinate-wise learning rates make these methods sensitive to orthonormal reparameterization, distinguishing them from equivariant methods like gradient descent, whose performance does not change under orthonormal reparameterization.

Orthonormal reparameterizations correspond to full-dimensional changes of basis, and include seemingly benign transformations such as rotations of the loss landscape about the origin. The sensitivity of adaptive algorithms to rotation means that changes of basis can affect the rate of convergence to local minima, and even impact the generalization properties of the obtained solutions in the presence of non-convex landscapes. Given the sensitivity of these ubiquitous optimization methods to choice of basis, we pose the research question:

When using an adaptive algorithm, in what settings and to what extent can change-of-basis improve optimization?

We address this question by identifying geometric properties of loss functions that govern

¹We focus on adaptive algorithms with diagonal preconditioners, corresponding to the standard implementation of Adagrad, Adam, and variants in popular software packages. Full-matrix Adagrad, proposed by Duchi et al. [14], is equivariant but computes a matrix root on every iteration and is thus rarely used in practice.

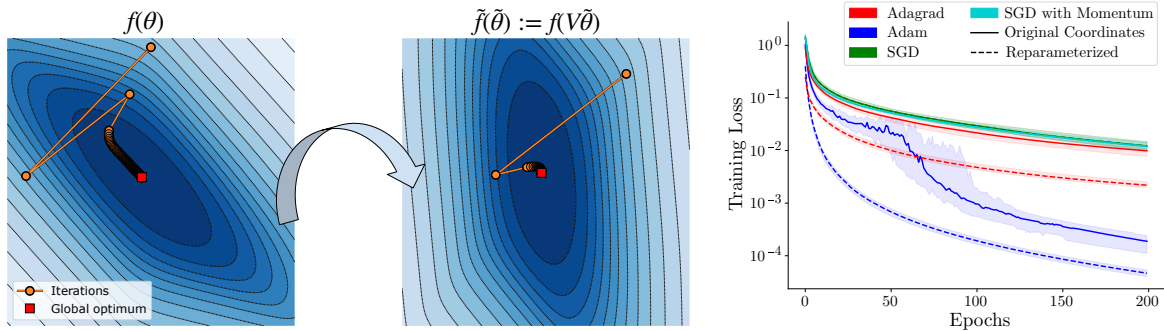


Figure 1: (Left) Using Adagrad in original coordinates and under EGOP reparameterization to optimize a two-dimensional log-sum-exp objective (defined in Section 6). In the EGOP eigenbasis, the primary directions of function variation are axis-aligned. Details in Section C.4. (Right) Cross-entropy loss from a 2-layer ReLU network in 2.4k dimensions trained on image classification using Adam, Adagrad, SGD, and SGD with momentum, in both original coordinates and under reparameterization. Equivariant methods (e.g. SGD) exhibit no change under reparameterization. Discussion and details in Section 6.

39 the sensitivity of adaptive algorithms to change-of-basis. We propose a reparameterization
 40 procedure based the *expected gradient outer product* (EGOP) matrix and show this reparamete-
 41 rization can improve convergence of adaptive methods. The geometric properties identified
 42 in this work—namely, strong decay of the EGOP eigenvalues—have been observed in a variety
 43 of machine learning objectives [29, 34, 42]. We include both empirical evidence and theoretical
 44 arguments suggesting that these properties arise when using natural data.

45 **Contributions.** We show that for a large class of objectives, the proposed reparameteri-
 46 zation procedure can improve the convergence of adaptive optimization algorithms. [AD :
 47 Added new sentence:] Our analysis of EGOP spectral decay and low-rank structure yields a
 48 novel and specific hypothesis for why adaptive algorithms are particularly sensitive to changes
 49 of basis in machine learning settings. Our main contributions are as follows: **(1)** We char-
 50 acterize a class of objective functions for which reparameterization can reduce the number of
 51 iterations required for adaptive algorithms to converge to minima and first-order stationary
 52 points. **(2)** For these functions, we identify a choice of basis in which adaptive algorithms
 53 will perform well, and we propose an approximation procedure that only requires access to a
 54 (stochastic) gradient oracle, rather than analytical knowledge about the loss function. This
 55 procedure is defined in Section 2. **(3)** We develop theory that proves that in neighborhoods
 56 of local minima, the proposed reparameterization endows adaptive algorithms with improved
 57 convergence guarantees, quantified in terms of the spectrum of the EGOP matrix. Our main
 58 results are discussed in Section 3. **(4)** We empirically investigate this procedure and find
 59 that the proposed reparameterization improves the convergence of adaptive algorithms on a
 60 diverse suite of optimization problems. These experiments include results with large-scale
 61 contemporary neural network architectures. We present empirical results in Section 6.

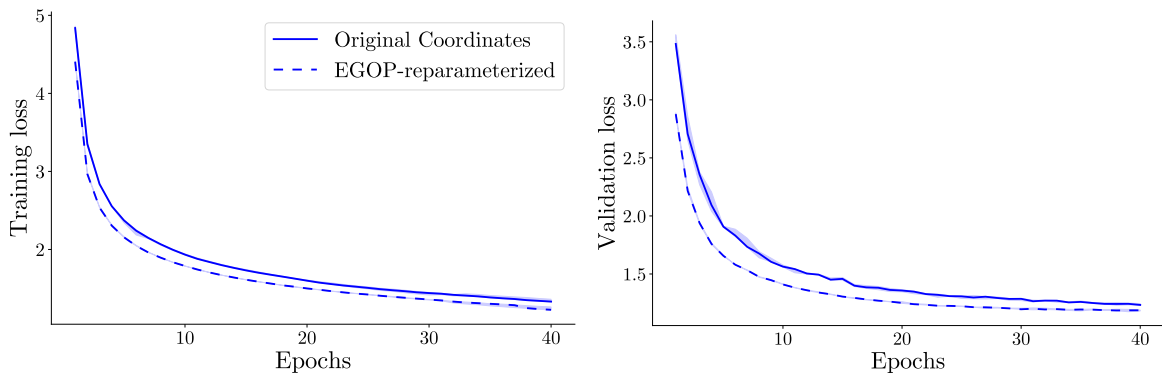


Figure 2: [AD : New figure.] EGOP reparameterization can scale to large, modern machine learning settings. Training and validation loss over epochs for a ResNet architecture trained to perform image classification on ImageNet. We compare optimizing the model with AdamW in original coordinates against AdamW under EGOP reparameterization. EGOP reparameterization improves convergence of both training and validation loss on this large-scale task. Results aggregated over five independent trials. Traces indicate medians, shaded regions indicate 25th-75th percentiles. Discussion and details in Section 6.

62 **1.1. Related Work.** Our work intersects with research on guarantees for adaptive algo-
 63 rithms, the role of orthogonal rotations in algorithmic performance, and geometry of machine
 64 learning objectives with natural data. Here we concisely survey related works, and we include
 65 an expanded discussion in Section F.

66 *Geometric Sensitivity of Adaptive Methods.* Recently, there has been renewed interest in dis-
 67 tinguishing the properties of adaptive algorithms versus stochastic gradient descent (SGD),
 68 arising in part from several empirical studies suggesting that adaptive methods outperform
 69 SGD when training transformer models [23, 42]. Traditional analyses of adaptive algorithms
 70 establish regret bounds for online convex optimization [14, 19], and more recent work estab-
 71 lishes convergence rates for smooth, non-convex objectives [11, 38]. However, because SGD is
 72 known to have optimal convergence rates in these settings, these theoretical results only show
 73 that adaptive algorithms achieve rates matching those of SGD.

74 In order to understand when adaptive algorithms enjoy provably stronger guarantees than
 75 SGD, recent work studies convergence under refined geometric assumptions, with particular
 76 emphasis on assumptions that are *not* rotationally invariant [21, 25, 40]. Xie et al. [40] estab-
 77 lish convergence guarantees in terms of the ℓ_∞ smoothness constant and show experimentally
 78 that rotationally invariant geometric assumptions do not suffice to capture settings when
 79 Adam outperforms SGD. Jiang et al. [21] and Liu et al. [25] study convergence of Adagrad on
 80 objectives that are *coordinate-wise smooth*, defined in Section 1.2. Our analysis builds on these
 81 results; Jiang et al. [21] and Liu et al. [25] show that the sum of the coordinate-wise smooth-
 82 ness constants governs Adagrad convergence, and we prove that EGOP reparameterization
 83 decreases this value by a factor as large as $1/d$.

84 *Change-of-Basis for Adaptive Algorithms.* Recent works propose that using different or-
 85 thonormal transformations with Adam and its variants can reduce computational costs and
 86 improve performance in neural networks [27, 37, 43]. Vyas et al. [37] proposed a method
 87 called SOAP, which computes an orthonormal reparameterization based on the singular vec-
 88 tors of the matrix-valued network gradients and performs optimization in this basis. Vyas
 89 et al. [37] empirically examine the performance of SOAP and find that it outperforms both
 90 Adam and Shampoo in LLM pre-training. Zhao et al. [43] propose GaLore, a method that
 91 simultaneously performs reparameterization and dimensionality reduction. GaLore computes
 92 a similar orthogonal basis to that used in SOAP, but instead of a full-dimensional change-of-
 93 basis GaLore retains only leading basis vectors in order to reduce dimension [43]. Maes et al.
 94 [27] empirically study Adam’s rotational sensitivity and propose an orthonormal reparamete-
 95 rization, similar to those used by SOAP and GaLore; they show empirically that this can
 96 improve Adam’s performance [27].

97 **[AD : New paragraph:]** EGOP reparameterization is related to SOAP, Shampoo, and Ga-
 98 Lore because all three methods use spectral information of some notion of gradient covariance,
 99 but the properties of the EGOP reparameterization proposed in this work are fundamentally
 100 distinct from those of Shampoo/SOAP/GaLore. As we discuss in Section 6, the change of
 101 basis leveraged by these methods is constrained to have a Kronecker product structure, while
 102 our approach does not impose this constraint; in Figure 8, we demonstrate empirically that
 103 as a result, EGOP reparameterization outperforms Shampoo and SOAP in settings where the
 104 Hessian is not well-approximated by a Kronecker product. Furthermore, as shown in Figure 9,
 105 we observe empirically that EGOP reparameterization can outperform SOAP and Shampoo
 106 even in the presence of exact Kronecker product structure, suggesting that EGOP spectral
 107 decay is a powerful tool even when additional structure is present.

108 Outside of training neural networks, several works have considered data-driven dimension-
 109 ality reduction methods for optimizing more general objectives with low-rank EGOP matrices
 110 [3, 7]. These procedures target objectives with exact low-rank structure, while our method can
 111 improve convergence of adaptive algorithms even when the EGOP matrix has strong spectral
 112 decay but is full-rank.

113 *EGOP Structure in Machine Learning.* Increasing empirical evidence across a wide range
 114 of applications, including language modeling and image classification, suggests that empirical
 115 loss functions used for training machine learning models are approximately low-rank—meaning
 116 these functions vary strongly in only a small subset of directions in parameter space [29, 34, 42].
 117 This approximate low-rank structure can be detected and analyzed using *active subspace*
 118 *methods*, which often leverage the *EGOP matrix*, defined in Section 2 [6]. Recently, there
 119 has been growing interest in the EGOP matrix in machine learning research [9, 28, 31, 44].
 120 Radhakrishnan et al. [31] provide theoretical and empirical evidence that the weights of neural
 121 networks trained with gradient descent correlate strongly with a kind of EGOP matrix.²

122 **1.2. Notation.** For a vector $\theta \in \mathbb{R}^d$, we denote its i^{th} entry by $\theta(i)$. We denote the inner
 123 product on \mathbb{R}^d by $\langle \cdot, \cdot \rangle$, and let $\|\cdot\|_p$ denote the vector p -norm on \mathbb{R}^d , with $\|\theta\|_\infty \stackrel{\text{def}}{=} \max_i |\theta(i)|$.

²In Mallinar et al. [28] Radhakrishnan et al. [31], and Zhu et al. [44], the EGOP is defined using the gradient with respect to the *input data* instead of the optimization parameters.

124 Given a matrix $A \in \mathbb{R}^{m \times n}$, we write $\|A\|_F$ for its *Frobenius* norm and $\|A\|_{\text{op}} \stackrel{\text{def}}{=} \sup_{\|\theta\|_2=1} \|A\theta\|_2$
 125 for its *operator* norm. Given a PSD matrix $H \in \mathbb{R}^{d \times d}$, we denote the norm $\|\theta\|_H \stackrel{\text{def}}{=} \sqrt{\langle \theta, H\theta \rangle}$.
 126 For a matrix $A \in \mathbb{R}^{n \times m}$, we denote the vectorization of A by $\text{vec}(A) \in \mathbb{R}^{nm}$.

127 We obtain guarantees in terms of the *coordinate-wise smoothness constants* of the objective
 128 $f(\cdot)$. Following Jiang et al. [21] and Liu et al. [25], we say that a function f is *coordinate-wise*
 129 *smooth* within a set $\Theta \subseteq \mathbb{R}^d$ with respect to constants $L_1, \dots, L_d > 0$ if $\forall \theta_1, \theta_2 \in \Theta$:

$$130 \quad (1.1) \quad |f(\theta_1) - f(\theta_2) - \langle \nabla f(\theta_2), \theta_1 - \theta_2 \rangle| \leq \frac{1}{2} \|\theta_1 - \theta_2\|_L^2,$$

131 where $L = \text{diag}(L_1, \dots, L_d)$.

132 **2. EGOP Reparameterization.** Given a function $f : \mathbb{R}^d \rightarrow \mathbb{R}$ and a sampling distribution
 133 ρ , the expected gradient outer product of $f(\cdot)$ with respect to ρ is defined as

$$134 \quad (2.1) \quad \text{EGOP}(f) \stackrel{\text{def}}{=} \mathbb{E}_{\theta \sim \rho} \left[\nabla f(\theta) \nabla f(\theta)^\top \right].$$

135 As $f : \mathbb{R}^d \rightarrow \mathbb{R}$, the EGOP is a $d \times d$ symmetric matrix, and it is positive semi-definite because it
 136 is an expectation over PSD matrices. We denote its eigendecomposition by $\text{EGOP}(f) = V\Lambda V^\top$
 137 where $V \in \mathbb{R}^{d \times d}$ is an orthonormal matrix. The EGOP eigenbasis captures key geometric qual-
 138 ities of the function $f(\cdot)$. When the sampling distribution ρ is isotropic, the leading eigen-
 139 vectors of the EGOP matrix capture the directions of greatest variation in $f(\cdot)$, whereas the
 140 eigenspaces of the smallest eigenvalues are directions along which $f(\cdot)$ does not vary strongly:
 141 this is reflected by the fact that for any eigenpair (λ_i, v_i) of $\text{EGOP}(f)$, $\lambda_i = \mathbb{E}_{\theta \sim \rho} [\langle \nabla f(\theta), v_i \rangle^2]$.

142 The EGOP matrix is defined with respect to a user-specified sampling distribution ρ . Our
 143 guarantees in Section 3 require that ρ be isotropic and that its scale is large enough with
 144 respect to the norm of some local minimum of $f(\cdot)$. In Section 6, we present empirical results
 145 when the EGOP is estimated with ρ a standard Gaussian, and when ρ is defined by common
 146 neural network initialization distributions [15].

147 In this work, we compare how adaptive optimization algorithms perform when optimizing
 148 $f(\cdot)$ versus the reparameterized function $\tilde{f} : \mathbb{R}^d \rightarrow \mathbb{R}$ defined $\tilde{f}(\tilde{\theta}) \stackrel{\text{def}}{=} f(V\tilde{\theta})$. For any f , the
 149 objective $\tilde{f}(\cdot)$ can be approximated by empirically estimating the EGOP via Monte Carlo
 150 sampling and forming the eigenvectors of the empirical EGOP matrix, as summarized in Al-
 151 gorithm 2.1. Note that as V is orthogonal, any solution $\tilde{\theta}$ obtained by optimizing $\tilde{f}(\cdot)$ can be
 152 transformed into a solution in original coordinates as $\theta \stackrel{\text{def}}{=} V\tilde{\theta}$, which satisfies $f(\theta) = \tilde{f}(\tilde{\theta})$.

Algorithm 2.1 Reparameterization by EGOP Eigenbasis

- 1: **input:** M number of gradient samples, distribution ρ .
 - 2: Generate $\{\theta_i\}_{i=1}^M \sim \rho$ i.i.d.
 - 3: Form empirical EGOP $\hat{P} = \frac{1}{M} \sum_{i=1}^M \nabla f(\theta_i) \nabla f(\theta_i)^\top$
 - 4: Form eigendecomposition $V\Lambda V^\top = \hat{P}$
 - 5: Define function $\tilde{f}(\cdot) = f \circ V$
 - 6: Optimize $\tilde{f}(\cdot)$ with adaptive algorithm of choice.
 - 7: **return** $\tilde{\theta} \in \mathbb{R}^d$ the result of optimizing $\tilde{f}(\cdot)$.
-

153 **3. Convergence Guarantees under EGOP Reparameterization.** In this section, we show
 154 that for objectives with strong EGOP spectral decay and Lipschitz Hessians, there is some
 155 neighborhood around each local minima such that EGOP reparameterization improves Ada-
 156 grad’s convergence guarantees within that neighborhood. Our results show that the radius
 157 of this neighborhood scales inversely with the Lipschitz constant of the Hessian. We analyze
 158 Adagrad’s convergence in both convex and nonconvex settings, focusing on the case of exact
 159 gradients for simplicity (see [Algorithm A.1](#) for a precise statement). In the convex setting, we
 160 study constrained Adagrad, which uses the following update rule:

$$161 \text{ (Adagrad)} \quad \theta_{t+1} = \Pi_{\Theta}^{H_t}(\theta_t - \eta H_t^{-1} \nabla f(\theta_t)), \quad H_t := \text{diag} \left(\epsilon^2 I_d + \sum_{j \leq t} \nabla f(\theta_j) \nabla f(\theta_j)^\top \right)^{1/2},$$

162 where Π_{Θ}^H denotes the orthogonal projection onto set Θ under the metric induced by $H \succ 0$.

163 Our results relate the improvement obtained through reparameterization by the EGOP
 164 eigenbasis to the *stable rank*³ of f :

$$165 \text{ (3.1)} \quad \text{sr}(f) \stackrel{\text{def}}{=} \frac{\sum_{i=1}^d \sqrt{\lambda_i(\text{EGOP}(f))}}{\sqrt{\lambda_{\max}(\text{EGOP}(f))}}.$$

166 Functions with strong EGOP spectral decay will have constant stable rank (tending towards
 167 1 as spectral decay increases), while functions without EGOP spectral decay will have stable
 168 rank scaling with d .

169 We now introduce the main assumptions for our theoretical results. Setting the stage,
 170 we consider a twice-differentiable objective $f : \mathbb{R}^d \rightarrow \mathbb{R}$, a sampling distribution $\rho(\cdot)$ and fix
 171 a local minimum θ^* . We also consider a set $\Theta \subseteq \mathbb{R}^d$, and denote the ℓ_p diameter of Θ by
 172 $D_p \stackrel{\text{def}}{=} \max_{\theta, \theta' \in \Theta} \|\theta - \theta'\|_p$.

173 **Assumption 1.** *The set Θ is convex, and the sampling distribution $\rho(\cdot)$ has support con-*
 174 *tained in Θ and is an isotropic distribution centered at some $\theta_c \in \Theta$ with scale c . This implies*
 175 $\mathbb{E}_{\rho}[\theta] = \theta_c$ and

$$176 \mathbb{E}_{\rho}[(\theta - \theta_c)(\theta - \theta_c)^\top] = c^2 \mathbb{I}.$$

177 *Moreover we assume Θ contains θ^* some local minimum of $f(\cdot)$ such that $\|\theta_c - \theta^*\|_2 \leq c$.*

178 Following the definitions introduced in Section 2, we let V denote the eigenbasis of the
 179 EGOP matrix and let $\tilde{f} \stackrel{\text{def}}{=} f \circ V$. Let $\tilde{\Theta}$ denote the corresponding transformation of set
 180 Θ : $\tilde{\Theta} \stackrel{\text{def}}{=} \{V^\top \theta \mid \theta \in \Theta\}$. Let $\{L_i\}_{i=1}^d$ denote the values for which $f(\cdot)$ is coordinate-wise
 181 smooth within Θ , following the definition in Eq. 1.1, and let $\{\tilde{L}_i\}_{i=1}^d$ denote the analogous
 182 coordinate-wise smoothness constants of $\tilde{f}(\cdot)$ within $\tilde{\Theta}$. We denote the vectors of coordinate-
 183 wise smoothness constants in original and reparameterized coordinates respectively by $\vec{L} \in \mathbb{R}^d$
 184 and $\vec{\tilde{L}} \in \mathbb{R}^d$.

³We call this quantity the stable rank because of the connection to the stable rank considered in numerical linear algebra; the ratio in (3.1) is related to $\|G\|_* / \|G\|_{\text{op}}$, where G denotes the empirical gradient bundle matrix $[\nabla f(\theta_1), \dots, \nabla f(\theta_M)] \in \mathbb{R}^{d \times M}$ and $\|\cdot\|_*$ denotes the nuclear norm. This ratio is often referred to in literature as the *stable* or *effective rank* of G [5, 33].

185 To measure the density of the eigenvectors of $\text{EGOP}(f)$, we introduce the values β_k : for
 186 v_k the k -th eigenvector of $\text{EGOP}(f)$, let $\beta_k \stackrel{\text{def}}{=} \|v_k\|_1^2/d$. Note that $\forall k \in [d]$, $\beta_k \in [1/d, 1]$ with
 187 $\beta_k = 1/d$ if v_k is one-hot, and $\beta_k = 1$ if v_k is uniformly dense.

188 Our second assumption requires that the Hessian of $f(\cdot)$ be Lipschitz within Θ and that
 189 its Lipschitz constant be suitably bounded.

190 **Assumption 2.** *The Hessian of $f(\cdot)$ is H -Lipschitz within Θ : $\|\nabla^2 f(\theta_1) - \nabla^2 f(\theta_2)\|_{\text{op}} \leq$
 191 $H\|\theta_1 - \theta_2\|_2 \forall \theta_1, \theta_2 \in \Theta$, and $\exists \delta \in [0, \beta_1^2)$ such that H satisfies*

$$192 \quad (3.2) \quad H \leq \frac{\sqrt{\delta \lambda_1(\text{EGOP}) + D_2^2 \lambda_1^2(\nabla^2 f(\theta^*))} - D_2 \lambda_1(\nabla^2 f(\theta^*))}{D_2^2}$$

193 **[AD : New paragraph:]** Assumption 2 quantifies how small the Lipschitz constant of the
 194 Hessian must be in order for the results in [Theorem 3.1](#) to apply, and its particular form
 195 derives from the analysis. As long as the function has similar curvature within the constraint
 196 set Θ (as quantified by [Eq. 3.2](#)), the proposed reparameterization is guaranteed to lead to
 197 an improvement. Critically, H is not an input to the algorithm and Adagrad in the new
 198 coordinate system is guaranteed to converge even when [Eq. 3.2](#) does not hold, albeit at a rate
 199 commensurate with that achieved in the original coordinate system.

200 Under the above assumptions, we show that EGOP spectral decay governs the improve-
 201 ments conferred by EGOP reparameterization. We first consider the convex setting. For
 202 simplicity, we instantiate these guarantees for constraint set Θ a ball; in practice, adaptive
 203 algorithms are often deployed with weight decay, which bounds the norm of the optimization
 204 parameters and effectively constrains θ to some ball. In [Section A.1.1](#), we state guarantees
 205 for generic convex constraint sets.

206 **Theorem 3.1 (Convergence for convex objectives with noise-free gradients).** *Consider $f(\cdot)$ a*
 207 *convex objective, constraint set Θ a ball, and sampling distribution $\rho(\cdot)$ satisfying [Assumptions](#)*
 208 *1 and 2, and $\tilde{f}(\cdot)$, $\tilde{\Theta}$ the EGOP-reparameterized objective and constraint set, respectively.*
 209 *Running [Algorithm A.1](#) with constrained updates and inputs $(\tilde{f}(\cdot), \nabla \tilde{f}(\cdot), \tilde{\theta}_0, T, \epsilon, \tilde{\Theta})$ for any*
 210 *initial condition $\tilde{\theta}_0 \in \tilde{\Theta}$ produces iterates $\{\tilde{\theta}_t\}_{t=1}^T$ satisfying*

$$211 \quad (3.3) \quad \frac{1}{T} \sum_{t=0}^{T-1} (\tilde{f}(\tilde{\theta}_t) - \tilde{f}(\tilde{\theta}^*)) = O \left(\left(\eta + \frac{D_2^2}{\eta} \right)^2 \frac{\|\vec{L}\|_1}{T} \left(\frac{\text{sr}(f)\sqrt{1+\delta}}{d(\beta_1 - \delta)} + \frac{\sqrt{\delta(1+\delta)}}{\beta_1 - \delta} \right) + \frac{\epsilon D_2^2}{\eta T} \right).$$

212 We compare the above bound with the convergence guarantee for Adagrad in original coordi-
 213 nates initialized at any $\theta_0 \in \Theta$:

$$214 \quad (3.4) \quad \frac{1}{T} \sum_{t=0}^{T-1} (f(\theta_t) - f(\theta^*)) \leq \left(\eta + \frac{D_2^2}{\eta} \right)^2 \cdot \frac{\|\vec{L}\|_1}{T} + \frac{\epsilon D_2^2}{\eta T}.$$

215 When $f(\cdot)$ has strong EGOP spectral decay and dense leading eigenvectors, [Theorem 3.1](#)
 216 implies reparameterization can significantly improve Adagrad's convergence bounds. In these
 217 settings, $\text{sr}(f)/(\beta_1 d) = O(1/d)$, implying stronger guarantees over original coordinates by up
 218 to a factor of d ; the term $\|\vec{L}\|_1$ appearing in both [Eq. 3.3](#) and [3.4](#) is the same, and denotes

219 the sum of the coordinate-wise smoothness constants of $f(\cdot)$ in original coordinates. Proof
 220 deferred to Section A.

221 In Section 4, we show that the two properties of the EGOP emphasized in this result,
 222 namely low stable rank and dense leading eigenvectors, are satisfied empirically for benchmark
 223 objectives. The factor $\beta = \|v_1\|_1^2/d$ tends towards 1 as the leading EGOP eigenvector gets
 224 denser. For the guarantee to reflect a benefit from reparameterization, it suffices to have
 225 $\beta \gg 1/d$ and small stable rank. We note that the density condition $\beta \gg 1/d$ is satisfied with
 226 high probability for random unit vectors in high dimensions: in particular, for $v \sim \text{Unif}(S^{d-1})$
 227 where S^{d-1} is the unit sphere in \mathbb{R}^d , with high probability $\|v\|_1^2/d > 0.6$. For simplicity,
 228 the above guarantee only considers the density of v_1 , but generalized guarantees in terms of
 229 β_k can be obtained. Similarly, one generalize [Theorem 3.1](#) for approximate versions of the
 230 EGOP eigenbasis; such guarantees will scale with the subspace distance between the EGOP
 231 eigenbasis and its approximation. We note that numerical stability parameter ϵ is typically
 232 chosen small enough that the second term, $\epsilon D_2^2/(\eta T)$, is not the dominant term.

233 Many naturally-motivated objectives in machine learning have locally Lipschitz Hessians,
 234 including loss functions used in logistic regression, over-parameterized matrix factorization,
 235 and training of multilayer linear networks; see [Section E.1.1](#) for examples.

236 In the non-convex setting, we study the convergence of unconstrained Adagrad ([Algo-](#)
 237 [rithm A.1](#) with unconstrained updates). We show there is some neighborhood around each
 238 local minima such that while iterates remain in this neighborhood, EGOP-reparameterized
 239 Adagrad enjoys local convergence bounds that are stronger than those in original coordinates
 240 by a factor of $\text{sr}(f)/d\beta_1$ (see [Theorem A.9](#)). Similarly to the convex setting, the radius of
 241 these neighborhoods grows as the Lipschitz constant of the Hessian decreases.

242 **4. EGOP Spectral Decay in Machine Learning.** Our analysis in [Section 3](#) shows EGOP
 243 spectral decay and dense leading EGOP eigenvectors are sufficient for EGOP reparameteriza-
 244 tion to improve convergence guarantees for Adagrad. Here we present empirical evidence
 245 that these conditions occur in benchmark machine learning objectives and discuss why natural
 246 data may produce EGOP spectral decay in real-world problems.

247 [Figure 3](#) shows the EGOP eigenspectrum for the objective function $f(\cdot)$ from an image
 248 classification problem. We use a 2-layer ReLU neural network to predict 10-class probabilities
 249 for handwritten digit images from the UCI ML digits dataset ([\[1\]](#)). Here $f(\cdot)$ denotes the
 250 negative log-likelihood loss on the training data. [Figure 3](#) shows strong EGOP spectral decay.
 251 We plot λ_k/λ_1 for all EGOP eigenvalues λ_k on a logarithmic scale, while the inset zooms in on
 252 the leading 100 indices on a linear scale. These plots illustrate roughly exponential eigenvalue
 253 decay. In [Section B.1](#), we show these trends are robust to the choice of sampling distribution
 254 (see [Figure 10](#)), that spectral decay persists when using *block* EGOP matrices, defined below
 255 in [Section 5](#), and that decay occurs in other benchmark datasets. In [Section B.1](#) we also
 256 examine the density of the leading EGOP eigenvectors for a 2-layer ReLU network on the
 257 UCI digits dataset and show that $\beta_k \gg 1/d$ for the leading eigenvectors (see [Figure 12](#)).

258 **4.1. Natural Data Induces EGOP Spectral Decay.** In addition to empirical evidence,
 259 simple gradient calculations suggest natural data may induce EGOP spectral decay in machine
 260 learning problems. Many common objectives can be expressed as $f(\theta) = h(A\theta)$, where $h(\cdot)$ is
 261 a loss function and $A \in \mathbb{R}^{n \times d}$ is a data matrix whose rows comprise samples $a_i \in \mathbb{R}^d$. By the

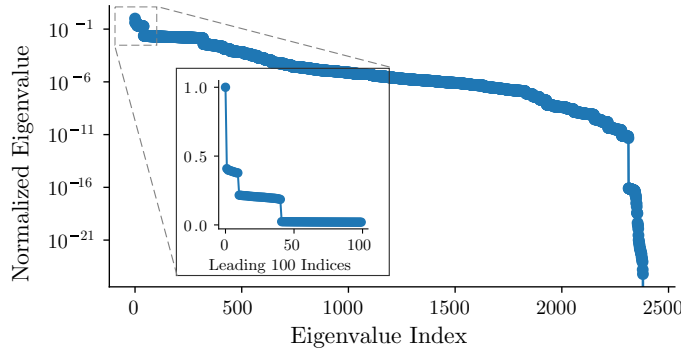


Figure 3: The EGOP eigenspectrum of a 2-layer ReLU network on the UCI handwritten digits dataset. Plot shows ratio λ_k/λ_1 as a function of eigenvalue index k , indexed in decreasing order. The EGOP eigenvalues decay sharply. In Figure 10, we demonstrate that these results are robust to the choice of sampling distribution ρ .

262 chain rule, the EGOP for such objectives satisfies

263 (4.1)
$$\text{EGOP}(f) = A^\top \mathbb{E}_{\theta \sim \rho} \left[\nabla_{\theta} h(A\theta) \nabla_{\theta} h(A\theta)^\top \right] A$$

264 where $\nabla_{\theta} h(A\theta)$ denotes the gradient $\nabla h(\cdot)$ evaluated at $A\theta$. This expression shows the EGOP
 265 is the transformation of some PSD matrix M by $A^\top M A$. It suggests that strong spectral decay
 266 in A may induce eigenvalue decay in the EGOP matrix. For many naturally occurring data
 267 distributions, the singular values of A exhibit strong decay [35]. The inner PSD matrix in the
 268 right hand side of Eq. 4.1 depends on both A and $h(\cdot)$, so without further assumptions on $h(\cdot)$
 269 it is difficult to precisely characterize the spectral decay induced by the composition with A^\top
 270 and A , but this can be quantified for specific choices of $h(\cdot)$; see Section A.2 for examples.

271 **5. Heuristics for Scalability.** Reparameterization with the EGOP eigenbasis incurs three
 272 main sources of additional computation: (1) sampling gradients to estimate the EGOP, (2)
 273 forming the EGOP eigenbasis, and (3) storing and applying the change-of-basis matrix to
 274 compute values and gradients of $f \circ V$. We outline some implementation details and heuristics
 275 that reduce the computational cost and enhance the scalability of the proposed framework.
 276 For a more detailed discussion, see Section D.

277 [AD : New paragraph] Empirically, we find that for functions with strong spectral decay,
 278 it suffices to accurately estimate only the leading EGOP eigenvectors. Based on this finding,
 279 one can use techniques like randomized SVD [18] to form $V_r \in \mathbb{R}^{d \times r}$, a matrix whose columns
 280 contain the estimated r leading eigenvectors of the EGOP. To enable optimization over the
 281 full parameter space, one can append a random orthogonal complementary basis to V_r , or
 282 introduce auxiliary variables restricted to the span of V_r^\perp . We refer to the latter approach as
 283 *auxiliary variable* EGOP reparameterization. We report empirical results using this heuristic
 284 in Figure 7, and we present full details and pseudocode in Section D.1.1.

285 For functions with strong spectral decay, a conservative number of gradient samples ($M \leq$
 286 d , where M is the number of samples and d is the number of problem parameters) is empirically

287 sufficient to yield change-of-basis matrices. In the presence of large EGOP spectral gaps,
 288 existing results formally establish that the number of gradient samples required to estimate
 289 the leading eigenspace grows logarithmically with dimension; see e.g., Corollary 3.8 in [6].

290 Structured or factorized approximations of V can also reduce the cost of storing and
 291 applying a change of basis. One can partition the variables in θ into subsets $\{S_i\}_{i=1}^L$ and
 292 perform *block reparameterization*. For each subset, we obtain a separate change-of-basis matrix
 293 $V^{(i)} \in \mathbb{R}^{|S_i| \times |S_i|}$ via the eigenvectors of the *block* EGOP matrix,

$$294 \quad \text{EGOP}^{(i)} \stackrel{\text{def}}{=} \frac{1}{M} \sum_{k=1}^M \nabla_{S_i} f(\theta_k) \nabla_{S_i} f(\theta_k)^\top$$

295 where $\nabla_{S_i} f(\theta_k) \in \mathbb{R}^{|S_i|}$ is the vector of partial derivatives of f w.r.t. the entries in S_i , and the
 296 points $\{\theta_k\}_{k=1}^M$ are sampled i.i.d. from ρ . Block reparameterization is well-suited to multilayer
 297 neural networks, where each layer forms a parameter block; such choices are well-motivated by
 298 empirical observations about the block-diagonal structure of the Hessian in neural networks
 299 [13]. Cost can be further reduced by only reparameterizing a subset of blocks; see Figure 14c
 300 for results employing this heuristic.

301 **[AD : New paragraph]** When employing block-reparameterization, a related cost-reducing
 302 heuristic is to share reparameterization matrices V across different blocks of compatible di-
 303 mensions; such a construction is well-motivated in convolutional layers, for example, where
 304 one can form blocks of parameters for each filter in a single layer. The key observation is
 305 that at initialization, the distribution of gradients from each filter is identical. Thus, the
 306 ground-truth EGOP matrix for each parameter block is identical, and so by estimating a
 307 single change-of-basis matrix V to be applied to all blocks in the convolutional layer, one can
 308 reduce both the number of forward/backward passes required to estimate the EGOP, as well
 309 as the storage cost of maintaining V ⁴. Figure 2 demonstrates this heuristic efficiently scales
 310 to large architectures.

311 **[AD : New paragraph]** One natural extension of Algorithm 2.1 is to consider periodically
 312 re-estimating the EGOP, using a modified sampling distribution ρ which has been re-centered
 313 at the latest optimization iterate, θ_t . This allows one to recompute a change-of-basis ma-
 314 trix based on local geometry, and may be well-suited to problems with highly heterogeneous
 315 structure. In Figure 4, we report experimental results using periodic EGOP reparameteriza-
 316 tion for a large-scale, non-convex optimization problem and find that it empirically improves
 317 convergence. However, we do find that throughout our experiments, the up-front, single-time
 318 reparameterization proposed in Algorithm 2.1 is sufficient to improve convergence, even for
 319 objectives where the Hessian is *not* Lipschitz. See Section 6 for details.

320 **6. Experimental Results.** We examine the impact of EGOP reparameterization in a va-
 321 riety of settings. We compare how the seminal adaptive algorithms Adagrad and Adam
 322 perform when optimizing functions in original coordinates versus under reparameterization

⁴We note that while such a construction is equivalent to employing a Kronecker-product reparameterization matrix V , the effective Kronecker product imposed by this construction is different from that employed by SOAP/Shampoo. Moreover, the effective EGOP change-of-basis matrix applied to each convolutional filter does *not* have any structural constraint; see Section D.1.2 for details.

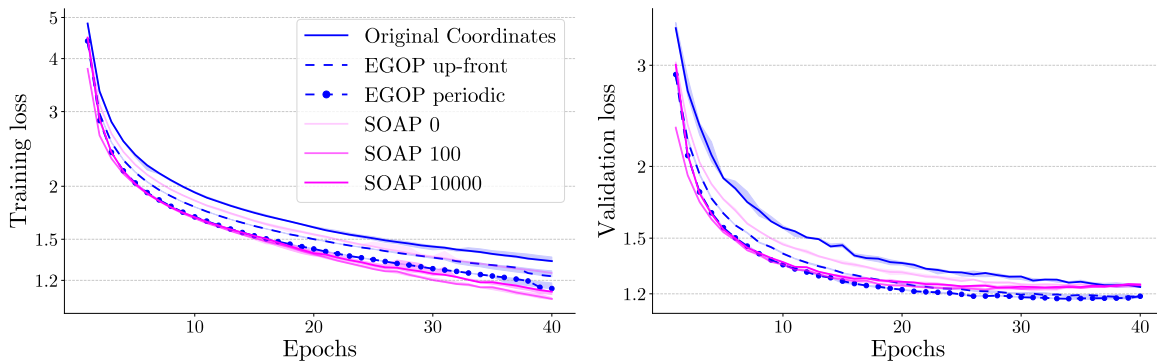


Figure 4: [AD : New figure] EGOP reparameterization can scale to large, modern machine learning settings. Training and validation loss over epochs for a ResNet architecture trained to perform image classification on ImageNet. We compare optimizing the model with Adam in original coordinates (blue, solid trace), under up-front EGOP reparameterization (blue, dashed trace), and with periodic EGOP reparameterization (blue, dashed trace with dots). We compare performance with SOAP under different update cadences; see discussion in the paragraph on *Large-scale Neural Networks* for descriptions of SOAP 0, SOAP 100, and SOAP 10000. EGOP reparameterization is competitive with SOAP at training on this large-scale task. Figure 5 demonstrates that EGOP reparameterization is competitive in wall-clock time as well. The results for Original Coordinates and EGOP up-front are identical to those in Figure 2, but in these plots we use logarithmic y-axis scaling for ease of comparing different traces.

323 by the empirical EGOP eigenbasis, both as described in Algorithm 2.1 and with the block
 324 reparameterization described in Section 5. We include comparisons to equivariant iterative
 325 optimization methods, namely (S)GD and (S)GD with momentum. When comparing meth-
 326 ods, we always choose equivalent initializations: when a method in original coordinates is
 327 initialized at θ_0 , chosen randomly, the reparameterized method is initialized at $V^T\theta_0$. Full
 328 experimental details are in Section C.

329 *Large-scale Neural Networks.* [AD : New section (through line 354)] In Figure 4, we pres-
 330 ent results using EGOP reparameterization to improve convergence in large-scale residual
 331 networks (ResNets). We emulate the well-tuned architecture established in He et al. [20]
 332 and consider the task of image classification using the ImageNet dataset [12]. We consider
 333 training with AdamW, an adaptive algorithm and variant of Adam which implements weight
 334 decay decoupled from the momentum computation [26]. We compare AdamW in original
 335 coordinates with EGOP-reparameterized AdamW using a single up-front computation of the
 336 EGOP reparameterization matrix, as proposed in Algorithm 2.1. We find that up-front EGOP
 337 reparameterization improves both minimization of the training and validation loss compared
 338 to AdamW in original coordinates. We also consider results using EGOP reparameterization
 339 with periodic re-computation of the EGOP eigenbasis, as discussed in Section 3, where the
 340 EGOP eigenbasis is re-estimated every epoch. We find that periodic reparameterization offers

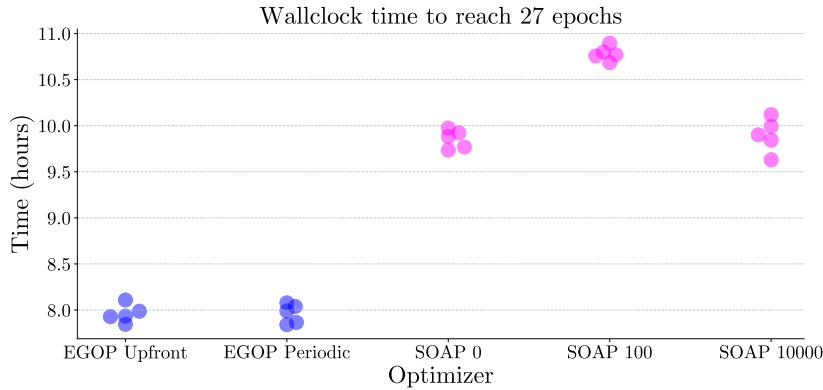


Figure 5: [AD : New figure] Wall-clock timing measures for training a residual network on the ImageNet dataset (counterpart to Figure 4). We report the wall-clock time required to train for 27 epochs, which is the number of epochs at which SOAP 100 minimizes its validation loss. As shown in Figure 4, at epoch 27 both EGOP methods achieve *lower* validation loss than SOAP. Both EGOP reparameterization methods achieve competitive wall-clock training times with SOAP and Shampoo. Horizontal jitter added for ease of visualization.

341 further benefits over up-front EGOP reparameterization in minimizing training loss, and that
 342 both variants achieve similar minimum validation loss. These results demonstrate the scala-
 343 bility and effectiveness of the heuristics discussed in Section 5, as well as demonstrating that
 344 EGOP reparameterization generalizes to a variety of modern neural network architectures.

345 We use this large-scale task to compare EGOP reparameterization with SOAP, which
 346 achieves state-of-the-art performance in training large neural networks [37]. We find that
 347 EGOP reparameterization is competitive both in performance and wall-clock time with this
 348 state-of-the-art method. In Figure 4, we present training and validation loss over epochs for
 349 both methods, and in Figure 5 we present wall-clock timing results. We find that EGOP-
 350 reparameterized AdamW achieves a lower minimum validation loss than SOAP, in addition
 351 to attaining competitive wall-clock time.

352 Our experiments include comparisons with SOAP at a variety of update cadences. In
 353 their original work proposing SOAP, Vyas et al. [37] reported that SOAP performance im-
 354 proves when the change-of-basis matrix is periodically updated, and they suggest an update
 355 cadence of every 10 to 100 batches. We include results with SOAP 0, which performs a single,
 356 up-front calculation of the change-of-basis matrix, comparable with up-front EGOP reparam-
 357 eterization; SOAP 100, which lies within the range of updates suggested by Vyas et al.; and
 358 SOAP 10k, which updates the change-of-basis matrix once per epoch, comparable with our
 359 implementation of periodic EGOP reparameterization.

360 *Linear Feedforward Networks.* We examine the impact of EGOP reparameterization on
 361 training 3-layer fully-connected linear feedforward networks using synthetic data. We consider
 362 parameters $\theta = [\text{vec}(W_1), \text{vec}(W_2), \text{vec}(W_3)]$ and train by minimizing loss function

$$363 \quad (6.1) \quad f(\theta) = \|W_3 W_2 W_1 A - Y\|_{\mathbb{F}}^2 / n_{\text{samples}}$$

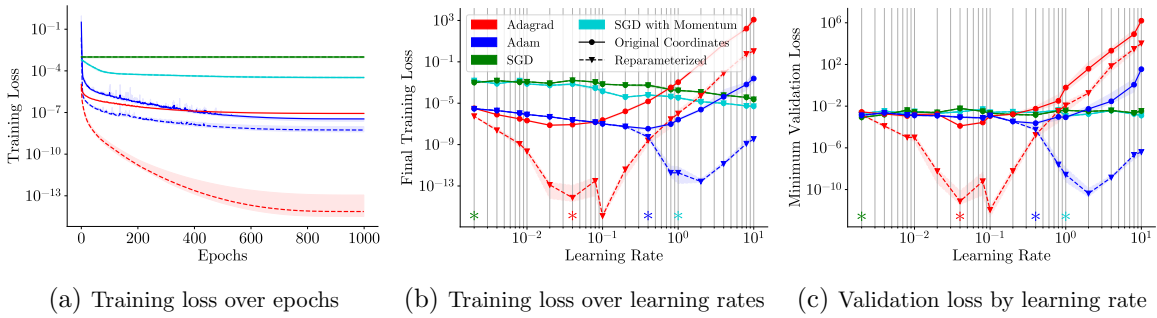


Figure 6: Training multilayer linear networks (6.1). Both SGD and SGD with momentum are equivariant optimization methods, so their results in original and reparameterized coordinates are exactly superimposed. In Figure 6c we consider the minimum validation loss achieved over epochs during training. Results are aggregated over 10 independent trials, with traces showing medians and shading indicating 25th-75th quartile. Asterisks indicate the learning rate used for each method in Figure 6a. Learning rates chosen to minimize validation loss of the algorithm in *original* coordinates.

364 where $A \in \mathbb{R}^{10 \times n_{\text{samples}}}$, and $Y = M^* A$ for $M^* \in \mathbb{R}^{10 \times 10}$ drawn from a standard Gaussian
 365 distribution. We use $W_1 \in \mathbb{R}^{50 \times 10}$, $W_2 \in \mathbb{R}^{30 \times 50}$, $W_3 \in \mathbb{R}^{10 \times 30}$. We induce spectral decay in
 366 $\text{EGOP}(f)$ by generating A with singular values $\sigma_k(A) = k^{-2}$ and random singular vectors. We
 367 use minibatched stochastic gradient samples throughout. Full details are deferred to section C.

368 In Figure 6 we study the impact of global EGOP reparameterization (Alg. 2.1). Figure 6a
 369 shows training loss over epochs in original coordinates and under reparameterization. Repa-
 370 rameterized Adagrad and reparameterized Adam achieve significantly lower final training loss
 371 and faster convergence than their counterparts in original coordinates. The adaptive methods
 372 also outperform the equivariant methods (SGD and SGD with momentum) in both coordinate
 373 settings. Figure 6b demonstrates the robustness of these results across learning rates. Aster-
 374 isks along the x-axis indicate the learning rates used in Figure 6a. Figure 6c confirms that
 375 the improved minimization of the training loss enabled by reparameterization does not lead
 376 to over-fitting; it shows that reparameterization enables adaptive methods to achieve better
 377 performance on validation data.

378 In Section B.2, we compare the above results to those obtained on the same task when
 379 employing several heuristics described in Section 5, including block reparameterization and
 380 reparameterizing only the first layer of the network. We find that both block-reparameterized
 381 Adagrad and Adam achieve lower final training loss than their counterparts in original co-
 382 ordinates. For Adagrad, reparameterizing only the first layer confers a benefit comparable
 383 to reparameterizing all layers, while for Adam reparameterizing the first layer alone provides
 384 only a marginal benefit over original coordinates.

385 **ReLU Networks for Image Classification.** We examine the impact of EGOP reparameter-
 386 ization when training networks to perform image classification using real-world data. We
 387 consider two benchmark image classification datasets: the UCI hand-written digits dataset

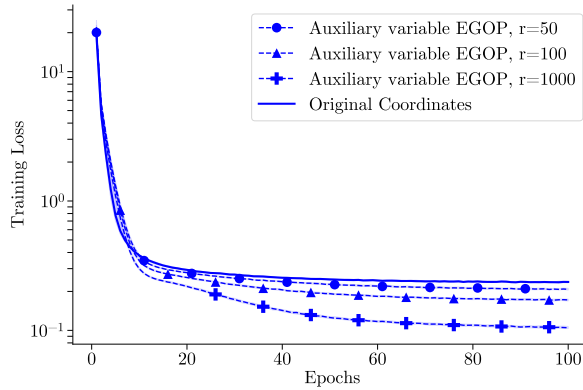


Figure 7: Scaling EGOP reparameterization to large fully-connected layers using auxiliary variable EGOP reparameterization, as described in Section 5. We perform reparameterization using varying numbers of leading EGOP eigenvectors, denoted by r . Even small values of r suffice to improve convergence, and the benefit increases as r increases. Results from training a ReLU network to classify the fashionMNIST dataset, in which the largest layer contains 78k weights. In Figure 16, we validate that this translates into improved wallclock time to convergence as well as improved accuracy on held-out data.

388 (8×8 pixel images), and the fashionMNIST dataset (28×28 pixel images of clothing) [1, 39].
 389 We train multilayer fully-connected ReLU networks to perform image classification on each
 390 dataset by minimizing the cross-entropy loss.

391 We perform block EGOP reparameterization, as described in Section 5. For the UCI
 392 digits, the number of gradient samples drawn equals the number of model weight param-
 393 eters. [AD : New text.] For the fashionMNIST results in Figure 7, we additionally employ
 394 heuristics described in Section 5. We perform auxiliary variable EGOP reparameterization,
 395 using randomized SVD to approximate the leading r EGOP eigenvectors, where $r \ll d$ the
 396 number of optimization variables. We use only a small number of stochastic minibatch gra-
 397 dients ($M = 0.01d$). These heuristics effectively scale EGOP reparameterization to large
 398 fully-connected layers. These heuristics also lead to faster convergence in terms of wall-clock
 399 time as reported in Figure 16.

400 Figure 1 (right) and Figure 7 plot training loss by epoch for the UCI digits and fash-
 401 ionMNIST respectively. On both datasets, reparameterized adaptive methods out-perform
 402 both their counterparts in original coordinates and the equivariant methods. In Section B.2
 403 we report results on hold-out data, showing reparameterized methods generalize as well as
 404 or better than original coordinates, and show that these results are robust to the choice of
 405 learning rate. For full experimental details, see Section C.

406 *Regression with Errors-in-Variables.* [AD : New section (through lines 432)] In this section,
 407 we empirically demonstrate one distinction between our proposed reparameterization and
 408 related methods of SOAP and Shampoo; EGOP reparameterization exploits approximate low-
 409 rank loss function geometry, while SOAP and Shampoo target Kronecker-product structure
 410 in the Hessian.

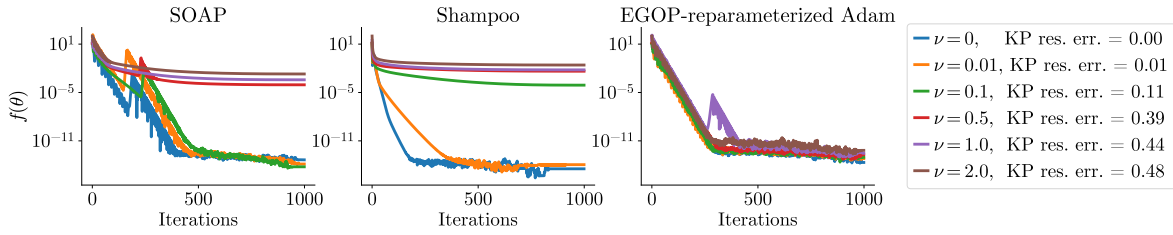


Figure 8: [AD : New figure] Comparing the impact of Hessian Kronecker-product (KP) structure on EGOP reparameterization, SOAP, and Shampoo, for objective $f(\theta)$ defined in (6.2). SOAP and Shampoo both exhibit worse convergence as the Hessian KP approximation worsens. We quantify Hessian KP structure via the KP residual error, defined in (6.3). EGOP-reparameterized Adam’s convergence is robust to this structure.

411 Shampoo and SOAP are designed for optimization problems where the variable has a
 412 natural matrix structure (e.g., a neural network weight matrix), and thus the gradient is also
 413 matrix-valued. Their change-of-basis derives from the eigenvectors of the Gram matrices of
 414 these gradient matrices. When mapped back to a vector representation (i.e., the vectorization
 415 of the matrix parameter and the setting of our EGOP-based approach), the corresponding
 416 SOAP/Shampoo change-of-basis matrix is restricted to a Kronecker product form. As a
 417 result, Shampoo and SOAP are hypothesized to perform best when the Hessian admits a close
 418 Kronecker-product approximation. Figure 8 presents evidence supporting this hypothesis.

419 We consider an objective that allows smooth interpolation between regimes when the
 420 Hessian admits an exact Kronecker-product structure, and cases where the Hessian is poorly
 421 approximated by such a factorization. For parameters $\theta = \text{vec}(W)$, $W \in \mathbb{R}^{n \times m}$, we consider

422 (6.2)
$$f(\theta) = \frac{1}{2} \|A(M \odot W) - Y\|_F^2.$$

423 Here, $M \in \mathbb{R}^{n \times m}$ is a fixed, known multiplicative noisy mask, $A \in \mathbb{R}^{n \times n}$ is the measurement
 424 map, and $Y = A(M \odot W^*)$ for some unknown $W^* \in \mathbb{R}^{n \times m}$. This models, for instance,
 425 regression problems based on data from sensors with well-characterized but potentially non-
 426 homogeneous gain.

427 Figure 8 displays $f(\cdot)$ by iteration for different instances of M . We draw the entries of
 428 M i.i.d. from $\mathcal{N}(1, \nu^2)$, where different problem instances use different ν . When $\nu = 0$, M
 429 is all ones, and the Hessian is exactly Kronecker-structured. As ν increases, the Kronecker
 430 structure of the Hessian degrades. We quantify this via Kronecker product residual error:

431 (6.3)
$$\text{KP residual}(f) = \min_{A \in \mathbb{R}^{n \times n}, B \in \mathbb{R}^{n \times n}} \|\nabla^2 f(\theta) - A \otimes B\|_F / \|\nabla^2 f(\theta)\|_F.$$

432 Figure 8 (left, center) shows that SOAP and Shampoo performance degrade as the Hessian
 433 becomes farther from Kronecker-product structured, while EGOP-reparameterized Adam Fig-
 434 ure 8 (right) maintains linear convergence across all problem instances. These experiments al-
 435 low SOAP and Shampoo to recompute their change-of-basis matrix on *every iteration*, whereas

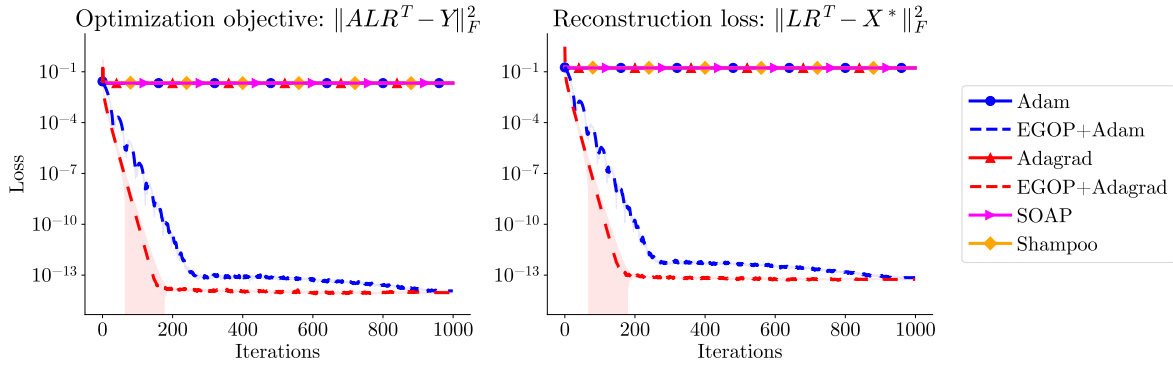


Figure 9: [AD : New figure] Comparing optimization algorithms for low-rank matrix factorization, as defined in (6.4). We note that base Adam, base Adagrad, Shampoo, and SOAP all have similar performance, and thus **their traces are superimposed**. On the right, we validate that the improved optimization exhibited by EGOP-reparameterized Adam and Adagrad selects a good minimizer, and does not overfit to noise. We emphasize that for this objective, the Hessian blocks $\nabla_{L,L}^2 f(\theta)$ and $\nabla_{R,R}^2 f(\theta)$ always admit exact Kronecker product factorizations, making this setting well-suited to the constraints imposed by SOAP and Shampoo.

436 for EGOP reparameterization we compute a single, up-front change-of-basis matrix following
 437 Algorithm 2.1. Even with this advantage, SOAP and Shampoo do not outperform EGOP-
 438 reparameterized Adam when the Hessian does not admit a good Kronecker approximation.

439 *Low-rank matrix factorization.* [AD : New section (through lines 451)] We complement
 440 the above results by showing EGOP reparameterization can offer advantages over SOAP
 441 and Shampoo even when the Hessian admits an exact Kronecker product factorization. We
 442 consider low-rank matrix factorization; for parameters $\theta = [\text{vec}(L), \text{vec}(R)]$ for $L \in \mathbb{R}^{n \times r}$, $R \in$
 443 $\mathbb{R}^{n \times r}$, we minimize the objective

$$444 \quad (6.4) \quad f(\theta) = \frac{1}{2} \|A(LR^T) - Y\|_F^2$$

445 where $Y = AX^* + \zeta$ is a matrix of observations with noise $\zeta_{i,j} \sim \mathcal{N}(0, 0.001)$, $X^* \in \mathbb{R}^{n \times n}$ is
 446 an unknown ground-truth matrix of (known) rank r , and $A \in \mathbb{R}^{n \times n}$ is a measurement map
 447 with singular value decay $\sigma_k(A) = k^{-1/2}$. We use standard spectral initialization.

448 In Figure 9, we plot the value of $f(\cdot)$ over iterations for varying algorithms, including
 449 SOAP and Shampoo, base Adam and Adagrad, and EGOP-reparameterized Adam and Ada-
 450 grad. EGOP-reparameterized Adam and Adagrad both achieve better minimization of $f(\cdot)$,
 451 and also yield lower reconstruction error, as quantified by $\|LR^T - X^*\|_F^2$. As before, we
 452 allow both SOAP and Shampoo to update their change-of-basis at every iteration, whereas
 453 EGOP reparameterization performs only a single up-front computation of the change-of-basis
 454 matrix. For this objective, the Hessian blocks $\nabla_{L,L}^2 f(\theta)$ and $\nabla_{R,R}^2 f(\theta)$ always admit exact
 455 Kronecker product decompositions, making this an objective well-suited to the constraints
 456 imposed by SOAP and Shampoo. Nonetheless, EGOP-reparameterized algorithms outper-
 457 form SOAP/Shampoo even in this favorable regime, indicating the need for further analysis

458 in order to predict when each approach will be optimal.

459 **7. Conclusions and Limitations.** In this work, we have shown through both analysis and
460 experiments that EGOP reparameterization can improve the convergence of adaptive algo-
461 rithms. There may be opportunities for future work to explore other reparameterizations and
462 additional hyperparameter tuning, which we have not examined in this work. One limitation of
463 this work is the engineering required for scalability, as discussed in Section 5. Fully character-
464 izing the trade-off between improved convergence and the up-front cost of reparameterization
465 remains a direction for future research. [AD : New text (through end of conclusion)] Another
466 interesting direction for future work would be studying whether the guarantees in Theorem 3.1
467 can be strengthened under a periodic EGOP reparameterization scheme, as described in Sec-
468 tion 5. Our experiments in Figure 4 demonstrate empirically that such periodic re-estimation
469 of the EGOP eigenbasis can offer further improvements over the up-front reparameterization
470 analyzed in this work.

471 Further work examining the tradeoffs between EGOP reparameterization and the repa-
472 rameterizations employed by SOAP and Shampoo would benefit practitioners by providing
473 greater guidance on when each method may perform best. Our empirical results in Figure 9
474 demonstrate that EGOP reparameterization can outperform SOAP and Shampoo even in the
475 presence of exact Kronecker product structure, suggesting that approximate low-rank geome-
476 try, as quantified by EGOP spectral decay, is a powerful tool even when additional geometric
477 structure is present.

478 **Appendix A. Proofs.** We briefly introduce additional notation used in our analysis. We
 479 denote the unit sphere in d dimensions by \mathbb{S}^{d-1} and the Euclidean ball of radius r centered
 480 at \bar{x} by $\mathcal{B}(\bar{x}; r)$. For matrices $A, B \in \mathbb{R}^{d \times d}$, we denote the Löwner ordering by $A \preceq B$. We
 481 write $A \lesssim B$ to indicate the existence of a dimension-independent positive constant $c > 0$
 482 such that $A \leq cB$. We write $\mathcal{O}(d, r) := \{X \in \mathbb{R}^{d \times r} \mid X^\top X = I_r\}$ and $\mathcal{O}(d) \equiv \mathcal{O}(d, d)$ for the
 483 set of matrices with orthogonal columns. We write $\langle X, Y \rangle \stackrel{\text{def}}{=} \text{tr}(X^\top Y)$ for the Euclidean inner
 484 product and $\|X\| = \sqrt{\langle X, X \rangle}$ for its induced norm.] We write $\mathcal{O}(d)$ for the set of matrices
 485 with orthogonal columns. Given a set of scalars $\{L_i\}_{i=1}^d$, we let $\text{diag}(L_1, \dots, L_d)$ denote the
 486 $d \times d$ diagonal matrix with values L_1, \dots, L_d along the diagonal. We let $\text{mat}(\theta) \in \mathbb{R}^{m \times n}$
 487 denote the reshaping of θ into a matrix for some m and n such that $mn = d$, and similarly
 488 we let $\text{vec}(M)$ denote the reshaping of a matrix into a vector.

489 Following the definitions introduced in Section 2, we let V denote the eigenbasis of the
 490 EGOP matrix and let $\tilde{f} \stackrel{\text{def}}{=} f \circ V$. Let $\tilde{\Theta}$ denote the corresponding transformation of set
 491 Θ : $\tilde{\Theta} \stackrel{\text{def}}{=} \{V^\top \theta \mid \theta \in \Theta\}$. We denote the corresponding diameter measurements as $\tilde{D}_p \stackrel{\text{def}}{=} \max_{\tilde{\theta}_1, \tilde{\theta}_2 \in \tilde{\Theta}} \|\tilde{\theta}_1 - \tilde{\theta}_2\|_p$. Let $\{L_i\}_{i=1}^d$ denote the values for which $f(\cdot)$ is coordinate-wise smooth
 492 within Θ , following Eq. 1.1, and let $\{\tilde{L}_i\}_{i=1}^d$ denote the analogous coordinate-wise smoothness
 493 constants of $\tilde{f}(\cdot)$ within $\tilde{\Theta}$. We denote the vectors of coordinate-wise smoothness constants in
 494 original and reparameterized coordinates respectively by $\vec{L} \in \mathbb{R}^d$ and $\vec{\tilde{L}} \in \mathbb{R}^d$.

495 We analyze the convergence of Adagrad in both convex and nonconvex settings. For
 496 completeness, we recall the full Adagrad algorithm in Algorithm A.1.
 497

Algorithm A.1 Adagrad($f, g, \theta_0, T, \epsilon : \text{default} = 10^{-8}, \Theta : \text{Optional}$)

1: **Input:** Objective f , oracle gradient g , $\theta_0 \in \mathbb{R}^d$, η step size, T number of iterations, ϵ
 numerical stability parameter, $\Theta \subseteq \mathbb{R}^d$ constraint set.
 2: $v_0 \leftarrow \epsilon^2 \mathbb{I}_d$
 3: **for** $t \in [1, T]$ **do**
 4: $g_t \leftarrow g(f, \theta_t)$
 5: $v_t(i) \leftarrow v_{t-1}(i) + g_t(i)^2 \forall i \in [d]$
 6: $H_t \leftarrow \text{diag}(v_t)^{1/2}$
 7: $w_t \leftarrow \theta_t - \eta H_t^{-1} g_t$
 8: **Constrained update:** $\theta_{t+1} = \Pi_{\Theta}^{H_t}(w_t)$
 9: **Unconstrained update:** $\theta_{t+1} = w_t$
 10: **end for**
 11: **Output:** $\theta_T \in \mathbb{R}^d$

498 **A.1. Proofs from Section 3.** Several recent works have proved convergence guarantees
 499 for Adagrad in terms of $\|\vec{L}\|_1$, showing that Adagrad enjoys better performance guarantees in
 500 settings when $\|\vec{L}\|_1$ is small [21, 25, 40]. Our main results establish that when $\text{EGOP}(f)$ has
 501 strong spectral decay and dense leading eigenvectors, then reparameterization can produce
 502 $\|\vec{\tilde{L}}\|_1$ significantly smaller than $\|\vec{L}\|_1$, implying better convergence bounds.

503 **Theorem A.1.** *If $\exists \delta \in [0, \beta_1^2)$ such that the Hessian of $f(\cdot)$ is H -Lipschitz within Θ for H*

504 *satisfying*

$$505 \quad (\text{A.1}) \quad H \leq \frac{\sqrt{\delta \lambda_1(\text{EGOP}) + D_2^2 \lambda_1^2(\nabla^2 f(\theta^*))} - D_2 \lambda_1(\nabla^2 f(\theta^*))}{D_2^2}$$

506 *Then within Θ and $\tilde{\Theta}$ respectively, the smoothness constants of $f(\cdot)$ and $\tilde{f}(\cdot)$ satisfy*

$$507 \quad \frac{\|\tilde{L}\|_1}{\|\tilde{L}\|_1} = O\left(\sqrt{1+\delta} \left(\frac{\text{sr}(f)}{d(\beta_1 - \delta)} + \frac{\sqrt{\delta}}{\beta_1 - \delta}\right)\right)$$

508 *where $\text{sr}(f)$ denotes the stable rank of $f(\cdot)$.*

509 **Thm. A.1** implies that when the EGOP leading eigenvectors have constant density, i.e. when
 510 $\beta_1 \gg 1/d$, then for small δ the ratio $\|\tilde{L}\|_1/\|\tilde{L}\|_1$ scales as $\text{sr}(f)/d\beta_1$. We note that the
 511 density condition $\beta \gg 1/d$ is satisfied with high probability for random unit vectors in high
 512 dimensions: in particular, for $v \sim \text{Unif}(S^{d-1})$ where S^{d-1} is the unit sphere in \mathbb{R}^d , with high
 513 probability $\|v\|_1^2/d \approx 2/\pi > 0.6$.

514 We prove **Theorem A.1** by establishing the following pair of claims. The first lower bounds
 515 the sum of the coordinate-wise smoothness constants of $f(\cdot)$ in original coordinates:

516 **Lemma A.2.** *Consider $f(\cdot)$, Θ , ρ satisfying Assumptions 1 and 2 and consider the set of*
 517 *dense unit vectors $\nu \in \{\pm d^{-1/2}\}^d$, the collection of vectors whose entries all have magni-*
 518 *tude $|\nu(i)| = d^{1/2}$. Then within Θ , $f(\cdot)$ is coordinate-wise smooth with respect to constants*
 519 *satisfying*

$$520 \quad \|\tilde{L}\|_1 \geq \frac{d}{2c} \cdot \max_{\nu \in \{\pm d^{-1/2}\}^d} \frac{\langle \nu, \text{EGOP } \nu \rangle - \gamma}{\sqrt{\lambda_{\max}(\text{EGOP}) + \gamma}}$$

521 *where*

$$522 \quad (\text{A.2}) \quad \gamma \stackrel{\text{def}}{=} 2H\lambda_{\max}(\nabla^2 f(\theta^*))M_3 + H^2M_4$$

523 *where H is the Lipschitz constant of the Hessian of $f(\cdot)$ in Assumption 2 and $M_p \stackrel{\text{def}}{=} \mathbb{E}_{\theta \sim \rho}[\|\theta -$
 524 $\theta^*\|_2^p]$.*

525 If the EGOP has dense leading eigenvectors, then the term $\langle \nu, \text{EGOP } \nu \rangle$ is large. In particular,
 526 **Lemma A.2** implies the following:

527 **Corollary A.3.** *For $f(\cdot)$, Θ, ρ satisfying Assumptions 1 and 2, the smoothness constants of*
 528 *f satisfy*

$$529 \quad \|\tilde{L}\|_1 \geq \frac{d}{2c} \frac{\beta_k \lambda_k(\text{EGOP}) - \gamma}{\sqrt{\lambda_{\max}(\text{EGOP}) + \gamma}}.$$

530 *where $\beta_k \stackrel{\text{def}}{=} \|v_k\|_1^2/d$ for $\lambda_k(\text{EGOP}), v_k$ the k^{th} eigenvalue and eigenvector of $\text{EGOP}(f)$, and*
 531 *γ defined in Eq. A.2.*

532 In contrast, we show that under reparameterization by the EGOP eigenbasis, the smooth-
 533 ness constants can be upper bounded by the following:

534 **Lemma A.4.** Let V be the eigenbasis of the EGOP of $f(\cdot)$ with respect to ρ , and define
 535 $\tilde{f}(\tilde{\theta}) \stackrel{\text{def}}{=} f(V\theta)$. Let the function $f(\cdot)$ satisfy Assumptions 1 and 2. Then within $\tilde{\Theta}$, $\tilde{f}(\cdot)$
 536 satisfies Eq. 1.1 with respect to smoothness coordinates whose sum is bounded by

$$537 \quad \|\tilde{\mathbf{L}}\|_1 \leq d \left(\frac{\sqrt{\gamma}}{c} + HD_2 \right) + \frac{1}{c} \sum_{i=1}^d \sqrt{\lambda_i}$$

538 where γ is defined in Eq. A.2. and λ_i denotes the i th eigenvalue of the EGOP of f w.r.t. ρ .

539 Note that as Lipschitz constant of Hessian in Assumption 2 goes to zero, so does the value
 540 of γ in both Lemmas A.3 and A.4. Theorem A.1 follows from Lemmas A.3 and A.4.

541 **Dense EGOP eigenvectors lead to large $\|\tilde{\mathbf{L}}\|_1$.** We defer the full proof of Lemma A.3 to
 542 Section E.1. We note here the main intermediate results used to establish Lemma A.3. The
 543 first intermediate result gives an alternative characterization of the coordinate wise smoothness
 544 constants defined in Eq. 1.1.

545 **Lemma A.5.** A twice-differentiable function $f : \mathbb{R}^d \rightarrow \mathbb{R}$ satisfies Eq. 1.1 with respect to
 546 smoothness constants L_1, \dots, L_d within Θ if and only if $\forall \theta \in \Theta, \forall v \in \mathbb{R}^d$,

$$547 \quad |\langle v, \nabla^2 f(\theta)v \rangle| \leq \langle v, \text{diag}(L)v \rangle$$

548 where $\text{diag}(L) \in \mathbb{R}^{d \times d}$ is the diagonal matrix with diagonal entries L_i .

549 The second intermediate result relates the EGOP to the Hessian.

550 **Lemma A.6.** Consider twice-differentiable function $f : \mathbb{R}^d \rightarrow \mathbb{R}$ and sampling distribution
 551 ρ satisfying Assumptions 2, and 1 with respect to local minimum θ^* . Then the EGOP of $f(\cdot)$
 552 with respect to ρ satisfies

$$553 \quad \mathbb{E}_{\theta \sim \rho}[\nabla f(\theta)\nabla f(\theta)^\top] = G(\theta^*) + E_f(\theta^*)$$

554 where $G(\theta^*)$ is a PSD matrix satisfying

$$555 \quad c^2 \nabla^2 f(\theta^*) \nabla^2 f(\theta^*)^\top \preceq G(\theta^*) \preceq 2c^2 \nabla^2 f(\theta^*) \nabla^2 f(\theta^*)^\top$$

556 and the matrix $E_f(\theta^*)$ satisfies

$$557 \quad |\langle v, E_f(\theta^*)v \rangle| \leq (2H\lambda_{\max}(\nabla^2 f(\theta^*))M_3 + H^2M_4) \|v\|_2^2$$

558 where

$$559 \quad M_p \stackrel{\text{def}}{=} \mathbb{E}_{\theta \sim \rho}[\|\theta - \theta^*\|_2^p].$$

560 We defer the proof of these intermediate results to Section E.1.

561 **Reparameterization by EGOP eigenbasis produces small $\|\tilde{\mathbf{L}}\|_1$.** We defer the full proof of
 562 Lemma A.4 to Section E. Here we present results establishing the equivariance of the EGOP
 563 matrix, which is the key step in proving Lemma A.4. Recall that Assumption 1 implies ρ is
 564 isotropic about $\theta_c \in \Theta$. It is convenient to express ρ as

$$565 \quad \rho(\theta) = \rho_0(\theta - \theta_c)$$

566 for ρ_0 mean-zero and isotropic (about the origin). For reparameterized set

$$567 \quad \tilde{\Theta} \stackrel{\text{def}}{=} \{V^\top \theta \mid \theta \in \Theta\}$$

568 and element $V^\top \theta_c \in \tilde{\Theta}$, we consider the reparameterized sampling distribution

$$569 \quad (\text{A.3}) \quad \tilde{\rho}(\tilde{\theta}) \stackrel{\text{def}}{=} \rho(\tilde{\theta} - V^\top \theta_c).$$

570 The distribution $\tilde{\rho}$ is then isotropic about $V^\top \theta_c$. Moreover, because ρ_0 is isotropic,

$$571 \quad \tilde{\rho}(\tilde{\theta}) = \rho_0(\tilde{\theta} - V^\top \theta_c) = \rho_0(V(\tilde{\theta} - V^\top \theta_c)) = \rho_0(V\tilde{\theta} - \theta_c) = \rho(V\tilde{\theta})$$

572 and similarly $\rho(\theta) = \tilde{\rho}(V^\top \theta)$.

573 **Lemma A.7.** *Given $f(\cdot)$ whose EGOP with respect to isometric distribution ρ has eigen-*
574 *value decomposition*

$$575 \quad \mathbb{E}_{\theta \sim \rho}[\nabla f(\theta) \nabla f(\theta)^\top] = V \Lambda V^\top,$$

576 *the EGOP of the reparameterized function $\tilde{f}(\theta) \stackrel{\text{def}}{=} f(V\theta)$ with respect to reparameterized*
577 *sampling distribution $\tilde{\rho}$ (defined in Eq. A.3), we have*

$$578 \quad \mathbb{E}_{\theta \sim \tilde{\rho}}[\nabla \tilde{f}(\theta) \nabla \tilde{f}(\theta)^\top] = \Lambda$$

579 We defer the proof of Lemma A.7 to Section E.

580 *Proof of Theorem A.1.* Combining Lemmas A.3 and A.4 implies Theorem A.1:

581 *Proof of Theorem A.1.* By Lemma A.3,

$$582 \quad \|\vec{L}\|_1 \geq \max_k \frac{d \beta_k \lambda_k(\text{EGOP}(f)) - \gamma}{2c \sqrt{\lambda_1(\text{EGOP}(f)) + \gamma}} \geq \frac{d \beta_1 \lambda_1(\text{EGOP}(f)) - \gamma}{2c \sqrt{\lambda_1(\text{EGOP}(f)) + \gamma}}$$

583 for

$$584 \quad \gamma \stackrel{\text{def}}{=} 2H \lambda_{\max}(\nabla^2 f(\theta^*)) M_3 + H^2 M_4.$$

585 Recall that

$$586 \quad M_p \stackrel{\text{def}}{=} \mathbb{E}_{\theta \sim \rho}[\|\theta - \theta^*\|_2^p]$$

587 and thus, because we assume $\text{supp}(\rho) \subseteq \Theta$, $M_p \leq D_2^p$ where D_2 is the ℓ_2 diameter of Θ .

588 Additionally, Lemma A.4 states that

$$589 \quad \|\vec{L}\|_1 \leq d \left(\frac{\sqrt{\gamma}}{c} + H D_2 \right) + \frac{1}{c} \sum_{i=1}^d \sqrt{\lambda_i(\text{EGOP}(f))}$$

590 where γ is the same as defined above.

591 Combining these two results implies

$$592 \quad \frac{\|\vec{L}\|_1}{\|\vec{L}\|_1} \leq \left(d \left(\frac{\sqrt{\gamma}}{c} + H D_2 \right) + \frac{1}{c} \sum_{i=1}^d \sqrt{\lambda_i} \right) \left(\frac{d \beta_1 \lambda_1 - \gamma}{2c \sqrt{\lambda_1 + \gamma}} \right)^{-1}$$

593 where, for conciseness, we let $\lambda_i \stackrel{\text{def}}{=} \lambda_i(\text{EGOP}(f))$ the i th eigenvalue of $\text{EGOP}(f)$, indexed in
 594 decreasing order. Expanding the above yields

$$595 \quad (\text{A.4}) \quad \frac{\|\vec{L}\|_1}{\|\vec{L}\|_1} \leq d \left(\frac{\sqrt{\gamma}}{c} + HD_2 \right) \cdot \frac{2c \sqrt{\lambda_1 + \gamma}}{d \beta_1 \lambda_1 - \gamma} + \frac{1}{c} \sum_{i=1}^d \sqrt{\lambda_i} \cdot \frac{2c \sqrt{\lambda_1 + \gamma}}{d \beta_1 \lambda_1 - \gamma}$$

$$596 \quad (\text{A.5}) \quad = 2(\sqrt{\gamma} + HcD_2) \cdot \frac{\sqrt{\lambda_1 + \gamma}}{\beta_1 \lambda_1 - \gamma} + 2 \sum_{i=1}^d \sqrt{\lambda_i} \cdot \frac{1 \sqrt{\lambda_1 + \gamma}}{d \beta_1 \lambda_1 - \gamma}.$$

597 Because $M_p \leq D_2^p$, by definition γ satisfies

$$598 \quad \gamma \leq 2H\lambda_{\max}(\nabla^2 f(\theta^*))D_2^3 + H^2D_2^4.$$

599 Thus

$$600 \quad H \leq \frac{\sqrt{\delta\lambda_1(\text{EGOP}) + D_2^2\lambda_1^2(\nabla^2 f(\theta^*))} - D_2\lambda_1(\nabla^2 f(\theta^*))}{D_2^2},$$

601 implies $\gamma \leq \delta\lambda_1$ for $\lambda_1 \stackrel{\text{def}}{=} \lambda_1(\text{EGOP}(f))$, the largest EGOP eigenvalue. Additionally, H
 602 admits a simplified upper bound:

$$603 \quad H \leq \frac{\sqrt{\delta\lambda_1(\text{EGOP}) + D_2^2\lambda_1^2(\nabla^2 f(\theta^*))} - D_2\lambda_1(\nabla^2 f(\theta^*))}{D_2^2} \leq \frac{\sqrt{\delta\lambda_1}}{D_2^2}.$$

604 Using these bounds on H and γ , the first term in Eq. A.5 can be bounded as

$$605 \quad 2(\sqrt{\gamma} + HcD_2) \cdot \frac{\sqrt{\lambda_1 + \gamma}}{\beta_1 \lambda_1 - \gamma} \leq 2 \left(\sqrt{\delta\lambda_1} + \frac{\sqrt{\delta\lambda_1}}{D_2^2} \cdot cD_2 \right) \cdot \frac{\sqrt{1 + \delta}}{\beta_1 - \delta} \cdot \frac{1}{\sqrt{\lambda_1}}$$

$$606 \quad = 2 \left(1 + \frac{cD_2}{D_2^2} \right) \frac{\sqrt{\delta(1 + \delta)}}{\beta_1 - \delta}$$

$$607 \quad \leq 6 \frac{\sqrt{\delta(1 + \delta)}}{\beta_1 - \delta}$$

608 where the last inequality follows from the fact that because $\text{supp}(\rho) \subseteq \Theta$, $c \leq D_2$.

609 The second term in Eq. A.5 can be bounded as

$$610 \quad 2 \sum_{i=1}^d \sqrt{\lambda_i} \cdot \frac{1 \sqrt{\lambda_1 + \gamma}}{d \beta_1 \lambda_1 - \gamma} \leq 2 \sum_{i=1}^d \sqrt{\lambda_i} \cdot \frac{1 \sqrt{1 + \delta}}{d(\beta_1 - \delta)\sqrt{\lambda_1}} = \frac{2 \text{sr}(f)}{d(\beta_1 - \delta)} \sqrt{1 + \delta}$$

611 where the equality follows from the definition of the stable rank. Combining these bounds
 612 yields the result. ■

613 **A.1.1. Convergence guarantees in the convex setting.** Here we prove [Theorem 3.1](#) by
 614 proving a more general statement in terms of any convex constraint set Θ .

615 **Lemma A.8 (Convergence for convex objectives with noise-free gradients).** Consider $f(\cdot)$
 616 a convex objective, constraint set Θ , and sampling distribution $\rho(\cdot)$ satisfying Assumptions

617 **1 and 2.** Let $\tilde{f}(\cdot)$ and $\tilde{\Theta}$ denote the reparameterized objective and constraint set respectively
 618 corresponding to the eigenbasis of EGOP(f). Then running Algorithm A.1 with constrained
 619 updates and with inputs $(\tilde{f}(\cdot), \nabla \tilde{f}(\cdot), \tilde{\theta}_0, T, \epsilon, \tilde{\Theta})$ for any initial condition $\tilde{\theta}_0 \in \tilde{\Theta}$ produces
 620 iterates $\{\tilde{\theta}_t\}_{t=1}^T$ satisfying

(A.6)

$$621 \quad \frac{1}{T} \sum_{t=0}^{T-1} (\tilde{f}(\tilde{\theta}_t) - \tilde{f}(\tilde{\theta}^*)) = O \left(\left(\eta + \frac{\tilde{D}_\infty^2}{\eta} \right)^2 \frac{\|\vec{L}\|_1}{T} \left(\frac{\text{sr}(f)\sqrt{1+\delta}}{d(\beta_1 - \delta)} + \frac{\sqrt{\delta(1+\delta)}}{\beta_1 - \delta} \right) + \frac{\epsilon D_2^2}{\eta T} \right).$$

622 Setting $\eta = \tilde{D}_\infty$ to minimize this upper bound yields

$$623 \quad \frac{1}{T} \sum_{t=0}^{T-1} (\tilde{f}(\tilde{\theta}_t) - \tilde{f}(\tilde{\theta}^*)) = O \left(\tilde{D}_\infty^2 \frac{\|\vec{L}\|_1}{T} \left(\frac{\text{sr}(f)\sqrt{1+\delta}}{d(\beta_1 - \delta)} + \frac{\sqrt{\delta(1+\delta)}}{\beta_1 - \delta} \right) + \frac{\epsilon D_2^2}{\eta T} \right).$$

624 The proof follows by applying Thm. A.1 to an intermediate result in the proof of Thm. 4.1
 625 in Liu et al. [25]. We defer a complete proof to Section E.

626 We note that the bound on reparameterized Adagrad is a function of \tilde{D}_∞ while the bound
 627 on Adagrad in original coordinates is a function of D_∞ . While \tilde{D}_∞^2 can be smaller or larger
 628 than D_∞^2 by up to a factor of d , in many typical settings $\tilde{D}_\infty \leq D_\infty$: if the constraint set
 629 Θ is chosen to be a ball, $D_\infty = \tilde{D}_\infty$. Optimizing convex objectives with L2 regularization–
 630 sometimes referred to as with optimization with weight decay in machine learning literature–
 631 corresponds to optimizing with Θ some ball centered at the origin. In Section B.3, we present
 632 experiments using L2 regularization with adaptive algorithms, and confirm that the empirical
 633 benefit from reparameterization persists in this setting.

634 **A.1.2. Convergence guarantees for non-convex objectives.** In the non-convex setting
 635 we can establish related local improved convergence guarantee for unconstrained Adagrad.
 636 For a given initial point θ_0 , let $\Delta_f(\theta_0) \stackrel{\text{def}}{=} f(\theta_0) - \min_\theta f(\theta)$. Observe that $\Delta_{\tilde{f}}(\tilde{\theta}_0) = \Delta_f(\theta_0)$
 637 for $\tilde{\theta}_0 \stackrel{\text{def}}{=} V^T \theta_0$. For unconstrained optimization, we add the following additional assumption:

638 **Assumption 3.** The objective $f(\cdot)$ is bounded below: $\inf_\theta f(\theta) > -\infty$.

639 **Theorem A.9 (Convergence for non-convex, unconstrained optimization).** Consider $f(\cdot)$, Θ ,
 640 and $\rho(\cdot)$ satisfying Assumptions 1, 2, and 3. Let $\tilde{f}(\cdot)$ and $\tilde{\Theta}$ denote the reparameterization
 641 of $f(\cdot)$ and Θ by V the eigenbasis of EGOP(f). Consider any initialization $\theta_0 \in \Theta$, and let
 642 $\tilde{\theta}_0 \stackrel{\text{def}}{=} V^T \theta_0$. Assume $\exists \delta \in [0, \beta_1^2)$ such that $\nabla^2 f(\cdot)$ has Lipschitz constant H satisfying the
 643 bound in Eq. 3.2. Let $\{\tilde{\theta}_t\}_{t=1}^T$ denote the iterates produced by running Algorithm A.1 using
 644 unconstrained updates and with inputs $(\tilde{f}(\cdot), \nabla \tilde{f}(\cdot), \tilde{\theta}_0, T, \epsilon)$ for $\epsilon < 1/d$. Then for all T such
 645 that $\{\tilde{\theta}_t\}_{t=1}^T \subseteq \tilde{\Theta}$, we have

$$646 \quad \frac{1}{T} \sum_{t=1}^T \|\nabla f(\theta_t)\|_1 = O \left(\frac{\Delta_f(\theta_0)}{\eta \sqrt{T}} + \frac{\eta \|\vec{L}\|_1}{\sqrt{T}} \left(\frac{\text{sr}(f)\sqrt{1+\delta}}{d(\beta_1 - \delta)} + \frac{\sqrt{\delta(1+\delta)}}{\beta_1 - \delta} \right) \log(p) \right)$$

647 where $p(T, \|\vec{L}\|_1, \nabla \tilde{f}(\tilde{\theta}_0))$ is polynomial in T , $\|\vec{L}\|_1$, and $\|\nabla \tilde{f}(\tilde{\theta}_0)\|_\infty$.

648 We compare the result in [Theorem A.9](#) to the convergence guarantee for Adagrad in original
649 coordinates initialized at θ_0 , as established in Jiang et al. [21]:

$$650 \quad \frac{1}{T} \sum_{t=1}^T \|\nabla f(\theta_t)\|_1 = O\left(\frac{\Delta_f(\theta_0)}{\eta\sqrt{T}} + \frac{\eta\|\vec{L}\|_1}{\sqrt{T}} \log(p(T, \|\vec{L}\|_1, \nabla f(\theta_0)))\right).$$

651 As with [Lemma A.8](#), this result implies that for $f(\cdot)$ with strong EGOP spectral decay and
652 dense leading EGOP eigenvectors, the local convergence bounds for reparameterized can be
653 smaller with by a factor of $1/d$.

654 We emphasize that the above is a *local* guarantee, and only holds for time horizons T such
655 that $\{\tilde{\theta}_t\}_{t=1}^T \subseteq \tilde{\Theta}$. If iterates move arbitrarily far away from the region in which EGOP samples
656 were concentrated, then without stronger assumptions about the landscape (i.e. vanishingly
657 small H) one cannot hope that the EGOP reparameterization will be informative for new
658 local geometry. For quadratic functions, $H = 0$ and thus the result holds for any T ; we note
659 that related works have restricted their analysis to quadratic functions for analyzing adaptive
660 algorithms [42].

661 *Proof of Thm. A.9.* The result follows immediately from combining the convergence result
662 in Thm. 3.1 of Jiang et al. [21] with the bound on $\|\tilde{L}\|_1/\|\vec{L}\|_1$ from [Thm. A.1](#). ■

663 **A.2. Proofs from Section 4.** In [Section 4](#), we noted that for objectives of the form
664 $f(\theta) = h(A\theta)$, for some loss function $h(\cdot)$ and data matrix $A \in \mathbb{R}^{d \times n}$, the EGOP of $f(\cdot)$ is

$$665 \quad (A.7) \quad \text{EGOP}(f) = A^\top \mathbb{E}_{\theta \sim \rho} \left[\nabla_\theta h(A\theta) \nabla_\theta h(A\theta)^\top \right] A$$

666 For general loss functions $h(\cdot)$, one can establish the following upper bounds showing how
667 the singular values of A control the EGOP eigepsectrum of $f(\cdot)$.

668 **Lemma A.10.** Consider $f : \mathbb{R}^n \rightarrow \mathbb{R}$ satisfying $f(\theta) = h(A\theta)$ for some loss function $h :$
669 $\mathbb{R}^n \rightarrow \mathbb{R}$ and nonsingular data matrix $A \in \mathbb{R}^{n \times n}$. Denote by $\sigma_i(\cdot)$ and $\lambda_i(\cdot)$ the i th singular
670 value and eigenvalue of a matrix respectively, indexed by decreasing value. Then all nonzero
671 eigenvalues of the EGOP of $f(\cdot)$ satisfy

$$672 \quad \frac{\lambda_k(\text{EGOP}(f))}{\lambda_1(\text{EGOP}(f))} \leq \left(\frac{\sigma_k(A)}{\sigma_1(A)} \right)^2 \frac{\lambda_1(M)}{\lambda_n(M)}.$$

673 where

$$674 \quad M \stackrel{\text{def}}{=} \mathbb{E}_{\theta \sim \rho} [\nabla_\theta h(A\theta) \nabla_\theta h(A\theta)^\top].$$

675 If M has some finite condition number that does not go to infinity as the spectral decay in A
676 increases, then [Lemma A.10](#) shows that increasing spectral decay in the data matrix A induces
677 spectral decay in $\text{EGOP}(f)$. For specific choices of $h(\cdot)$, one can more precisely characterize
678 how spectral decay in the matrix A induces decay in the EGOP of $f \stackrel{\text{def}}{=} h \circ A$:

679 **Lemma A.11.** Consider $A \in \mathbb{R}^{d \times n}$ and let $f(\theta) = \frac{1}{2} \|A\theta - y\|_2^2$ where $y = A\theta^* + \eta$, for
680 η some mean-zero measurement noise. Assume sampling density ρ is a standard Gaussian
681 distribution. Then the eigenvalues of $\text{EGOP}(f)$, $\{\lambda_k\}_{k=1}^d$ indexed in decreasing order, satisfy

$$682 \quad (A.8) \quad \frac{\lambda_k}{\lambda_1} \leq \left(\frac{\sigma_{k-1}(A)}{\sigma_1(A)} \right)^4 \quad \forall k \in [2, \dots, n]$$

683 where $\sigma_i(A)$ denotes the i^{th} singular value of A , indexed in decreasing order.

684 We defer the proof of Lemma A.11 to Section E.2.

685

REFERENCES

- 686 [1] E. ALPAYDIN AND C. KAYNAK, *Optical Recognition of Handwritten Digits*. UCI Machine Learning Repos-
687 itory, 1998. DOI: <https://doi.org/10.24432/C50P49>.
- 688 [2] J. BRADBURY, R. FROSTIG, P. HAWKINS, M. J. JOHNSON, C. LEARY, D. MACLAURIN, G. NECULA,
689 A. PASZKE, J. VANDERPLAS, S. WANDERMAN-MILNE, AND Q. ZHANG, *JAX: composable transforma-*
690 *tions of Python+NumPy programs*, 2018, <http://github.com/jax-ml/jax>.
- 691 [3] C. CARTIS, X. LIANG, E. MASSART, AND A. OTEMISOV, *Learning the subspace of variation for global*
692 *optimization of functions with low effective dimension*, arXiv preprint arXiv:2401.17825, (2024).
- 693 [4] X. CHEN AND E. HAZAN, *Open problem: Black-box reductions and adaptive gradient methods for noncon-*
694 *convex optimization*, in Proceedings of Thirty Seventh Conference on Learning Theory, S. Agrawal and
695 A. Roth, eds., vol. 247 of Proceedings of Machine Learning Research, PMLR, 30 Jun–03 Jul 2024,
696 pp. 5317–5324, <https://proceedings.mlr.press/v247/chen24e.html>.
- 697 [5] H.-H. CHOU, C. GIESHOFF, J. MALY, AND H. RAUHUT, *Gradient descent for deep matrix factorization:*
698 *Dynamics and implicit bias towards low rank*, Applied and Computational Harmonic Analysis, 68
699 (2024), p. 101595.
- 700 [6] P. G. CONSTANTINE, *Active subspaces: Emerging ideas for dimension reduction in parameter studies*,
701 SIAM, 2015.
- 702 [7] R. COSSON, A. JADBABAIE, A. MAKUR, A. REISIZADEH, AND D. SHAH, *Low-Rank Gradient Descent*,
703 IEEE Open Journal of Control Systems, (2023).
- 704 [8] B. CREW, *Google Scholar reveals its most influential papers for 2020*, Jul 2020.
- 705 [9] C. CUI, K. ZHANG, T. DAULBAEV, J. GUSAK, I. OSELEDETS, AND Z. ZHANG, *Active subspace of neural*
706 *networks: Structural analysis and universal attacks*, SIAM Journal on Mathematics of Data Science,
707 2 (2020), pp. 1096–1122.
- 708 [10] DEEPMIND, I. BABUSCHKIN, K. BAUMLI, A. BELL, S. BHUPATIRAJU, J. BRUCE, P. BUCHLOVSKY,
709 D. BUDDEN, T. CAI, A. CLARK, I. DANIELKA, A. DEDIEU, C. FANTACCI, J. GODWIN, C. JONES,
710 R. HEMSLEY, T. HENNIGAN, M. HESSEL, S. HOU, S. KAPTUROWSKI, T. KECK, I. KEMAEV,
711 M. KING, M. KUNESCH, L. MARTENS, H. MERZIC, V. MIKULIK, T. NORMAN, G. PAPAMAKARIOS,
712 J. QUAN, R. RING, F. RUIZ, A. SANCHEZ, L. SARTRAN, R. SCHNEIDER, E. SEZENER, S. SPENCER,
713 S. SRINIVASAN, M. STANOJEVIĆ, W. STOKOWIEC, L. WANG, G. ZHOU, AND F. VIOLA, *The DeepMind*
714 *JAX Ecosystem*, 2020, <http://github.com/google-deeppmind>.
- 715 [11] A. DÉFOSSEZ, L. BOTTOU, F. BACH, AND N. USUNIER, *A Simple Convergence Proof of Adam and*
716 *Adagrad*, arXiv preprint arXiv:2003.02395, (2020).
- 717 [12] J. DENG, W. DONG, R. SOCHER, L.-J. LI, K. LI, AND L. FEI-FEI, *Imagenet: A large-scale hierarchi-*
718 *cal image database*, in 2009 IEEE Conference on Computer Vision and Pattern Recognition, 2009,
719 pp. 248–255, <https://doi.org/10.1109/CVPR.2009.5206848>.
- 720 [13] Z. DONG, Y. ZHANG, J. YAO, AND R. SUN, *Towards quantifying the hessian structure of neural networks*,
721 arXiv preprint arXiv:2505.02809, (2025).
- 722 [14] J. DUCHI, E. HAZAN, AND Y. SINGER, *Adaptive subgradient methods for online learning and stochastic*
723 *optimization.*, Journal of Machine Learning Research, 12 (2011).
- 724 [15] X. GLOROT AND Y. BENGIO, *Understanding the difficulty of training deep feedforward neural networks*, in
725 Proceedings of the thirteenth international conference on artificial intelligence and statistics, JMLR
726 Workshop and Conference Proceedings, 2010, pp. 249–256.
- 727 [16] G. H. GOLUB, *Some modified matrix eigenvalue problems*, SIAM Review, 15 (1973), pp. 318–334.
- 728 [17] V. GUPTA, T. KOREN, AND Y. SINGER, *Shampoo: Preconditioned stochastic tensor optimization*, in
729 International Conference on Machine Learning, PMLR, 2018, pp. 1842–1850.
- 730 [18] N. HALKO, P.-G. MARTINSSON, AND J. A. TROPP, *Finding structure with randomness: Probabilistic*
731 *algorithms for constructing approximate matrix decompositions*, SIAM Review, 53 (2011), pp. 217–
732 288.

- 733 [19] E. HAZAN ET AL., *Introduction to online convex optimization*, Foundations and Trends® in Optimization,
734 2 (2016), pp. 157–325.
- 735 [20] K. HE, X. ZHANG, S. REN, AND J. SUN, *Deep residual learning for image recognition*, in Proceedings of
736 the IEEE Conference on Computer Vision and Pattern Recognition (CVPR), June 2016.
- 737 [21] R. JIANG, D. MALADKAR, AND A. MOKHTARI, *Convergence analysis of adaptive gradient methods under
738 refined smoothness and noise assumptions*, arXiv preprint arXiv:2406.04592, (2024).
- 739 [22] D. P. KINGMA AND J. BA, *Adam: A method for stochastic optimization*, 2017, [https://arxiv.org/abs/
740 1412.6980](https://arxiv.org/abs/1412.6980), <https://arxiv.org/abs/1412.6980>.
- 741 [23] F. KUNSTNER, R. YADAV, A. MILLIGAN, M. SCHMIDT, AND A. BIETTI, *Heavy-tailed class imbalance
742 and why Adam outperforms gradient descent on language models*, arXiv preprint arXiv:2402.19449,
743 (2024).
- 744 [24] S. Z. LING, N. SHARP, AND A. JACOBSON, *Vectoradam for rotation equivariant geometry optimization*,
745 Advances in Neural Information Processing Systems, 35 (2022), pp. 4111–4122.
- 746 [25] Y. LIU, R. PAN, AND T. ZHANG, *AdaGrad under Anisotropic Smoothness*, arXiv preprint
747 arXiv:2406.15244, (2024).
- 748 [26] I. LOSHCHILOV AND F. HUTTER, *Decoupled weight decay regularization*, 2019, [https://arxiv.org/abs/1711.
749 05101](https://arxiv.org/abs/1711.05101), <https://arxiv.org/abs/1711.05101>.
- 750 [27] L. MAES, T. H. ZHANG, A. JOLICOEUR-MARTINEAU, I. MITLIAGKAS, D. SCIEUR, S. LACOSTE-JULIEN,
751 AND C. GUILLE-ESCURET, *Understanding Adam Requires Better Rotation Dependent Assumptions*,
752 arXiv preprint arXiv:2410.19964, (2024).
- 753 [28] N. MALLINAR, D. BEAGLEHOLE, L. ZHU, A. RADHAKRISHNAN, P. PANDIT, AND M. BELKIN, *Emergence
754 in non-neural models: grokking modular arithmetic via average gradient outer product*, arXiv preprint
755 arXiv:2407.20199, (2024).
- 756 [29] V. POPYAN, *The full spectrum of deepnet Hessians at scale: Dynamics with SGD training and sample size*,
757 arXiv preprint arXiv:1811.07062, (2018).
- 758 [30] A. PASZKE, S. GROSS, S. CHINTALA, G. CHANAN, E. YANG, Z. DEVITO, Z. LIN, A. DESMAISON,
759 L. ANTIGA, AND A. LERER, *Automatic differentiation in PyTorch*, in NIPS-W, 2017.
- 760 [31] A. RADHAKRISHNAN, D. BEAGLEHOLE, P. PANDIT, AND M. BELKIN, *Mechanism of feature learning
761 in deep fully connected networks and kernel machines that recursively learn features*, arXiv preprint
762 arXiv:2212.13881, (2022).
- 763 [32] S. J. REDDI, S. KALE, AND S. KUMAR, *On the Convergence of Adam and Beyond*, 2019, [https://arxiv.
764 org/abs/1904.09237](https://arxiv.org/abs/1904.09237), <https://arxiv.org/abs/1904.09237>.
- 765 [33] M. RUDELSON AND R. VERSHYNIN, *Sampling from large matrices: An approach through geometric func-
766 tional analysis*, Journal of the ACM (JACM), 54 (2007), pp. 21–es.
- 767 [34] L. SAGUN, U. EVCI, V. U. GUNEY, Y. DAUPHIN, AND L. BOTTOU, *Empirical analysis of the hessian of
768 over-parametrized neural networks*, arXiv preprint arXiv:1706.04454, (2017).
- 769 [35] M. UDELL AND A. TOWNSEND, *Why are big data matrices approximately low rank?*, SIAM Journal on
770 Mathematics of Data Science, 1 (2019), pp. 144–160.
- 771 [36] P. VIRTANEN, R. GOMMERS, T. E. OLIPHANT, M. HABERLAND, T. REDDY, D. COURNAPEAU,
772 E. BUROVSKI, P. PETERSON, W. WECKESSER, J. BRIGHT, S. J. VAN DER WALT, M. BRETT,
773 J. WILSON, K. J. MILLMAN, N. MAYOROV, A. R. J. NELSON, E. JONES, R. KERN, E. LAR-
774 SON, C. J. CAREY, İ. POLAT, Y. FENG, E. W. MOORE, J. VANDERPLAS, D. LAXALDE, J. PERK-
775 TOLD, R. CIMRMAN, I. HENRIKSEN, E. A. QUINTERO, C. R. HARRIS, A. M. ARCHIBALD, A. H.
776 RIBEIRO, F. PEDREGOSA, P. VAN MULBREGT, AND SciPy 1.0 CONTRIBUTORS, *SciPy 1.0: Fun-
777 damental Algorithms for Scientific Computing in Python*, Nature Methods, 17 (2020), pp. 261–272,
778 <https://doi.org/10.1038/s41592-019-0686-2>.
- 779 [37] N. VYAS, D. MORWANI, R. ZHAO, I. SHAPIRA, D. BRANDFONBRENER, L. JANSON, AND S. KAKADE,
780 *SOAP: Improving and stabilizing Shampoo using Adam*, arXiv preprint arXiv:2409.11321, (2024).
- 781 [38] R. WARD, X. WU, AND L. BOTTOU, *Adagrad stepsizes: Sharp convergence over nonconvex landscapes*,
782 Journal of Machine Learning Research, 21 (2020), pp. 1–30.
- 783 [39] H. XIAO, K. RASUL, AND R. VOLLGRAF, *Fashion-mnist: a novel image dataset for benchmarking machine
784 learning algorithms*, arXiv preprint arXiv:1708.07747, (2017).
- 785 [40] S. XIE, M. A. MOHAMADI, AND Z. LI, *Adam Exploits ℓ_∞ -geometry of Loss Landscape via Coordinate-wise
786 Adaptivity*, 2024, <https://arxiv.org/abs/2410.08198>, <https://arxiv.org/abs/2410.08198>.

- 787 [41] M. D. ZEILER, *Adadelta: An adaptive learning rate method*, 2012, <https://arxiv.org/abs/1212.5701>, <https://arxiv.org/abs/1212.5701>.
788
- 789 [42] Y. ZHANG, C. CHEN, T. DING, Z. LI, R. SUN, AND Z.-Q. LUO, *Why transformers need Adam: A Hessian*
790 *perspective*, arXiv preprint arXiv:2402.16788, (2024).
- 791 [43] J. ZHAO, Z. ZHANG, B. CHEN, Z. WANG, A. ANANDKUMAR, AND Y. TIAN, *Galore: Memory-efficient*
792 *llm training by gradient low-rank projection*, arXiv preprint arXiv:2403.03507, (2024).
- 793 [44] L. ZHU, C. LIU, A. RADHAKRISHNAN, AND M. BELKIN, *Catapults in SGD: spikes in the training loss*
794 *and their impact on generalization through feature learning*, arXiv preprint arXiv:2306.04815, (2023).

795 **Supplementary material begins here.** _____

796 **Appendix B. Supplementary Figures.**

797 **B.1. Supplementary Figures for Section 4.**

798 *Spectral decay is robust to choice of sampling distribution.* In Figure 10, we present evidence
 799 that spectral decay visualized in Figure 3 is robust to choice of sampling distribution. The
 800 EGOP whose eigenspectrum is displayed in Figure 3 was generated from gradient samples
 801 $\nabla f(\theta_i)$ for $\theta_i \sim \mathcal{N}(0, \mathbb{I})$, but as visualized in later figures (Figure 10), the level of spectral
 802 decay is comparable to that obtained with Gaussian distributions of differing scales. We
 803 also compare to the spectral decay exhibited by the EGOP matrix when ρ is taken to be an
 804 initialization distribution used in practice [15]. For details on this distribution, see Section C.6.

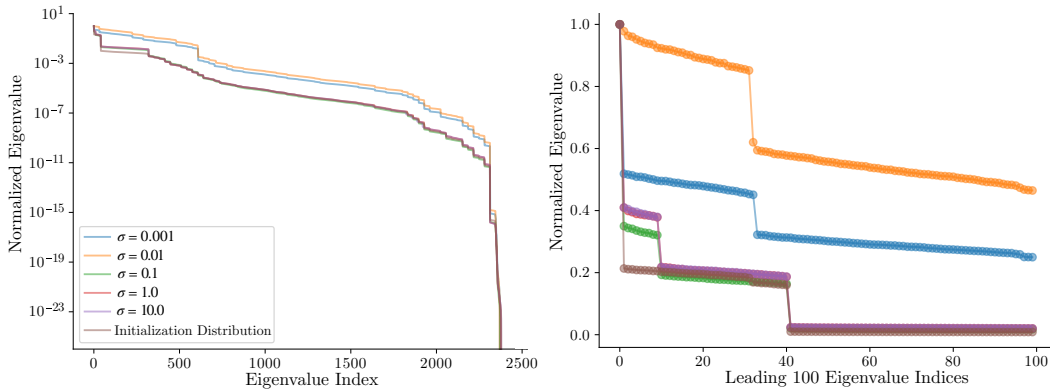


Figure 10: Comparing EGOP spectral decay of a 2-layer ReLU network on tinyMNIST dataset. Plot displays the ratio λ_k/λ_1 as a function of eigenvalue index k , for eigenvalues indexed in decreasing order. The blue, orange, green, red, and purple colored traces display the eigenspectrum of the EGOP with respect to a mean-zero Gaussian with covariance $\sigma^2\mathbb{I}$, for varying values of σ . The brown trace displays the eigenspectrum of the EGOP with respect to a realistic initialization distribution for this architecture: weights for each layer are drawn from a scaled Xavier normal distribution, and biases are initialized from a scaled uniform distribution (see Section C.6). We observe that under all sampling distributions, the eigenspectrum exhibits spectral decay, and that the realistic initialization distribution has spectral decay very comparable to that of the standard Gaussian, displayed in Figure 3 of the main body.

805 We also note that the interesting shelf structure of the leading eigenvalues is also robust
 806 to choice of sampling distribution, as illustrated in Figure 10. Full details for the objective
 807 and EGOP estimation procedure for this figure are detailed in Section C.5.

808 *Spectral decay persists in block EGOP matrices.* Figure 11 plots the normalized eigenspectra
 809 of the block matrices corresponding to the first and second layers of ReLU networks on the
 810 UCI digits dataset and fashionMNIST dataset respectively. For full details on these datasets
 811 and the architectures used, see Section C.5. Interestingly, both datasets exhibit shared char-
 812 acteristics: the normalized spectral decay in the first layer is strikingly similar, and in both
 813 networks the spectral decay in the first layer is more pronounced than in the second layer.

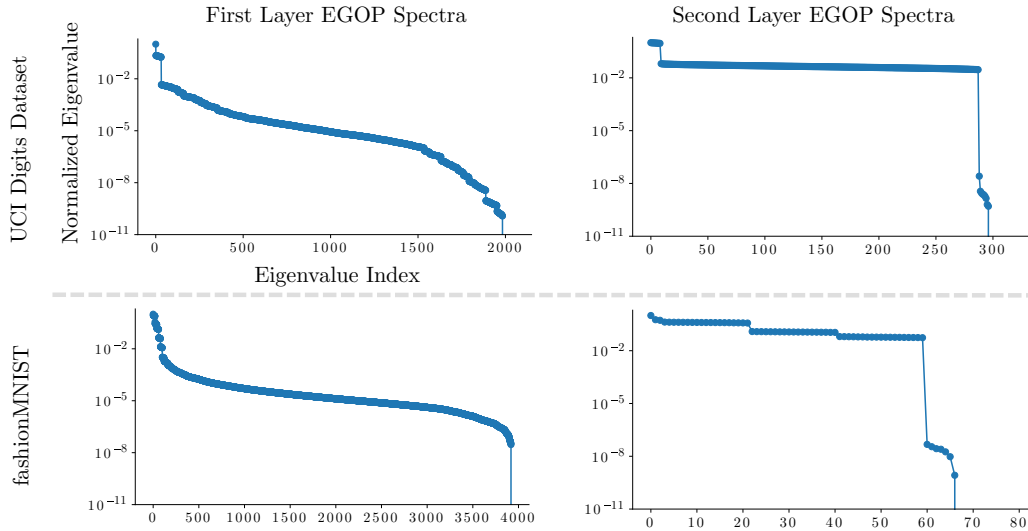


Figure 11: Eigenspectra of the layerwise EGOP matrices of neural networks on `tinyMNIST` and `fashionMNIST`. The spectral decay observed in Figure 3 persists for layer EGOP matrices, defined in Section 5, and across datasets. Y-axes for all figures display identical ranges.

814 *Density of leading EGOP eigenvectors.* Thms. A.8 and A.9 show that when the lead-
 815 ing EGOP eigenvector is dense, reparameterized Adagrad enjoys much stronger convergence
 816 guarantees. Density of the k^{th} eigenvector is measured by $\beta_k \stackrel{\text{def}}{=} \|v_k\|_1^2$. Our results state a
 817 guarantee in terms of $\beta \stackrel{\text{def}}{=} \beta_1$, the density of the leading eigenvector.

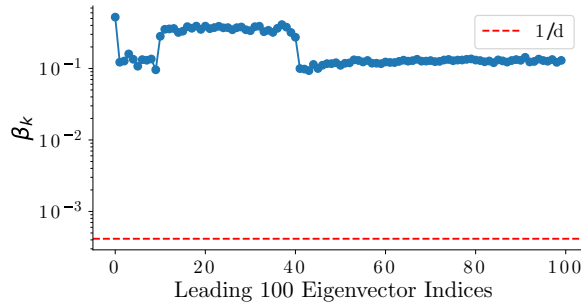


Figure 12: Plotting the density measure $\beta_k \stackrel{\text{def}}{=} \|v_k\|_1^2$ for the leading 100 eigenvectors of $\text{EGOP}(f)$, where $f(\cdot)$ is the cross-entropy loss of a 2-layer ReLU neural network on the UCI digits training dataset. The leading eigenvector satisfies $\beta_1 > 0.5$ and several have density $\beta_k > 0.3$. We visualize the value $1/d$ in red (for this example, $d = 2,410$) to verify that for the leading eigenvectors, $\beta_k \gg 1/d$.

818 In order for the factors $\text{sr}(f)/d(\beta_1 - \delta)$ in the bound for reparameterized Adagrad from
 819 Thms. A.8 and A.9 to reflect an improvement, it is necessary that $\beta_1 \gg 1/d$. In Figure 12,

820 we show that for the UCI digits dataset, this condition is satisfied. Specifically, we show that
 821 $\beta_1 > 0.5$, while for this network $1/d < 5e - 4$. Moreover, we show that not only does the
 822 leading eigenvector satisfy this density assumption, but several of the leading eigenvectors
 823 satisfy $\beta_k \gg 1/d$.

824 B.2. Supplementary Figures for Section 6.

825 *Residual Networks.* In Figure 13, we supplement the results shown in Figures 2 and 4. In
 826 Figure 13, we plot training and validation accuracy over epochs, confirming that the improved
 827 training loss convergence under reparameterization leads to improved classification accuracy.

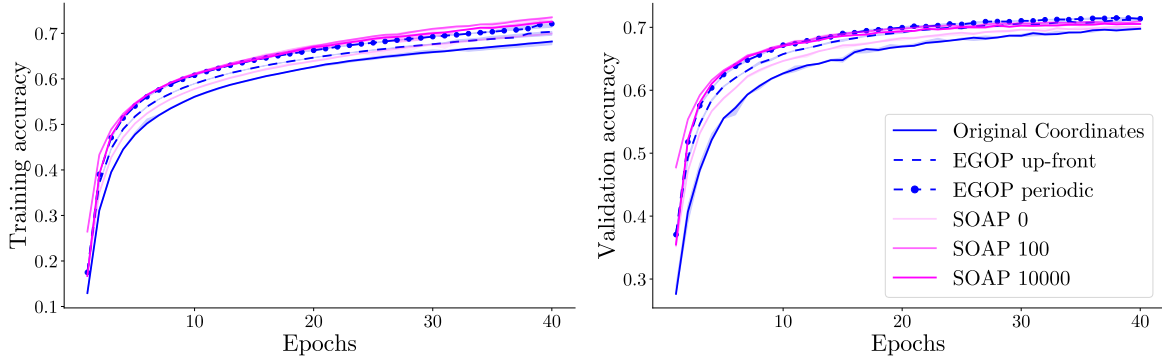


Figure 13: Training and validation accuracy for image classification with residual networks. The improved training and validation loss convergence displayed in Figures 2 and 4 leads to improved training and validation accuracy. EGOP reparameterization is competitive with SOAP in accuracy metrics as well as loss metrics.

828 *Multilayer linear networks.* We compare three methods for EGOP reparameterization in
 829 order to examine some heuristics proposed in Section 5. In Figure 14, we consider training a
 830 multilayer linear network (6.1) under global EGOP reparameterization, wherein all parameters
 831 are reparameterized simultaneously as in Algorithm 2.1 (Figure 14a); block reparameterization
 832 for all layers, following the procedure defined in Section 5 (Figure 14b); and block reparam-
 833 eterization of only the parameters in the first layer (Figure 14c). For all three methods, we
 834 estimate the EGOP using the same number of gradient samples: $M = 2d$, where d is the total
 835 number of network parameters.

836 Comparing Figure 14a with Figure 14b shows that for this problem, EGOP reparameter-
 837 ization offers comparable benefit when using block reparameterization as when using global
 838 reparameterization. This suggests that block EGOP reparameterization, which has a reduced
 839 computational cost compared to global EGOP reparameterization, may be an effective way to
 840 accelerate adaptive methods when problem instances are too large to permit global reparam-
 841 eterization. Reparameterizing only the first layer (Figure 14c) improves Adagrad’s performance
 842 by a margin comparable to that of global and block reparameterization of all layers, but the
 843 benefit to Adam under reparameterization of only the first layer is much less pronounced.

844 *Image Classification with ReLU Networks.* Figure 15 expands on the results shown in Fig-
 845 ure 1 (right) for 2-layer ReLU networks on the UCI digits dataset. Figure 15a plots final

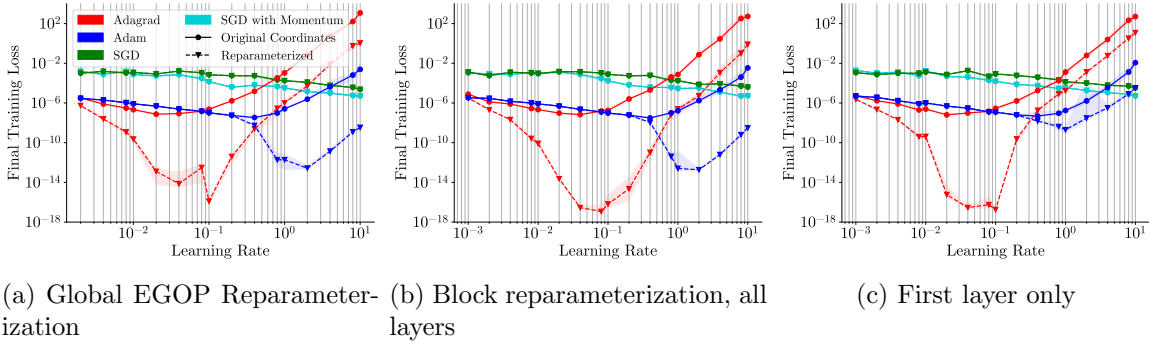


Figure 14: Comparing three EGOP reparameterization methods for training a multilayer linear network (6.1). Figure 14a is a reproduction of Figure 6b in the main body, and shows results for performing EGOP reparameterization of all parameters simultaneously. Figure 14b shows results when performing block EGOP reparameterization, where each network layer forms a block. Figure 14c shows results when block-reparameterizing only the parameters in the first layer.

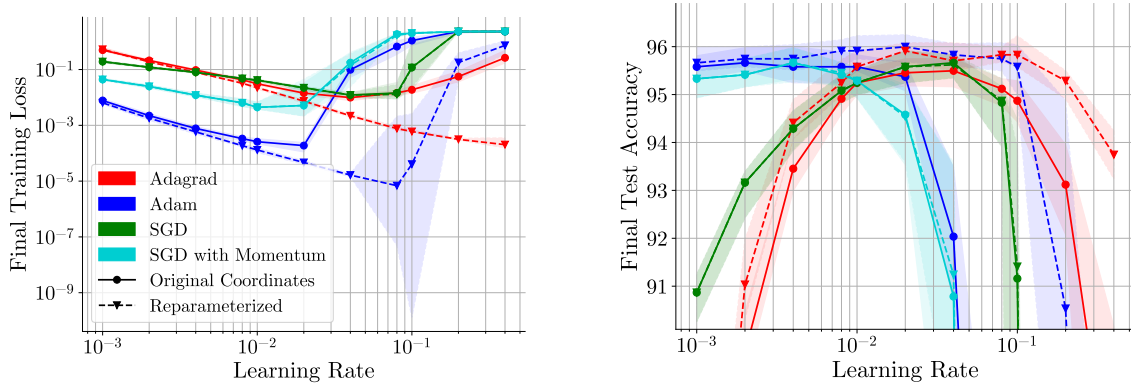
846 training loss versus learning rate, and shows that benefit of reparameterization shown in Fig-
 847 ure 1 is robust to choice of learning rate. Figure 15b plots final test accuracy versus learning
 848 rate, and shows that the improved training offered by reparameterization does not lead to
 849 over-fitting. For full experimental details on the architecture, dataset, and loss function used
 850 for these experiments, see Section C.

851 Similarly, Figure 16 expands on the results presented in Figure 7. Figure 16 (left) shows
 852 that the improved training does not lead to over-fitting, but rather that reparameterization
 853 leads to improved accuracy on hold-out data. Figure 16 (right) reports wallclock time to con-
 854 vergence, demonstrating that the improved convergence in epochs shown in Figure 7 translates
 855 to faster convergence in wallclock time. For this experiment, we define convergence as the first
 856 epoch at which the model attains 99.5% of its maximum validation accuracy. In Figure 16
 857 (left), we use square a circular markers to indicate the epochs at which the model in original
 858 coordinates and reparameterized models respectively achieve this threshold.

859 **Convex Objectives.** In addition to training neural networks, which is a primary application
 860 of adaptive optimization algorithms, we also study reparameterization for convex optimization.
 861 Figure 17a shows the result of minimizing

$$862 \text{ (B.1)} \quad f(\theta) = \log \left(\sum_{i=1}^n \exp(\langle a_i, \theta \rangle - y_i)^2 \right),$$

863 where $a_i \in \mathbb{R}^d$ denote vectors of observations and $y_i \stackrel{\text{def}}{=} \langle a_i, \theta^* \rangle$ for ground truth θ^* . We
 864 also plot results from minimizing logistic regression objectives (Figure 17b) and linear least-
 865 squares objectives (Figure 17c) arising from problems with data matrices $A \in \mathbb{R}^{n \times d}$. For all
 866 convex objectives, we induce EGOP spectral decay by choosing matrices A with singular value
 867 decay. For these objectives, we estimate the EGOP and perform optimization using noiseless



(a) Median final training loss versus learning rate

(b) Median final test accuracy versus learning rate

Figure 15: Training loss and test accuracy results for the UCI digits dataset image classification task. Results are aggregated over independent trials corresponding to different random initializations. Medians are plotted as traces, and shaded regions indicate the 25th-75th percentiles.

868 (full-batch) gradients.

869 Figure 17 demonstrates that EGOP reparameterization can improve convergence of adap-
 870 tive algorithms to global minima of convex objectives. EGOP-reparameterized Adagrad out-
 871 performs its counterpart in original coordinates for all three objectives. Notably, methods with
 872 momentum excel for this logistic regression, and reparameterization boosts the performance of
 873 Adagrad (which does not use momentum), making it competitive with the momentum-based
 874 optimizers. EGOP reparameterization improves Adam’s convergence on the log-sum-exp ob-
 875 jective and linear least-squares objectives (Figures 17a and 17c), but has no impact on Adam’s
 876 performance for logistic regression (Figure 17b).

877 **B.3. Optimizing with Weight Decay.** In Section 3, we instantiated guarantees for a ball
 878 constraint set, and remarked that in practice adaptive optimization algorithms are often em-
 879 ployed along with weight decay/L2 regularization, which implicitly constrains the algorithms
 880 to some ball in parameter space centered at the origin. We note that our experiments on
 881 residual networks and on image classification using fashionMNIST both employ AdamW, an
 882 adaptive algorithm with explicit weight decay, thus validating that EGOP reparameterization
 883 can also improve convergence when weights are implicitly constrained to some ball.

884 Appendix C. Experimental details.

885 Experiments were implemented in python, and both code and data will be made publicly
 886 available upon publication. Experiments with neural networks were implemented in Pytorch,
 887 and experiments with convex objectives were implemented using auto-differentiation in Jax
 888 and optimizers from Optax [30, 2, 10]. With the exception of experiments performed on
 889 residual networks, all experiments can be performed on CPU with a MacBook Pro with an

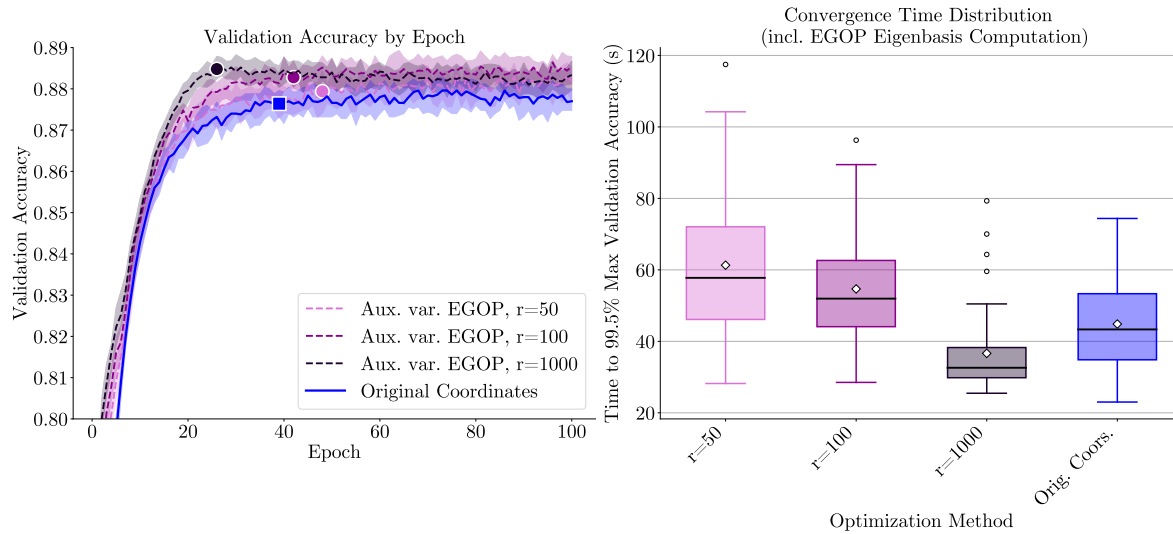


Figure 16: Validation accuracies and timing results for image classification on fashionMNIST. Counterpart to Figure 7. (Left) Plotting validation accuracy demonstrates that the improved minimization of training loss leads to improved generalization. (Right) Plotting wallclock time to convergence for the model in original coordinates, as well as under auxiliary variable EGOP reparameterization, as described in Section 5, for varying values of r . We find that the cost-benefit tradeoff favors intermediate values of r , e.g. $r = 1000$ in this figure. Recall that for this architecture, $r = 1000$ is much smaller than the number of optimization variables, $d \approx 78k$.

890 Intel Core i5 processor and 16GB memory. However, all experiments can be run more quickly
 891 on GPU, and GPU support is included in the codebase. Experiments with residual networks
 892 were performed on an internal cluster, using 1 RTX 6000 Ada GPU, and 8 CPU cores on an
 893 AMD EPYC 74F3 processor. The timing results with auxiliary variable reparameterization
 894 for ReLU networks in Figure 16 were performed on an internal cluster using 1 RTX A6000
 895 GPU and 2 CPU cores on an Intel Xeon Silver 4114 processor.

896 C.1. Dataset Details.

897 *UCI handwritten digits dataset.* The UCI digits dataset [1] contains 5620 instances, each
 898 an 8×8 pixel grayscale image of a handwritten digit of values $0, \dots, 9$. Each instance has
 899 an integer label, $0, \dots, 9$. We split the dataset into a training dataset with 3823 instances, a
 900 validation dataset with 598 instances, and a test dataset with 1199 instances.

901 *fashionMNIST dataset.* The fashionMNIST dataset consists of 28×28 pixel grayscale im-
 902 ages of clothing, each having a label from one of 10 classes. The full dataset contains 60k
 903 training samples and 10k test samples; for the results in Figure 7, we use the full training
 904 set and we subdivide the test set into a validation set of 2.5k samples and a test set of 7.5k
 905 samples [39]. For the experiments in Figure 11, we restrict to only instances corresponding
 906 to the first four classes (labeled “t-shirt,” “trouser,” “pullover,” and “dress”). This yields

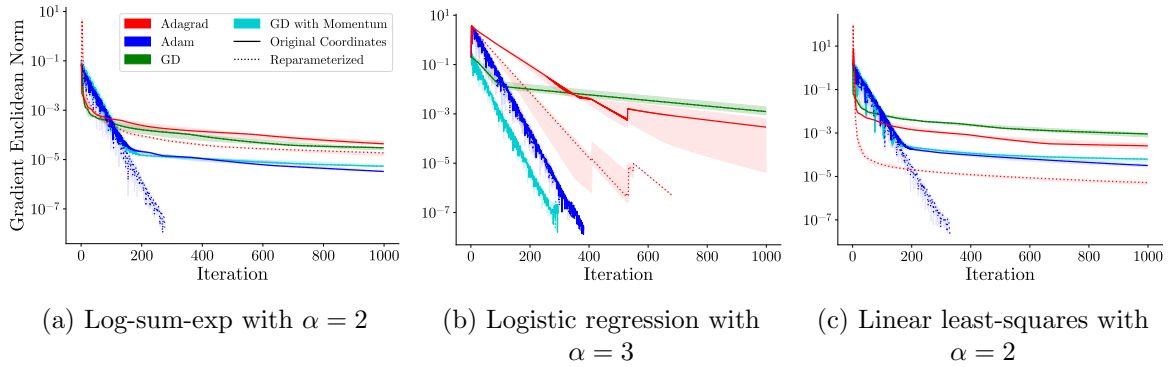


Figure 17: Gradient Euclidean norm of solution at t^{th} iterate. Learning rates tuned separately for methods in original coordinates and under reparameterization. We induce EGOP spectral decay by choice of data matrix A with singular values $\sigma_k(A) = k^{-\alpha}$. As noted in the prose, in some plots the dotted traces coincide with the solid and are thus not visible (Adam in Figure 17b, and equivariant methods GD and GD with momentum).

907 a dataset of 28,000 instances, which we subdivide into a training set of 24,000 instances, a
 908 validation set of 1,000 instances, and a test set of 3,00 instances. We do this in order to reduce
 909 the problem size to a setting where we can visualize the full EGOP eigenspectrum obtained
 910 without using heuristics.

911 *ImageNet dataset.* For our large-scale experiments, we use the ImageNet classification
 912 dataset [12]. ImageNet consists of approximately 1.28 million training images and 50,000
 913 validation images spanning 1,000 object categories. Following standard practice, we train on
 914 the full training set and report results on the validation set. During training, all models use
 915 the conventional preprocessing pipeline of random resized crops to 224×224 and horizontal
 916 flips. For evaluation, we resize each image so that its shorter side is 256 pixels and use a
 917 ten-crop evaluation protocol (four corners, center crop, and their horizontal flips). ImageNet
 918 serves as a challenging benchmark to evaluate the scalability and efficiency of our method in
 919 deep convolutional architectures.

920 C.2. Architecture Details.

921 *Architecture for multilayer linear networks.* To implement the objective

$$922 \quad (\text{C.1}) \quad f(\theta) = \|W_3 W_2 W_1 A - Y\|_{\mathbb{F}}^2 / n_{\text{samples}}$$

923 we use a 3-layer fully connected network without any nonlinear activations. In accordance
 924 with the description in Section 6, the first layer has 10 input nodes and 30 output nodes, the
 925 second layer has 50 output nodes, and the last layer has 10 output nodes.

926 *Architecture for UCI handwritten digits.* For the UCI digits dataset we use a 2-layer ReLU
 927 network. Its first layer is a fully-connected linear layer with bias, with 64 input nodes and
 928 32 output nodes, followed by a ReLU nonlinearity. Its second layer is a fully-connected linear
 929 layer with bias, with 32 input nodes and 10 output nodes, followed by a log-softmax. We form

930 the function $f(\cdot)$ by taking the negative log-likelihood loss on the training dataset.

931 *Large architecture for fashionMNIST dataset.* For the results in Figure 7, we use a 2-layer
932 ReLU network. The first layer is a fully-connected linear layer with bias, with 784 input
933 nodes and 100 output nodes, followed by a ReLU activation function. The second layer is a
934 fully-connected linear layer with bias, with 100 input nodes and 10 output nodes. We form
935 the function $f(\cdot)$ by taking the cross entropy loss on the training dataset. We note that this
936 is equivalent to computing the log-softmax on the network output and then using negative
937 log-likelihood loss, which is how we compute $f(\cdot)$ on tinyMNIST.

938 *Small architecture for fashionMNIST dataset.* **Architecture details for auxilliary.**

939 *Architecture for ImageNet.* For the ImageNet experiments, we use the 34 layer residual
940 network (ResNet-34) architecture presented in the original paper [20]. The model begins with
941 a convolutional layer consisting of a 7×7 convolution with 64 channels and stride 2, followed
942 by batch normalization, a ReLU activation, and a 3×3 max pooling layer with stride 2. The
943 network then contains four stages of residual blocks: the first stage has three residual blocks
944 with 64 channels; the second stage has four residual blocks with 128 channels, with the first
945 block using stride 2 for spatial downsampling; the third stage has six residual blocks with 256
946 channels, again with stride 2 in the first block; and the final stage has three residual blocks
947 with 512 channels, with stride 2 in the first block. Each residual block consists of two 3×3
948 convolutional layers with batch normalization and ReLU activations, together with a skip
949 connection. After the residual stages, we apply global average pooling and a fully connected
950 layer mapping from 512 features to 1000 output classes. We define the function $f(\cdot)$ as the
951 cross entropy loss computed on the ImageNet training dataset.

952 **C.3. Algorithm Details.** We consider four optimization algorithms: Adagrad, Adam,
953 SGD, and SGD with momentum. We use standard settings for all hyperparameters except
954 learning rate: for Adam, we use $\beta_1 = 0.9$ and $\beta_2 = 0.999$, and for SGD with momentum
955 we use momentum parameter 0.9, matching that of Adam. We use zero weight decay for all
956 algorithms, and for all objectives except the linear feed-forward networks, we use a constant
957 learning rate schedule.

958 We tune learning rates for each algorithm. Unless otherwise noted, when choosing the
959 range of learning rates to sweep for tuning, we the following two-step procedure. First we
960 sweep a coarse sweep using powers of 10 (e.g. $1e-3, 1e-2, \dots, 100$) to find an upper and lower
961 bound on the learning rates that produce best performance. The metric by which we quantify
962 the “best performance” for each different experiment is discussed below in Section C.6. We
963 then perform a refined sweep, where we discretize each interval between subsequent power of
964 10 using doubling. For example, if the coarse sweep identified lower bound 0.01 and upper
965 bound 1.0, we would perform the refined sweep over values $[0.01, 0.02, 0.04, 0.08, 0.1, 0.2, 0.4,$
966 $0.8, 1.0]$.

967 *Auxiliary Variables.* **Explain that we expect best results when we use weight decay because**
968 **we’re adding overparameterization i.e. auxiliary variables. Explain what values we used in**
969 **our weight decay sweep.**

970 *ResNet for ImageNet.* Following the standard ResNet training protocol, we use hyperpa-
971 rameters comparable to those reported in the original work. Although the original ResNet
972 work [20] trains using SGD with momentum 0.9, our experiments use AdamW. To preserve

comparable smoothing of gradient information, we set $\beta_1 = 0.9$, which plays a role analogous to momentum in SGD, together with the standard choice $\beta_2 = 0.999$. We use the same weight decay parameter as the original ResNet work, namely 10^{-4} . For learning rate selection, we perform a logarithmic sweep over the interval $[10^{-5}, 1]$, training each candidate learning rate for 10 epochs and selecting the value that achieves the highest top-1 validation accuracy under ten-crop evaluation. After selecting the learning rate, we train the final model for 40 epochs using a batch size of 128 and the standard ImageNet data augmentations described in Section C.1. We chose 40 epochs because, under a fixed learning rate without decay, the validation losses consistently plateau by this point, indicating that the model has reached its first stage of convergence.

To compute the EGOP matrix for convolutional layers, we repeatedly reinitialize the network weights and record the gradients obtained from a single forward-backward pass at initialization. Because a single pass produces one gradient sample for each kernel in a convolutional layer, a layer with c output channels yields c gradient samples per pass. We therefore determine the required number of passes by ensuring that every convolutional layer accumulates at least 1000 gradient samples in total. Let c_{\min} denote the smallest number of output channels across all convolutional layers. Since this layer contributes only c_{\min} samples per pass, we perform $\lceil 1000/c_{\min} \rceil$ forward-backward passes. In each pass, we compute the gradient of the cross-entropy loss with respect to all convolutional and linear weights.

C.4. Additional Details for Figures in Section 1. For Figure 1 (left), we generate the loss landscape pictured by taking $d = 2$ dimensions and $n = 100$ samples. We generate a matrix A with singular values $\sigma_k = k^{-\alpha}$ with $\alpha = 1.5$ and random right- and left-singular vectors. We sampled $\theta^* \sim \mathcal{N}(0, \mathbb{I})$ and used A, θ^* to induce $f(\cdot)$ a log-sum-exp objective, as defined in Section 6. We generated problem instances at random, and chose a random seed that produced a problem where the primary directions of variation for $f(\cdot)$ were clearly un-aligned with the coordinate axes. We allow Adagrad to take 1000 iterations.

We selected an initial point θ_0 that would produce a clear visual distinction between the two coordinate systems in early iterates. To select the learning rate, we first swept over the values 1, 10, 20, \dots , 100, and examined the suboptimality of the solution produced after 1000 iterates. For each learning rate we conducted 5 random trials on log-sum-exp objectives generated with identical values of d, n, α , and with $\theta^* \sim \mathcal{N}(0, \mathbb{I})$; these trials were initialized at points drawn from $\mathbf{N}(0, \|\theta_0\|_2^2 \mathbb{I})$, where θ_0 is the initial point used in Figure 1 (left). We found that for learning rates 30 through 90, the suboptimality of the solution returned by the algorithm in original coordinates was very comparable across learning rates and very close to zero ($< 1e-10$). Thus we selected a learning rate with these properties that was on the lower end (hence 30) because at larger learning rates, the oscillations around the global minimum made it more difficult to visually assess the difference between trajectories before and after reparameterization.

C.5. Additional Details for Figures in Section 4. Figure 3 plots the eigenspectrum of an empirically estimated EGOP matrix. We use the UCI digits dataset and consider the training subset described in Section C.1. We use the 2-layer ReLU network detailed in Section C.2.

For Figure 3 we use full-batch gradients to estimate the EGOP matrix. We take ρ to be a standard normal distribution: for each gradient sample, we form θ by sampling the

1016 entries of all weights and biases i.i.d. from a standard Gaussian. We estimate $\text{EGOP}(f)$ using
 1017 $M = 10d$, where d is the total number of parameters (sum of all weight and bias entries) for
 1018 the network. This is a larger number of EGOP samples than we use in later experiments; this
 1019 is done intentionally with the goal of clearly resolving the spectral decay, rather than having
 1020 decay appear as an artifact of numerical estimation with few samples.

1021 For Figure 10, we use the same dataset, architecture, objective $f(\cdot)$, and procedure for
 1022 EGOP estimation, but we use different sampling distributions ρ . The “realistic initialization”
 1023 distribution is described in full detail below in Section C.6.

1024 For Figure 11, we use the same architecture and objective for the UCI digits dataset. For
 1025 fashionMNIST, we use the train-validation-test split, as described in Section C.1, and use the
 1026 architecture described in Section C.2.

1027 For Figure 11, we use minibatches of size 300 to estimate gradients. For each architecture,
 1028 we use $M = 5d$ gradient samples to estimate each block EGOP matrix, where d is the number
 1029 of parameters in the network. This is a larger number of EGOP samples than we use in later
 1030 experiments; this is done intentionally with the goal of clearly resolving the spectral decay,
 1031 rather than having decay appear as an artifact of numerical estimation with few samples. We
 1032 perform block EGOP reparameterization for both networks, where the blocks are defined by
 1033 the weight matrices of each network. For more details on block reparameterization for neural
 1034 network weights, see the example in Section D.

1035 When drawing gradients to estimate the block EGOP matrices, the distribution over
 1036 parameters ρ from which we draw is the same as the distribution we later use to initialize
 1037 the networks during training. For a full description of this initialization distribution, see
 1038 Section C.6.

1039 C.6. Additional Details for Figures in Section 6.

1040 *Linear Feedforward Networks.* As described in Section 6, for Figure 6 we consider parameters
 1041 $\theta = [\text{vec}(W_1), \text{vec}(W_2), \text{vec}(W_3)]$ where $W_1 \in \mathbb{R}^{50 \times 10}$, $W_2 \in \mathbb{R}^{30 \times 50}$, $W_3 \in \mathbb{R}^{10 \times 30}$. We seek to
 1042 minimize loss

$$1043 \quad f(\theta) = \|W_3 W_2 W_1 A - Y\|_F^2 / n_{\text{samples}}$$

1044 where $A \in \mathbb{R}^{10 \times n_{\text{samples}}}$, and $Y = M^* A$ for $M^* \in \mathbb{R}^{10 \times 10}$ drawn from a standard Gaussian
 1045 distribution. We induce spectral decay in $\text{EGOP}(f)$ by generating A with singular values
 1046 $\sigma_k(A) = k^{-2}$ and random right- and left-singular vectors. For each trial, we generate 20,000
 1047 data samples, which we split into a test set of 10,000 training data instances, 4,000 validation
 1048 instances, and 6,000 test instances. We define the training loss to be $f(\cdot)$ restricted to the
 1049 training instances and their labels, and use this function when estimating the EGOP matrix.
 1050 We use stochastic mini-batches of size 500 when estimating the EGOP and when performing
 1051 optimization.

1052 For each trial, we generate A , M^* , and Y as described above. We then estimate $\text{EGOP}(f)$
 1053 using $M = 2d$ samples, where d is the total number of parameters in the network ($d = 2, 300$).
 1054 When estimating the EGOP, we let ρ be the same distribution used when initializing networks
 1055 for training. Specifically, we initialize each weight matrix from a mean-zero Xavier normal
 1056 distribution, also called Glorot initialization, which is widespread in practice [15]. The Xavier
 1057 normal distribution is a Gaussian with standard deviation $\sqrt{2/\text{fan-in-fan-out}}$, where for a
 1058 matrix $W \in \mathbb{R}^{n \times m}$, $\text{fan-in-fan-out} = n + m$. We compute the full eigenvalue decomposition to

1059 find the change-of-basis matrix V .

1060 We use 1000 epochs during training. See Section C.3 for choice of hyperparameters and
 1061 choice of the range of learning rates to sweep for tuning. We measure performance using the
 1062 median minimum validation loss achieved during training. We perform 10 independent trials
 1063 at each learning rate. For each algorithm, we choose the learning rate at which the algorithm
 1064 in original achieved lowest median minimum validation loss. Results of this learning rate
 1065 sweep are in Figure 6c.

1066 Because Adam in original coordinates exhibited numerical instability due to the near-zero
 1067 gradient values encountered, we used cosine annealing to decay the learning rate over the
 1068 course of training. We use Pytorch’s default implementation of cosine annealing. For ease
 1069 of comparison, we apply the same learning rate decay schedule to all algorithms (Adagrad,
 1070 Adam, SGD, and SGD with momentum) in both original and reparameterized coordinates.
 1071 We use this learning rate decay schedule throughout learning rate tuning, and in all results
 1072 displayed in Figure 6.

1073 *ReLU Networks for Image Classification.* We use the UCI digits dataset and fashionMNIST,
 1074 with the train-validation-split described in Section C.1 and architectures described in Sec-
 1075 tion C.2. We use the same (random) partition of datapoints into train, validation, and test
 1076 datasets for all trials.

1077 For both networks and objectives, we use minibatches of size 300 to estimate gradients
 1078 during both EGOP estimation and optimization. We perform block EGOP reparameterization
 1079 for both networks, where the blocks are defined by the weight matrices of each network. For
 1080 more details on block reparameterization for neural network weights, see the example in
 1081 Section D. For the UCI digits dataset, we use $M = d$ minibatch gradient samples to estimate
 1082 the EGOP, where d is the number of parameters in the network, and compute the orthonormal
 1083 reparameterization matrices for each layer by taking the (full, deterministic) SVD of the matrix
 1084 of gradient samples for that layer, as this is equivalent to taking the eigenvalue decomposition
 1085 of the empirical block EGOP matrix.

1086 For the fashionMNIST experiments reported in Figure 7, we use pytorch’s `svd_lowrank`
 1087 function to implement randomized SVD, which in turn is based off Algorithms 4.1 and 5.1
 1088 in Halko et al. [18]. This yields $V_r \in \mathbb{R}^{d_1 \cdot d_2 \times r}$ a matrix with orthonormal columns. We use
 1089 this matrix to perform auxiliary variable reparameterization, as introduced in Section 5. See
 1090 Section D for full details. **The following are OLD experimental details. Update to describe
 1091 number of EGOP samples for FashionMNIST. Keep use of randomized SVD discussion.** For
 1092 fashionMNIST, we collect $M = 0.1d$ minibatch gradient samples, where d is the number of
 1093 parameters in the model. For each layer, of dimension $d_1 \times d_2$, we approximate the leading
 1094 $0.1d_1 \cdot d_2$ left-singular vectors of the matrix of gradient samples for that layer.

1095 For each objective and network, we use the same distribution over parameters for draw-
 1096 ing gradient samples to estimate the EGOP and for initializing the network at training
 1097 time. For the experiments in Figure 1 (right), we draw weight entries i.i.d. from a mean-
 1098 zero Xavier distribution, and bias entries i.i.d. from a uniform distribution with range
 1099 $[-(\text{fan-in-fan-out})^{-1/2}, (\text{fan-in-fan-out})^{-1/2}]$.

**Check that these details about initialization and sampling distributions are accurate for
 new auxilliary method. If not, update.** For the experiments in Figure 7, that distribution is
 Pytorch’s default method for re-setting parameters in a linear layer, namely drawing weight

entries i.i.d. from a uniform distribution with range

$$[-(\text{in-features})^{-1/2}, (\text{in-features})^{-1/2}]$$

and drawing bias entries i.i.d. from a uniform distribution with range

$$[-(\text{fan-in-fan-out})^{-1/2}, (\text{fan-in-fan-out})^{-1/2}],$$

1100 where fan-in-fan-out for a layer is defined above in Section C.5. All of these initialization
 1101 distributions stem from heuristics that are widespread in practice, and the distinction between
 1102 the weight distributions used with the digits dataset versus the fashionMNIST dataset was
 1103 the result of using legacy code; we did not tune initializations to different applications.

1104 We selected hyperparameters and tuned learning rates following the procedures discussed
 1105 in Section C.3. We used a constant learning rate for both datasets, as all methods appeared
 1106 stable. We measure performance using the median maximum validation classification accu-
 1107 racy achieved during training. For each algorithm, we choose the learning rate at which the
 1108 algorithm in original achieved highest median maximum validation accuracy. For the digits
 1109 dataset, we performed 50 independent trials at each learning rate. **Check and update param-**
 1110 **eter tuning description for fashionMNIST if needed.** For fashionMNIST, we performed a single
 1111 trial at each learning rate, and only trained for 20 epochs during tuning. Once learning rates
 1112 were selected, we set the number of epochs over which to train by examining loss and accuracy
 1113 curves by eye, and choosing a value large enough that all algorithms had converged. For the
 1114 digits dataset, we trained for 200 epochs, and for fashionMNIST we trained for 80 epochs.

1115 *Regression with Errors-in-Variables.* As described in Section 6, we optimize

$$1116 \quad f(\theta) = \frac{1}{2} \|A(M \odot W) - Y\|_F^2.$$

1117 For each trial, the data are generated as follows. We set the problem dimension to $n = d = 10$
 1118 and sample the design matrix A so that its singular values obey $\sigma_k = k^{-\alpha}$ with $\alpha = 2.0$. For
 1119 $\nu \geq 0$, each entry M_{ij} is drawn independently from $\mathcal{N}(1, \nu^2)$. The ground truth parameter
 1120 W^* is sampled entrywise from a standard normal distribution, and the target Y is set as
 1121 $Y = A(M \odot W^*)$.

1122 For optimization, we compare the performance of SOAP, Shampoo, base Adam and Ada-
 1123 grad, and EGOP-reparameterized Adam and Adagrad. All parameters are initialized at zero,
 1124 and each optimization run proceeds for 1000 steps using a cosine annealing learning rate
 1125 schedule. We employ this learning rate schedule in order to improve SOAP’s convergence. On
 1126 independent trials, a new random instance of A , M , W^* , and Y is generated, on which each
 1127 optimizer is evaluated.

1128 We perform a logarithmic learning rate sweep for each optimizer and each ν . In every
 1129 trial, the training loss trajectory is recorded. We measure performance for each learning rate
 1130 using the median final training loss across the five trials. The best learning rate for each
 1131 setting is taken to be the one with the lowest median final loss.

1132 For EGOP reparameterization, we sample 1000 gradients from parameters W , where the
 1133 entries of W are sampled from a standard normal distribution.

1134 *Low-rank Matrix Factorization.* As stated in Section 6, we consider minimizing

$$1135 \quad f(\theta) = \frac{1}{2} \|A(LR^T) - Y\|_F^2.$$

1136 where $L \in \mathbb{R}^{n \times r}$ and $R \in \mathbb{R}^{m \times r}$, and $A \in \mathbb{R}^{n \times n}$. On each trial, the synthetic ground truth
 1137 X^* is generated to have rank r and a prescribed condition number κ . Specifically, we factor
 1138 X^* as $X^* = U\Sigma V^\top$ with random orthonormal matrices $U \in \mathbb{R}^{n \times r}$, $V \in \mathbb{R}^{m \times r}$, and singular
 1139 values $\Sigma = \text{diag}(\sigma_1, \dots, \sigma_r)$, for values forming a linear grid from 1 to $1/\kappa$. The measurement
 1140 operator A is also constructed as a random matrix with controlled singular value decay: its
 1141 singular values follow $\sigma_k = k^{-\alpha}$. Both X^* and A are independently regenerated for each trial.
 1142 Each observation matrix Y was generated using $Y = AX^* + \mathcal{N}(0, \sigma_{\text{obs}}^2)$. All experiments are
 1143 performed with $n = m = 20$, $r = 5$, $\kappa = 2$, and $\alpha = 1/2$, and the standard deviation of the
 1144 observation noise was set to $\sigma_{\text{obs}} = 0.01$.

1145 We use spectral initialization: to form L_0, R_0 , we solve for $A^\dagger Y$ where A^\dagger indicates the
 1146 Moore-Penrose pseudoinverse, perform an SVD, and truncate to rank r to obtain U_r, Σ_r, V_r ,
 1147 and form left- and right- factors $L_0 = U_r \Sigma_r^{-1/2}$, $R_0 = V_r \Sigma_r^{-1/2}$. SOAP, Shampoo, and
 1148 base Adam and Adagrad are initialized at (L_0, R_0) , and EGOP-reparameterized methods are
 1149 initialized at the appropriate equivalent points in the coordinates of the EGOP-eigenbasis.

1150 For each optimizer, we perform logarithmic learning rate sweeps, performing five indepen-
 1151 dent trials on each. For each optimizer, we selected the learning rate achieving the lowest
 1152 median final loss. All parameters are optimized for 1000 steps starting from this initializa-
 1153 tion. Figure 9 reports results aggregated over 5 independent trials, corresponding to problem
 1154 instances *independent* from any problems generated during the learning rate sweep.

1155 *Convex objectives.* We study three objectives. For each, we generate a matrix $A \in$
 1156 $\mathbb{R}^{n_{\text{samples}} \times d}$ with singular value decay $\sigma_k = k^{-\alpha}$ and random orthonormal right- and left-
 1157 singular vectors. The values of α are specified in the caption of Figure 17. The log-sum-exp
 1158 objective is defined in Eq. B.1. The logistic regression objective is defined as

$$1159 \quad f(\theta) = \sum_{i=1}^{n_{\text{samples}}} -y_i \log \left(\frac{1}{1 + e^{-\langle a_i, \theta \rangle}} \right) - (1 - y_i) \log \left(\frac{e^{-\langle a_i, \theta \rangle}}{1 + e^{-\langle a_i, \theta \rangle}} \right)$$

1160 where $\{a_i \in \mathbb{R}^d\}_{i=1}^{n_{\text{samples}}}$ are the columns of A , and labels $y_i \sim \text{Bernoulli}(\pi(\theta^*)_i)$ where $\theta^* \in \mathbb{R}^d$
 1161 is drawn from $\mathcal{N}(0, \mathbb{I})$ and

$$1162 \quad \pi(\theta^*)_i \stackrel{\text{def}}{=} \frac{e^{-\langle a_i, \theta^* \rangle}}{1 + e^{-\langle a_i, \theta^* \rangle}}.$$

1163 The logistic regression objective is

$$1164 \quad f(\theta) = \frac{1}{2} \|A\theta - y\|_2^2$$

1165 where $y = A\theta^*$ and $\theta^* \in \mathbb{R}^d$ is drawn from $\mathcal{N}(0, \mathbb{I})$.

1166 For log-sum-exp and linear least squares, $n_{\text{samples}} = d = 100$. For logistic regression,
 1167 $n_{\text{samples}} = 100$ and $d = 3$. For each objective, on each trial we generate A, θ^* , and thus $f(\cdot)$
 1168 randomly and independently following the above procedure.

1169 For all convex objectives, use deterministic (full-batch) gradients to estimate the EGOP
 1170 and to optimize. For all objectives, we use $M = 5d$ gradient samples to estimate the EGOP
 1171 matrix. We take ρ to be a standard Gaussian distribution, and use the same distribution to
 1172 initialize when optimizing. See Section C.3 for choice of hyperparameters and choice of the
 1173 range of learning rates to sweep for tuning. We used a constant learning rate for all convex
 1174 objectives. We measure performance using the median final training loss achieved. For each
 1175 algorithm, we choose the learning rate at which the algorithm in original achieved lowest
 1176 median training loss. We perform 5 independent trials at each learning rate.

1177 For each objective, we optimize for 1000 iterations, or until the gradient Euclidean norm
 1178 drops below $1e - 8$. The latter termination condition is comparable to default termination
 1179 conditions employed by `scipy.optimize` [36].

1180 **Appendix D. Expanded Discussion of Heuristics for Scalability.** In this section, we
 1181 expand on the heuristics proposed in Section 5.

1182 **Block Reparameterization.** Here we instantiate block reparameterization for multilayer
 1183 neural networks. Given an L -layer neural network whose i th layer is parameterized by
 1184 weight matrix $W_i \in \mathbb{R}^{n_{\text{in}}^{(i)} \times n_{\text{out}}^{(i)}}$ and bias vector $b_i \in \mathbb{R}^{n_{\text{out}}^{(i)}}$, we consider optimizing $f(\theta) =$
 1185 $\text{loss}(\theta; X_{\text{train}}, y_{\text{train}})$ with respect to parameters $\theta = [(W_1, b_1), \dots, (W_L, b_L)]$.

1186 Rather than forming the full EGOP as described in Algorithm 2.1, we consider estimating
 1187 the *layer* EGOP matrices

$$1188 \quad \text{EGOP}^{(i)} \stackrel{\text{def}}{=} \frac{1}{M} \sum_{k=1}^M \nabla_{W_i} f(\theta_k) \nabla_{W_i} f(\theta_k)^\top$$

1189 where $\nabla_{W_i} f(\theta_k) \in \mathbb{R}^{n_{\text{in}}^{(i)} n_{\text{out}}^{(i)}}$ is the vector of partial derivatives of f w.r.t. the entries of W_i
 1190 evaluated at θ_k , and the points $\{\theta_k\}_{k=1}^M \sim \rho$ are drawn i.i.d..

1191 Given these L empirical layer EGOP matrices, one can obtain L change-of-basis matrices
 1192 $V^{(i)} \in \mathbb{R}^{n_{\text{in}}^{(i)} n_{\text{out}}^{(i)} \times n_{\text{in}}^{(i)} n_{\text{out}}^{(i)}}$, and form the objective

$$1193 \quad \tilde{f}(\theta) = f([(V^{(1)}W_1, b_1), \dots, (V^{(L)}W_L, b_L)]).$$

1194 Block EGOP reparameterization requires storing and applying L orthonormal matrices,
 1195 each of size $n_{\text{in}}^{(i)} n_{\text{out}}^{(i)} \times n_{\text{in}}^{(i)} n_{\text{out}}^{(i)}$. For deep neural networks, this can be considerably less ex-
 1196 pensive than storing and applying the global change-of-basis matrix, which would be of size
 1197 $\sum_{i=1}^L n_{\text{in}}^{(i)} n_{\text{out}}^{(i)} \times \sum_{i=1}^L n_{\text{in}}^{(i)} n_{\text{out}}^{(i)}$. This may also reduce the sampling cost; each layer EGOP is
 1198 of smaller dimension than the global EGOP and thus has a smaller number of eigenvectors
 1199 to estimate. Thus the layer EGOP eigenbases may be estimated with fewer total gradient
 1200 samples.

1201 Another benefit to block EGOP reparameterization is that it offers an easy way to reduce
 1202 cost by choosing to only reparameterize a subset of layers. Thus for example if a network
 1203 has a subset of layers which are too wide to efficiently reparameterize, one can choose to only
 1204 reparameterize the narrower layers.

1205 For networks where the first layer involves a linear transformation of the input data,
 1206 reparameterizing only the first layer corresponds to an orthogonal transformation of the data.

1207 Thus pre-computing the EGOP eigenbasis and applying this matrix to the input data up-front
 1208 would allow one to reparameterize the first layer, without requiring one to store and apply
 1209 the change-of-basis matrix during training.

1210 **D.1. Efficient Heuristics for Large-Scale Neural Networks.** One major application area
 1211 for adaptive gradient methods is training of massive multi-layer neural networks. In this
 1212 section, we present computational heuristics that enable EGOP-reparameterization to scale
 1213 to massive architectures.

1214 [AD : Rick and Jose: you can both access commands that allow you to leave notes like
 1215 this. You can change your color by going to TMLR/TMLR_custom_commands.sty and editing
 1216 your command.] Jose: This is an example of Jose’s note command Rick: This is an example
 1217 of Rick’s note command

1218 **D.1.1. Reparameterizing Massive Fully-Connected Layers.** For fully-connected layers
 1219 with a large number of parameters, storing and applying the full reparameterization matrix
 1220 $V \in \mathbb{R}^{d \times d}$ can be prohibitively expensive. We describe two memory-efficient variants that
 1221 compute and retain only the top- r eigenvectors of the EGOP matrix, forming a reduced basis
 1222 $V_r \in \mathbb{R}^{d \times r}$ where $r \ll d$. The first approach constrains parameters to the subspace spanned
 1223 by V_r , while the second uses auxiliary variables to maintain full expressivity while still storing
 1224 only V_r .

1225 *Reduced-dimension EGOP.* The primary motivation for reduced-dimension reparameteri-
 1226 zation is to reduce the memory cost of storing the reparameterization matrix V . Storing the
 1227 full orthonormal basis $V \in \mathbb{R}^{d \times d}$ requires $O(d^2)$ memory, while storing only the top- r eigen-
 1228 vectors $V_r \in \mathbb{R}^{d \times r}$ requires $O(dr)$ memory. Since the memory savings from reducing V scale as
 1229 $d^2 - dr = d(d-r)$ while the savings from reducing the parameter count scale as $d-r$, shrinking
 1230 V to V_r provides a factor of d larger benefit than shrinking θ to θ_r . Additionally, computing
 1231 only the top- r eigenvectors via randomized SVD is significantly faster than computing the full
 1232 eigendecomposition.

Rather than the full reparameterization $f(\theta) = f(V\theta)$ with $\theta \in \mathbb{R}^d$ and $V \in \mathbb{R}^{d \times d}$,
 reduced-dimension EGOP uses

$$f(\theta) = f(V_r \theta_r)$$

1233 where $\theta_r \in \mathbb{R}^r$ and $V_r \in \mathbb{R}^{d \times r}$ contains the top- r eigenvectors of the empirical EGOP matrix
 1234 \hat{P} . This constrains θ to lie in the r -dimensional subspace $\text{span}(V_r) \subset \mathbb{R}^d$.

1235 *Applicability and limitations.* In practice, we compute V_r by first sampling M gradients
 1236 $\{g_i\}_{i=1}^M$ at random initializations $\{\theta_i\}_{i=1}^M \sim \rho$, forming the gradient matrix $G \in \mathbb{R}^{d \times M}$ with
 1237 columns g_i . We then apply randomized SVD to G to obtain the top- r left singular vectors,
 1238 which form V_r . This approach avoids forming the full $\hat{P} = \frac{1}{M} G G^T$ matrix and is computa-
 1239 tionally cheaper than full eigendecomposition.

1240 *Computing the reduced basis.* If the EGOP matrix V has exact rank r , then V_r spans the
 1241 full column space of V and we lose no expressivity by switching to V_r . However, in practice,
 1242 the EGOP matrix is typically full-rank with rapidly decaying eigenvalues rather than having
 1243 exact rank r . Thus reduced-dimension reparameterization necessarily introduces a projection
 1244 error of $\|\theta - V_r V_r^T \theta\| > 0$ for general $\theta \in \mathbb{R}^d$, limiting the expressivity of the reparameterized
 1245 model to the r -dimensional subspace $\text{span}(V_r)$. We introduce auxiliary variables EGOP below

1246 to overcome this limitation.

Auxiliary variables EGOP. Rather than completing the basis by appending a random orthonormal basis for $\text{span}(V_r)^\perp$, we introduce auxiliary parameters to represent the orthogonal complement. The key idea is to decompose parameters into two orthogonal components via the reparameterized objective:

$$\tilde{f}(\theta_r, \theta_d) = f(V_r \theta_r + (I - V_r V_r^T) \theta_d)$$

where $\theta_r \in \mathbb{R}^r$ represents coordinates in the dominant r -dimensional EGOP subspace, and $\theta_d \in \mathbb{R}^d$ are auxiliary variables whose projection $(I - V_r V_r^T) \theta_d$ represents the residual component orthogonal to this subspace. The full parameter vector in \mathbb{R}^d is thus reconstructed as:

$$\theta = V_r \theta_r + (I - V_r V_r^T) \theta_d$$

1247 where $V_r V_r^T$ is the orthogonal projection onto $\text{span}(V_r)$. The explicit projection $(I - V_r V_r^T)$
 1248 applied to θ_d ensures the residual component remains in $\text{span}(V_r)^\perp$ throughout optimization.
 1249 Although the parameter count increases from r to $r + d$ relative to the reduced-dimension
 1250 approach, memory cost remains dominated by V_r storage (which scales as $O(dr)$) rather than
 1251 parameter storage (which scales as $O(d)$), making this trade-off favorable in practice.

1252 *Exact reconstruction.* To initialize the auxiliary variables model equivalently to an original
 1253 model at θ_0 , we set:

$$\begin{aligned} 1254 \quad \theta_r &= V_r^T \theta_0 \in \mathbb{R}^r \\ 1255 \quad \theta_d &= (I - V_r V_r^T) \theta_0 \in \mathbb{R}^d \end{aligned}$$

This initialization guarantees that the reparameterized model exactly reconstructs θ_0 :

$$V_r \theta_r + (I - V_r V_r^T) \theta_d = \theta_0$$

1256 ensuring functional equivalence: $f(\theta_0) = \tilde{f}(\theta_r, \theta_d)$ for all inputs.

Algorithm D.1 EGOP Auxiliary Variables Linear Layer Forward Pass

- 1: **Input:** $x \in \mathbb{R}^{n \times d_{\text{in}}}$ (batch input), $\theta_r \in \mathbb{R}^r$ (reduced params), $\theta_d \in \mathbb{R}^d$ (auxiliary params),
 $V \in \mathbb{R}^{d \times r}$ (EGOP basis)
 - 2: **Output:** $y \in \mathbb{R}^{n \times d_{\text{out}}}$
 - 3: **Layer dims:** $d_{\text{out}}, d_{\text{in}}$ where $d = d_{\text{out}} \times d_{\text{in}}$
 - 4: {Compute $\theta = V_r \theta_r + (I - V_r V_r^T) \theta_d$ using only 2 matrix multiplications}
 - 5: $\theta_{\text{full}} \leftarrow \theta_d + V(\theta_r - V^T \theta_d) \setminus [d] = [d] + [d, r] \times [r]$
 - 6:
 - 7: $\setminus \setminus$ Reshape to weight matrix
 - 8: $W \leftarrow \text{reshape}(\theta_{\text{full}}, (d_{\text{out}}, d_{\text{in}})) \setminus [d_{\text{out}}, d_{\text{in}}]$
 - 9:
 - 10: $\setminus \setminus$ Apply linear transformation
 - 11: $y \leftarrow x W^T + b \setminus [n, d_{\text{out}}] = [n, d_{\text{in}}] \times [d_{\text{in}}, d_{\text{out}}]$
 - 12:
 - 13: **Return:** y
-

1257 **D.1.2. Reparameterizing Convolutional Network Layers.** To evaluate the scalability and
 1258 generalization of EGOP reparameterization beyond fully connected architectures, we extend
 1259 the method to convolutional neural networks (CNNs) and residual networks (ResNets). We
 1260 study EGOP reparameterization on larger and deeper models such as AlexNet and ResNet
 1261 trained on ImageNet, and compare performance against well tuned baseline optimizers. This
 1262 extension allows us to assess whether EGOP can improve optimization in modern large scale
 1263 visual recognition pipelines.

1264 *Background on convolutional layers.* A convolutional layer consists of a collection of learn-
 1265 able filters (also called kernels). Each filter is a three-dimensional tensor of shape $k \times k \times C_{\text{in}}$,
 1266 where k is the spatial size and C_{in} is the number of input channels. A layer with C_{out} output
 1267 channels therefore contains C_{out} such filters, each producing one output channel by convolv-
 1268 ing the filter with the input feature map. Stacking the C_{out} resulting feature maps yields the
 1269 layer’s activation tensor.

1270 *Motivation for EGOP in convolutional networks.* A direct application of EGOP to convolu-
 1271 tional layers would require constructing and applying a separate reparameterization matrix
 1272 for every filter, which is prohibitive for modern CNNs. However, filters within the same layer
 1273 are drawn from the same initialization distribution, so their first-step gradient statistics are
 1274 identically distributed; consequently, their EGOP matrices at early stages of training may
 1275 also be similar. This motivates our shared reparameterization strategy: we compute a single
 1276 EGOP matrix per convolutional layer and reuse it across all filters within that layer. This
 1277 reduces computational cost while preserving the benefits of EGOP, enabling its use in deep
 1278 architectures such as AlexNet and ResNet.

1279 **Adela:** The ideas in these next two paragraph are correct, but it might be confusing for
 1280 a reader to digest. Let’s think about how we could make this more tractable; we could try to
 1281 workshop the writing. Maybe some little graphic/block diagram could help explain waht we
 1282 mean?

1283 *Shared Gradient Structure at Initialization.* A key observation that enables an efficient exten-
 1284 sion of EGOP to CNNs is that, for most practical initialization distributions (e.g. Kaiming,
 1285 Xavier, Glorot normal/uniform, filters within the same convolutional layer are identically
 1286 distributed and so are their gradients. This property allows us to substantially reduce the
 1287 computational cost of the reparameterization. During initialization, each forward/backward
 1288 pass yields multiple gradients simultaneously rather than a single one, which decreases the
 1289 number of required passes proportionally to the number of filters. Once these gradients are
 1290 collected, we compute a single orthogonal matrix V that captures the shared structure of the
 1291 layer. This V is then reused across all filters, leading to a dramatic reduction in memory
 1292 footprint and computational speedup, due to the low number of required forward/backward
 1293 passes required for gradient sampling.

1294 **Jose:** Should I show plots here? **Adela:** Time permitting. We could include a plot
 1295 comparing “shared V ” versus “different V ” on some small model (cifar10) to show how sharing
 1296 V across filters impacts performance.

1297 *Kronecker Product Structure from Reusing V .* As referenced in Section 5, reusing V across
 1298 all filters as described above is equivalent to employing a Kronecker-product reparameter-
 1299 ization matrix. Moreover, the effective Kronecker product imposed by this construction is
 1300 different from that employed by SOAP/Shampoo, and the effective EGOP change-of-basis

1301 matrix applied to each convolutional filter does *not* have any structural constraint. We clarify
 1302 this point below.

1303 In the EGOP reparameterization procedure described in the preceding paragraphs, for
 1304 each convolutional layer we form $V \in \mathbb{R}^{(k^2 C_{\text{in}}) \times (k^2 C_{\text{in}})}$ and apply this change-of-basis to each
 1305 (vectorized) filter in the convolutional layer, denoted θ_j for $j = 1 \dots C_{\text{out}}$, as

$$1306 \quad V \text{vec}(\theta_j).$$

1307 This is equivalent to employing the following Kronecker-product structured change-of-basis

$$1308 \quad (\mathbb{I}_{C_{\text{out}}} \otimes V) \begin{bmatrix} \text{vec}(\theta_1) \\ \vdots \\ \text{vec}(\theta_{C_{\text{out}}}) \end{bmatrix}.$$

1309 We now compare with the Kronecker-product structured change-of-basis employed by
 1310 SOAP/Shampoo. For a convolutional layer with the above-described dimensions, SOAP and
 1311 Shampoo form four orthogonal matrices

$$1312 \quad Q_h \in \mathbb{R}^{k \times k}, \quad Q_w \in \mathbb{R}^{k \times k}, \quad Q_{\text{in}} \in \mathbb{R}^{C_{\text{in}} \times C_{\text{in}}}, \quad Q_{\text{out}} \in \mathbb{R}^{C_{\text{out}} \times C_{\text{out}}}.$$

1313 The SOAP/Shampoo forward pass is then equivalent to performing

$$1314 \quad (Q_{\text{out}} \otimes Q_w \otimes Q_h \otimes Q_{\text{in}}) \begin{bmatrix} \text{vec}(\theta_1) \\ \vdots \\ \text{vec}(\theta_{C_{\text{out}}}) \end{bmatrix}.$$

1315 Comparing the transformation employed in EGOP reparameterization with that employed
 1316 by SOAP/Shampoo, we see that the effective transformation applied to each individual filter
 1317 by SOAP/Shampoo is constrained to have a Kronecker product structure, $Q_w \otimes Q_h \otimes Q_{\text{in}}$,
 1318 whereas the transformation V employed in EGOP reparameterization is a generic matrix.
 1319 SOAP/Shampoo employ a more general mapping over output channels, namely Q_{out} , com-
 1320 pared to $\mathbb{I}_{C_{\text{out}}}$ employed in EGOP reparameterization. However, as explained in the preceding
 1321 paragraphs, the distribution over different filters is identical at initialization, and in practice,
 1322 we find employing the identity mapping over output channels in EGOP reparameterization
 1323 yields performance competitive with that of SOAP/Shampoo.

1324 **Algorithmic Implementation.** The convolutional reparameterization follows the same princi-
 1325 ple as in the linear case, with additional tensor reshaping steps to accommodate convolutional
 1326 filters. For each layer, we simultaneously collect gradients at initialization and flatten them
 1327 into a single matrix. A full SVD is performed for each layer, and the resulting right singular
 1328 vectors V define a shared orthogonal basis that captures the dominant gradient directions.
 1329 Each filter is then flattened, preconditioned by V , reshaped back to its convolutional form,
 1330 and used in the standard forward pass. The detailed procedure is provided below.

Algorithm D.2 EGOP Reparameterization for Convolutional Layers

```

1: Input: CNN forward map  $f$  with convolutional layers  $\{W^{(l)}\}_{l=1}^L$ , number of initializations
    $K$ , sampling distribution  $\rho$ 
2: Output: Reparameterized forward map  $\tilde{f}$ 
3:
4: \1. Sample random initializations and collect gradients
5: Sample  $\{\theta_k\}_{k=1}^K \sim \rho$  i.i.d.
6: for  $k = 1, \dots, K$  do
7:   Initialize network parameters with  $\theta_k$ 
8:   Perform one forward/backward pass
9:   Store gradients  $\{\nabla W^{(l,k)}\}_{l=1}^L$ 
10: end for
11:
12: \2. Form per-layer gradient matrices
13: for each convolutional layer  $l = 1, \dots, L$  do
14:   \Flatten each kernel gradient  $\nabla W_j^{(l,k)} \in \mathbb{R}^{c_{\text{in}} \times k_h \times k_w}$  into a vector
15:   \Stack across initializations  $k$  and output channels  $j$ 
16:    $G^{(l)} \leftarrow \text{stack}\left(\text{flatten}(\nabla W_j^{(l,k)})\right)_{k=1, \dots, K; j=1, \dots, c_{\text{out}}}$  \  $G^{(l)} \in \mathbb{R}^{(c_{\text{in}} k_h k_w) \times (K c_{\text{out}})}$ 
17: end for
18:
19: \3. Compute EGOP eigenbasis per layer
20: for each layer  $l = 1, \dots, L$  do
21:    $V^{(l)} \leftarrow \text{left\_singular\_vectors}(G^{(l)})$ 
22: end for
23:
24: \4. Define reparameterized forward map
25: Define  $\tilde{f}(\{W_j^{(l)}\}_{j=1}^L) = f\left(\left\{\text{reshape}\left(V^{(l)} \text{flatten}(W_j^{(l)}), (c_{\text{in}}, k_h, k_w)\right)\right\}_{j=1}^L\right)$ 
26:
27: \5. Optimize the reparameterized model
28: Optimize  $\tilde{f}$  using an adaptive optimizer of choice
29: return  $\tilde{f}$ 

```

1331 **Appendix E. Deferred Proofs.**

1332 **E.1. Deferred Proofs from Section A.1.** In this section, we will leverage the following
1333 basic results from linear algebra, whose proofs we include for completeness.

1334 **Lemma E.1.** Given $A \in \mathbb{R}^{d \times d}$ a symmetric PSD matrix, for any vector $v \in \mathbb{R}^d$,

$$1335 \quad \langle v, A^\top A v \rangle \leq \lambda_{\max}(A) \langle v, A v \rangle$$

1336 where $\lambda_{\max}(A)$ denotes the largest eigenvalue of A .

1337 *Proof of Lemma E.1.* Since A is PSD, it admits a square root $A^{1/2}$ which is itself PSD.
 1338 Consequently,

$$1339 \quad \langle Av, Av \rangle = \|Av\|^2 \leq \|A^{1/2}\|_{\text{op}}^2 \|A^{1/2}v\|^2 = \lambda_{\max}(A) \langle v, Av \rangle. \quad \blacksquare$$

1340 **Lemma E.2 (Order-preserving square root).** For any positive-semidefinite matrices X, Y ,
 1341 the following holds:

$$1342 \quad (\text{E.1}) \quad X \preceq Y \implies X^{1/2} \preceq Y^{1/2}.$$

1343 *Proof of Lemma A.7.* Because $\tilde{\rho}(\tilde{\theta}) = \rho(V\tilde{\theta})$,

$$1344 \quad \mathbb{E}_{\tilde{\theta} \sim \tilde{\rho}}[q(\tilde{\theta})] = \mathbb{E}_{\theta \sim \rho}[q(V^\top \theta)]$$

1345 By the chain rule,

$$1346 \quad \nabla \tilde{f}(\tilde{\theta}) = V^\top \nabla f(V\tilde{\theta}).$$

1347 Thus

$$1348 \quad \mathbb{E}_{\tilde{\theta} \sim \tilde{\rho}}[\nabla \tilde{f}(\tilde{\theta}) \nabla \tilde{f}(\tilde{\theta})^\top] = V^\top \mathbb{E}_{\tilde{\theta} \sim \tilde{\rho}}[\nabla f(V\tilde{\theta}) \nabla f(V\tilde{\theta})^\top] V = V^\top \mathbb{E}_{\theta \sim \rho}[\nabla f(\theta) \nabla f(\theta)^\top] V = \Lambda. \quad \blacksquare$$

1349 *Proof of Lemma E.2.* Since $Y - X \succeq 0$, any unit vector v satisfies

$$1350 \quad 0 \leq \langle v, (Y - X)v \rangle = \langle v, (Y^{1/2} + X^{1/2})(Y^{1/2} - X^{1/2})v \rangle.$$

1351 Since X, Y are PSD, their square roots are well-defined and symmetric. In particular, the
 1352 matrix $Y^{1/2} - X^{1/2}$ is symmetric and, as such, has real eigenvalues. Therefore, it suffices to
 1353 prove all eigenvalues of this matrix are nonnegative.

1354 To that end, let (λ, v) be an eigenpair of $Y^{1/2} - X^{1/2}$. We have

$$1355 \quad (\text{E.2}) \quad 0 \leq \langle (Y^{1/2} - X^{1/2})v, (Y^{1/2} + X^{1/2})v \rangle = \lambda \langle v, (Y^{1/2} + X^{1/2})v \rangle.$$

1356 In particular, assume $\lambda < 0$ (since otherwise there is nothing to show). Consequently,

$$1357 \quad 0 > \lambda = \langle v, (Y^{1/2} - X^{1/2})v \rangle \implies \langle v, X^{1/2}v \rangle > \langle v, Y^{1/2}v \rangle \geq 0,$$

1358 where the last inequality follows from the fact that Y is PSD. In particular, this implies
 1359 that $\langle v, (Y^{1/2} + X^{1/2})v \rangle > 0$; combined with (E.2) and the assumption $\lambda < 0$, this yields
 1360 a contradiction. Therefore, $\lambda \geq 0$. Since the choice of eigenpair was arbitrary, the claim
 1361 follows. \blacksquare

1362 Given these lemmas, we are now equipped to present full proofs of the results presented
 1363 in Section A.1.

1364 *Proof of Lemma A.8.* An intermediate result in the proof of Thm. 4.1 in Liu et al. [25]
 1365 implies that in the case of noise-free gradients, i.e. oracle gradient $g(\cdot) = \nabla \tilde{f}(\cdot)$, Adagrad with
 1366 constraint set $\tilde{\Theta}$ produces iterates such that

$$1367 \quad \frac{1}{T} \sum_{t=0}^{T-1} (\tilde{f}(\tilde{\theta}_t) - \tilde{f}(\theta^*)) \leq \left(\eta + \frac{\tilde{D}_\infty^2}{\eta} \right)^2 \cdot \frac{\|\tilde{L}\|_1}{T} + \frac{\epsilon \tilde{D}_2^2}{\eta T}.$$

1368 We note $D_2^2 = \tilde{D}_2^2$ because the Euclidean norm is invariant under orthonormal transformation.
 1369 Because the constraint update implies $\{\tilde{\theta}_t\} \subseteq \tilde{\Theta}$, applying Thm. A.1 implies

$$1370 \quad \|\vec{\tilde{L}}\|_1 \leq \|\vec{L}\|_1 \cdot \sqrt{1 + \delta} \left(\frac{2 \text{sr}(f)}{d(\beta_1 - \delta)} + 2 \frac{\sqrt{\delta}}{\beta_1 - \delta} \left(1 + \frac{2}{\sqrt{d}} \right) \right)$$

1371 which immediately implies the result. ■

1372 *Proof of Lemma A.5.* By definition, f is twice-differentiable if and only if it satisfies

$$1373 \quad (\text{E.3}) \quad f(y) = f(x) + \langle \nabla f(x), y - x \rangle + \langle y - x, \nabla^2 f(x)(y - x) \rangle + o(\|y - x\|^2), \quad \text{for all } x, y \in \mathbb{R}^d.$$

1374 We first prove the ‘‘if’’ version: assume $|\langle v, \nabla^2 f(\theta)v \rangle| \leq \langle v, \text{diag}(L)v \rangle$. Consider any $x, y \in \Theta$.
 1375 By the mean value theorem,

$$\begin{aligned} 1376 \quad f(y) - f(x) - \langle \nabla f(x), y - x \rangle &= \int_0^1 t \langle y - x, \nabla^2 f(x + t(y - x))(y - x) \rangle dt \\ 1377 &\leq \int_0^1 t \langle y - x, \text{diag}(L)(y - x) \rangle dt \\ 1378 &= \langle y - x, \text{diag}(L)(y - x) \rangle, \end{aligned}$$

1379 where the inequality follows by assumption and the convexity of Θ . Taking absolute values
 1380 on both sides yields the claim.

1381 We now prove the ‘‘only if’’ part. In particular, for any $x, y \in \Theta$, writing $y = x + tv$ and
 1382 invoking (E.3) we obtain

$$\begin{aligned} 1383 \quad |\langle y - x, \nabla^2 f(x)(y - x) \rangle| &\leq |f(y) - f(x) - \langle \nabla f(x), y - x \rangle| + o(\|y - x\|^2) \\ 1384 &\leq \langle y - x, \text{diag}(L)(y - x) \rangle + o(\|y - x\|^2) \\ 1385 \quad \Leftrightarrow t^2 |\langle v, \nabla^2 f(x)v \rangle| &\leq t^2 \langle v, \text{diag}(L)v \rangle + o(t^2). \quad \blacksquare \end{aligned}$$

1386 Because t was chosen such that $y = x + tv$ for $x, y \in \Theta$, convexity of Θ implies that $x + \tau v \in \Theta$
 1387 for all $\tau \in [0, t]$. Thus dividing both sides by t and letting $t \downarrow 0$ implies the result.

1388 *Proof of Lemma A.6.* Given $f(\cdot)$ satisfying Assumption 2, $\forall \theta_1, \theta_2 \in \Theta$

$$1389 \quad \|\nabla f(\theta_1) - \nabla f(\theta_2) - \nabla^2 f(\theta_2)(\theta_1 - \theta_2)\|_2 \leq H \|\theta_1 - \theta_2\|_2^2.$$

1390 For fixed θ^* the local minimum in Assumption 1, this implies that $\forall \theta \in \Theta \exists R_f(\theta) \in \mathbb{R}^d$ such
 1391 that

$$1392 \quad \nabla f(\theta) = \nabla f(\theta^*) + \nabla^2 f(\theta^*)(\theta - \theta^*) + R_f(\theta) \quad \text{and} \quad \|R_f(\theta)\|_2 \leq H \|\theta - \theta^*\|_2^2.$$

1393 Given θ^* a stationary point of $f(\cdot)$, $\nabla f(\theta^*) = 0$, so the above simplifies to

$$1394 \quad \nabla f(\theta) = \nabla^2 f(\theta^*)(\theta - \theta^*) + R_f(\theta).$$

1395 We can substitute the above expression into the definition of the EGOP matrix:

$$\begin{aligned}
 1396 \quad \mathbb{E}_{\theta \sim \rho}[\nabla f(\theta) \nabla f(\theta)^\top] &= \mathbb{E}_{\theta \sim \rho} \left[(\nabla^2 f(\theta^*)(\theta - \theta^*) + R_f(\theta)) (\nabla^2 f(\theta^*)(\theta - \theta^*) + R_f(\theta))^\top \right] \\
 1397 \quad &= \mathbb{E}_{\theta \sim \rho} \left[\nabla^2 f(\theta^*)(\theta - \theta^*)(\theta - \theta^*)^\top \nabla^2 f(\theta^*)^\top \right] + E_f(\theta^*)
 \end{aligned}$$

1398 where

$$1399 \quad E_f(\theta^*) \stackrel{\text{def}}{=} \mathbb{E}_{\theta \sim \rho} [R_f(\theta)(\theta - \theta^*)^\top \nabla^2 f(\theta^*)^\top + \nabla^2 f(\theta^*)(\theta - \theta^*)R_f(\theta)^\top + R_f(\theta)R_f(\theta)^\top].$$

1400 We can simplify the first term by leveraging the fact that ρ is isotropic about θ_c , as
 1401 specified in Assumption 1:

$$\begin{aligned}
 1402 \quad \mathbb{E}_{\theta \sim \rho} \left[\nabla^2 f(\theta^*)(\theta - \theta^*)(\theta - \theta^*)^\top \nabla^2 f(\theta^*)^\top \right] &= \nabla^2 f(\theta^*) \mathbb{E}_{\theta \sim \rho} \left[(\theta - \theta^*)(\theta - \theta^*)^\top \right] \nabla^2 f(\theta^*)^\top \\
 1403 \quad &= \nabla^2 f(\theta^*) \mathbb{E}_{\theta \sim \rho} \left[(\theta - \theta_c + (\theta_c - \theta^*))(\theta - \theta_c + (\theta_c - \theta^*))^\top \right] \nabla^2 f(\theta^*)^\top \\
 1404 \quad &= \nabla^2 f(\theta^*) \left(\mathbb{E}_{\theta \sim \rho} \left[(\theta_c - \theta)(\theta_c - \theta)^\top \right] + (\theta_c - \theta^*)(\theta_c - \theta^*)^\top \right) \nabla^2 f(\theta^*)^\top \\
 1405 \quad &= \nabla^2 f(\theta^*) \left(c^2 \mathbb{I} + (\theta_c - \theta^*)(\theta_c - \theta^*)^\top \right) \nabla^2 f(\theta^*)^\top
 \end{aligned}$$

1406 where the penultimate equality follows from the assumption that ρ has mean θ_c and the
 1407 last equality follows from the assumption that ρ is isotropic about θ_c . Because θ^* is a local
 1408 minimum, $\nabla^2 f(\theta^*)$ is PSD, and further the rank-1 outer product is also PSD, so we have

$$1409 \quad \nabla^2 f(\theta^*) \left(c^2 \mathbb{I} + (\theta_c - \theta^*)(\theta_c - \theta^*)^\top \right) \nabla^2 f(\theta^*)^\top \succeq c^2 \nabla^2 f(\theta^*) \nabla^2 f(\theta^*)^\top.$$

1410 Moreover by Assumption 1, $\|\theta_c - \theta^*\|_2^2 \leq c^2$, and thus

$$1411 \quad (\theta_c - \theta^*)(\theta_c - \theta^*)^\top \preceq c^2 \mathbb{I}.$$

1412 This implies

$$1413 \quad \nabla^2 f(\theta^*) \left(c^2 \mathbb{I} + (\theta_c - \theta^*)(\theta_c - \theta^*)^\top \right) \nabla^2 f(\theta^*)^\top \preceq 2c^2 \nabla^2 f(\theta^*) \nabla^2 f(\theta^*)^\top.$$

1414 Combining these two matrix inequalities yields the matrix inequalities for

$$1415 \quad G(\theta^*) \stackrel{\text{def}}{=} \nabla^2 f(\theta^*) \left(c^2 \mathbb{I} + (\theta_c - \theta^*)(\theta_c - \theta^*)^\top \right) \nabla^2 f(\theta^*)^\top.$$

1416 Lastly, we bound $|\langle v, E_f(\theta^*)v \rangle|$ for $E_f(\theta^*)$ defined above. Applying triangle inequality,
 1417 Jensen's inequality, and Cauchy-Schwarz, to the definition of $E_f(\theta^*)$ yields

$$\begin{aligned}
 1418 \quad |\langle v, E_f(\theta^*)v \rangle| &\leq \mathbb{E}_{\theta \sim \rho} \left[2|\langle v, R_f(\theta) \rangle| \cdot |\langle \theta - \theta^*, \nabla^2 f(\theta^*)v \rangle| + |\langle v, R_f(\theta) \rangle|^2 \right] \\
 1419 \quad &\leq \mathbb{E}_{\theta \sim \rho} \left[2\|R_f(\theta)\|_2 \|\theta - \theta^*\|_2 \lambda_{\max}(\nabla^2 f(\theta^*)) \|v\|_2^2 + \|R_f(\theta)\|_2^2 \|v\|_2^2 \right].
 \end{aligned}$$

1420 Recall that under Assumption 2, $\|R_f(\theta)\|_2 \leq H\|\theta - \theta^*\|_2^2$. Substituting this into the above
 1421 bound yields

$$1422 \quad |\langle v, E_f(\theta^*)v \rangle| \leq (2H\lambda_{\max}(\nabla^2 f(\theta^*))M_3 + H^2M_4) \|v\|_2^2$$

1423 where

$$1424 \quad M_p \stackrel{\text{def}}{=} \mathbb{E}_{\theta \sim \rho} [\|\theta - \theta^*\|_2^p]. \quad \blacksquare$$

1425 *Proof of Lemma A.3.* By Lemma A.5, the smoothness constants L_1, \dots, L_d must satisfy

$$1426 \quad \forall \theta \in \Theta, \forall v, \quad \langle v, \text{diag}(L)v \rangle \geq |\langle v, \nabla^2 f(\theta)v \rangle|.$$

1427 where $\text{diag}(L) \in \mathbb{R}^{d \times d}$ is the diagonal matrix with entries $\text{diag}(L)_{i,i} = L_i$. In particular,
 1428 consider $\theta = \theta^*$, and let $\nu \in \{\pm d^{-1/2}\}^d$ be a dense unit vector whose entries all satisfy
 1429 $|\nu(i)| = d^{-1/2}$. Then the smoothness constants must satisfy

$$1430 \quad \langle \nu, \text{diag}(L)\nu \rangle = \frac{1}{d} \sum_{i=1}^d L_i \geq \langle \nu, \nabla^2 f(\theta^*)\nu \rangle.$$

1431 Given θ^* a local minimum, $\nabla^2 f(\theta^*)$ is PSD. Thus by Lemma E.1,

$$1432 \quad \langle \nu, \nabla^2 f(\theta^*)\nu \rangle \geq \frac{1}{\lambda_{\max}(\nabla^2 f(\theta^*))} \langle \nu, \nabla^2 f(\theta^*)\nabla^2 f(\theta^*)^\top \nu \rangle$$

1433 where $\lambda_{\max}(\nabla^2 f(\theta^*))$ denotes the leading eigenvalue of $\nabla^2 f(\theta^*)$. Moreover, Lemma A.6 im-
 1434 plies

$$1435 \quad \langle \nu, \nabla^2 f(\theta^*)\nabla f(\theta^*)^\top \nu \rangle \geq \frac{1}{2c^2} \langle \nu, G(\theta^*)\nu \rangle = \frac{1}{2c^2} \langle \nu, (\mathbb{E}_{\theta \sim \rho}[\nabla f(\theta)\nabla f(\theta)^\top] - E_f(\theta^*)) \nu \rangle.$$

$$1436 \quad \langle v_1, \nabla^2 f(\theta^*)\nabla f(\theta^*)^\top v_1 \rangle \geq \frac{\lambda_{\max}(\text{EGOP}) - \langle v_1, E_f(\theta^*)v_1 \rangle}{2c^2}.$$

1437 Chaining the preceding inequalities yields

$$1438 \quad (\text{E.4}) \quad \langle \nu, \text{diag}(L)\nu \rangle \geq \frac{\langle \nu, \text{EGOP} \nu \rangle - \langle \nu, E_f(\theta^*)\nu \rangle}{2c^2 \lambda_{\max}(\nabla^2 f(\theta^*))}.$$

1439 Given $\nabla^2 f(\theta^*)$ PSD, we have

$$1440 \quad \lambda_{\max}(\nabla^2 f(\theta^*)) = \sqrt{\lambda_{\max}(\nabla^2 f(\theta^*))^2} = \sqrt{\lambda_{\max}(\nabla^2 f(\theta^*)\nabla^2 f(\theta^*)^\top)}.$$

1441 Moreover, by Lemma A.6 we have the following upper bound:

$$\begin{aligned} 1442 \quad \lambda_{\max}(\nabla^2 f(\theta^*)\nabla^2 f(\theta^*)^\top) &\leq \frac{1}{c^2} \lambda_{\max}(G(\theta^*)) \\ 1443 &= \frac{1}{c^2} \max_{v \in \mathbb{R}^d} \langle v, (\mathbb{E}_{\theta \sim \rho}[\nabla f(\theta)\nabla f(\theta)^\top] - E_f(\theta^*))v \rangle \\ 1444 &\leq \frac{1}{c^2} (\lambda_{\max}(\text{EGOP}) + \lambda_{\max}(E_f(\theta^*))). \end{aligned}$$

1445 Combining this inequality with the preceding equation implies

$$1446 \quad \lambda_{\max}(\nabla^2 f(\theta^*)) \leq \frac{1}{c} \sqrt{\lambda_{\max}(\text{EGOP}) + \lambda_{\max}(E_f(\theta^*))}.$$

1447 Substituting this inequality into Eq. E.4 implies

$$1448 \quad \langle \nu, \text{diag}(L)\nu \rangle \geq \frac{\langle \nu, \text{EGOP } \nu \rangle - \langle \nu, E_f(\theta^*)\nu \rangle}{2c\sqrt{\lambda_{\max}(\text{EGOP})} + \lambda_{\max}(E_f(\theta^*))}.$$

1449 Recall that $\nu \in \{\pm d^{-1/2}\}^d$. Thus

$$1450 \quad \langle \nu, \text{diag}(L)\nu \rangle = \sum_{i=1}^d L_i v(i)^2 = \frac{1}{d} \sum_{i=1}^d L_i.$$

1451 Combining this with the above inequality implies

$$1452 \quad \sum_{i=1}^d L_i \geq \frac{d(\langle \nu, \text{EGOP } \nu \rangle - \langle \nu, E_f(\theta^*)\nu \rangle)}{2c\sqrt{\lambda_{\max}(\text{EGOP})} + \lambda_{\max}(E_f(\theta^*))}.$$

1453 By Lemma A.6, $\forall v \in \mathbb{R}^d$

$$1454 \quad |\langle v, E_f(\theta^*)v \rangle| \leq \gamma \|v\|_2^2$$

1455 where

$$1456 \quad \gamma \stackrel{\text{def}}{=} 2H\lambda_{\max}(\nabla^2 f(\theta^*))M_3 + H^2M_4$$

1457 Thus

$$1458 \quad \sum_{i=1}^d L_i \geq \frac{d}{2c} \frac{\langle \nu, \text{EGOP } \nu \rangle - \gamma}{\sqrt{\lambda_{\max}(\text{EGOP})} + \gamma}.$$

1459 As this holds for all $\nu \in \{\pm d^{-1/2}\}^d$, we conclude

$$1460 \quad \sum_{i=1}^d L_i \geq \frac{d}{2c} \cdot \max_{\nu \in \{\pm d^{-1/2}\}^d} \frac{\langle \nu, \text{EGOP } \nu \rangle - \gamma}{\sqrt{\lambda_{\max}(\text{EGOP})} + \gamma}.$$

1461 We observe that for any $v \in \mathbb{R}^d$,

$$1462 \quad \langle v, \text{EGOP } v \rangle = \sum_{i=1}^d \lambda_i \langle v, v_i \rangle^2 \geq \max_k \lambda_k \langle v, v_k \rangle^2$$

1463 where (λ_i, v_i) denote the i^{th} eigenvalue and eigenvector of the EGOP indexed in decreasing
 1464 order of magnitude. In particular, for any v_k , consider $\nu_k \in \mathbb{R}^d$ defined to have entries
 1465 $\nu_k(i) = \text{sign}(v_k(i))d^{-1/2}$. Observe that

$$1466 \quad \nu_k \in \text{argmax}\{\langle \nu, v_k \rangle^2 \mid \nu \in \{\pm d^{-1/2}\}^d\}.$$

1467 For this ν_k ,

$$1468 \quad \langle \nu, v_k \rangle^2 = \left(\sum_{i=1}^d \text{sign}(v_k(i))d^{-1/2}v_k(i) \right)^2 = \frac{1}{d} \left(\sum_{i=1}^d |v_k(i)| \right)^2 = \frac{1}{d} \|v_k\|_1^2.$$

1469 Thus in general,

$$1470 \quad \max_{\nu \in \{\pm d^{-1/2}\}^d} \langle \nu, \text{EGOP } \nu \rangle = \max_{\nu \in \{\pm d^{-1/2}\}^d} \sum_{k=1}^d \lambda_k \langle \nu, v_k \rangle^2 \geq \max_k \frac{\lambda_k}{d} \|v_k\|_1^2.$$

1471 This implies the lower bound

$$1472 \quad \sum_{i=1}^d L_i \geq \frac{d}{2c} \cdot \max_k \frac{\lambda_k \|v_k\|_1^2 / d - \gamma}{\sqrt{\lambda_{\max}(\text{EGOP}) + \gamma}}. \quad \blacksquare$$

1473 *Proof of Lemma A.4.* Consider reparameterized function $\tilde{f}(\theta) \stackrel{\text{def}}{=} f(V\theta)$ and correspond-
 1474 ing $\tilde{\Theta}$ and $\tilde{\rho}$ as defined in Eq. A.3. Observe that Assumption 2 implies $\nabla^2 f(\cdot)$ is also H -
 1475 Lipschitz within $\tilde{\Theta}$ and that $V^T \theta^* \in \tilde{\Theta}$ is a local minimum. Applying Lemma A.6 to $\tilde{f}(\cdot)$
 1476 implies

$$1477 \quad \mathbb{E}_{\tilde{\rho}}[\nabla \tilde{f}(\theta) \nabla \tilde{f}(\theta)^T] = G(V^T \theta^*) + E_{\tilde{f}}[V^T \theta^*]$$

1478 where

$$1479 \quad c^2 \nabla^2 \tilde{f}(V^T \theta^*) \nabla^2 \tilde{f}(V^T \theta^*)^T \preceq G(V^T \theta^*) \preceq 2c^2 \nabla^2 \tilde{f}(V^T \theta^*) \nabla^2 \tilde{f}(V^T \theta^*)^T$$

1480 and

$$1481 \quad |\langle v, E_{\tilde{f}}[V^T \theta^*] v \rangle| \leq \gamma \|v\|_2^2$$

1482 for

$$1483 \quad \gamma \stackrel{\text{def}}{=} 2H \lambda_{\max}(\nabla^2 f(\theta^*)) M_3 + H^2 M_4.$$

1484 We've used the fact that

$$1485 \quad M_p \stackrel{\text{def}}{=} \mathbb{E}_{\theta \sim \rho} [\|\theta - \theta^*\|_2^p] = \mathbb{E}_{\theta \sim \rho} [\|V^T(\theta - \theta^*)\|_2^p] = \mathbb{E}_{\tilde{\theta} \sim \tilde{\rho}} [\|\tilde{\theta} - V^T \theta^*\|_2^p]$$

1486 and

$$1487 \quad \lambda_{\max}(\nabla^2 \tilde{f}(V^T \theta^*))$$

1488 to bound $|\langle v, E_{\tilde{f}}[V^T \theta^*] v \rangle|$ by the same value γ .

1489 Combining the above matrix inequalities with the implication $E_{\tilde{f}}[V^T \theta^*] \preceq \gamma \mathbb{I}$ implies

$$1490 \quad c^2 \nabla^2 \tilde{f}(V^T \theta^*) \nabla^2 \tilde{f}(V^T \theta^*)^T \preceq \mathbb{E}_{\theta \sim \rho} [\nabla \tilde{f}(\theta) \nabla \tilde{f}(\theta)^T] + \gamma \mathbb{I}$$

1491 Recall that the PSD ordering $X \preceq Y$ implies $X^{1/2} \preceq Y^{1/2}$ (for completeness, see statement
 1492 and proof in Lemma E.2). Thus the preceding matrix inequality implies

$$1493 \quad \nabla^2 \tilde{f}(V^T \theta^*) \preceq \frac{1}{c} \left(\mathbb{E}_{\theta \sim \rho} [\nabla \tilde{f}(\theta) \nabla \tilde{f}(\theta)^T] + \gamma \mathbb{I} \right)^{1/2}$$

1494 Lemma A.7 implies $\mathbb{E}_{\tilde{\rho}}[\nabla \tilde{f}(\theta) \nabla \tilde{f}(\theta)^T]$ is diagonal, and its entries are given by the eigen-
 1495 values of $\text{EGOP}(f) \stackrel{\text{def}}{=} \mathbb{E}_{\theta \sim \rho} [\nabla f(\theta) \nabla f(\theta)^T]$. As the EGOP is an expectation of PSD matrices,
 1496 these eigenvalues are all non-negative, and thus

$$1497 \quad \mathbb{E}_{\theta \sim \tilde{\rho}} [\nabla \tilde{f}(\theta) \nabla \tilde{f}(\theta)^T]_{i,i}^{1/2} = \sqrt{\lambda_i(\text{EGOP}(f))}.$$

1498 Since $\mathbb{E}_{\theta \sim \rho}[\nabla \tilde{f}(\theta) \nabla \tilde{f}(\theta)^\top]$ is diagonal and $\gamma \mathbb{I}$ is diagonal,

$$1499 \quad \left(\mathbb{E}_{\theta \sim \rho}[\nabla \tilde{f}(\theta) \nabla \tilde{f}(\theta)^\top] + \gamma \mathbb{I} \right)_{i,i}^{1/2} = \sqrt{\mathbb{E}_{\theta \sim \rho}[\nabla \tilde{f}(\theta) \nabla \tilde{f}(\theta)^\top]_{i,i} + \gamma}.$$

1500 Thus for all $v \in \mathbb{R}^d$, we have

$$1501 \quad (\text{E.5}) \quad \langle v, \nabla^2 \tilde{f}(V^T \theta^*) v \rangle \leq \frac{1}{c} \sum_{i=1}^d v(i)^2 \sqrt{\mathbb{E}_{\theta \sim \rho}[\nabla \tilde{f}(\theta) \nabla \tilde{f}(\theta)^\top]_{i,i} + \gamma}.$$

1502 Employing reverse triangle inequality and the fact that $V^T \theta^*$ is a local minimum implies
 1503 $\nabla^2 \tilde{f}(V^T \theta^*)$ is PSD, we can show that Assumption 2 implies

$$1504 \quad |\langle v, \nabla^2 \tilde{f}(\theta) v \rangle| \leq \langle v, \nabla^2 \tilde{f}(V^T \theta^*) v \rangle + H \|\theta - V^T \theta^*\|_2 \|v\|_2^2$$

1505 that $\forall \theta \in \tilde{\Theta}, \forall v \in \mathbb{R}^d$. Observe that $\|\theta - V^T \theta^*\|_2 \leq \tilde{D}_2 = D_2$ for D_2 the ℓ_2 -diameter of Θ .
 1506 Thus combining the above with Eq. E.5 implies $\forall \theta \in \tilde{\Theta}, \forall v \in \mathbb{R}^d$,

$$1507 \quad |\langle v, \nabla^2 \tilde{f}(\theta) v \rangle| \leq \sum_{i=1}^d \sqrt{\lambda_i(\text{EGOP}(f)) + \gamma/c + HD_2} v(i)^2.$$

1508 Thus by Lemma A.5, this implies that $\tilde{f}(\cdot)$ satisfies coordinate-wise smoothness with respect
 1509 to constants

$$1510 \quad L_i \leq \frac{1}{c} \sqrt{\lambda_i(\text{EGOP}(f)) + \gamma} + HD_2.$$

1511 This implies the sum of the coordinate-wise smoothness constants is bounded above as

$$1512 \quad L_{\tilde{f}} \leq \sum_{i=1}^d \left(\frac{1}{c} \sqrt{\lambda_i(\text{EGOP}(f)) + \gamma} + HD_2 \right)$$

$$1513 \quad \leq d \left(\frac{\sqrt{\gamma}}{c} + HD_2 \right) + \frac{1}{c} \sum_{i=1}^d \sqrt{\lambda_i(\text{EGOP}(f))}. \quad \blacksquare$$

1514 **E.1.1. Lipschitz Constants of Hessians in Machine Learning.** In this section, we note
 1515 one family of naturally-motivated non-convex objectives with locally Lipschitz Hessians that
 1516 arise in machine learning problems. We consider the over-parameterized matrix factorization
 1517 problem: let $\theta \in \mathbb{R}^{(d_1+d_2)k}$ be parameters, whose entries can be grouped into two matrices as
 1518 $\theta = (L, R)$ for $L \in \mathbb{R}^{d_1 \times k}, R \in \mathbb{R}^{d_2 \times k}$. We define the objective

$$1519 \quad (\text{E.6}) \quad f(\theta) = \|\mathcal{A}(LR^\top) - b\|_2^2$$

1520 where $b \in \mathbb{R}^m$ correspond to measurements of some matrix under map $\mathcal{A}(\cdot)$, and where
 1521 $\mathcal{A} : \mathbb{R}^{d_1 \times d_2} \rightarrow \mathbb{R}^m$ denotes a map

$$1522 \quad \mathcal{A}(X) \stackrel{\text{def}}{=} (\langle A_1, X \rangle, \dots, \langle A_m, X \rangle) \in \mathbb{R}^m.$$

1523 We note that one a special case of Eq. E.6 is the set of objectives that arise when training a
 1524 two-layer linear feed-forward network using mean-squared-error loss.

1525 We can bound the Lipschitz constant of the Hessian of this objectives in this family:

1526 **Lemma E.3.** Consider the overparameterized matrix factorization objective with a linear
 1527 measurement map, as defined in Eq. E.6. Let $\theta \in \mathbb{R}^{(d_1+d_2)k}$ be parameters, whose entries
 1528 can be grouped into two matrices: $\theta = (L, R)$ for $L \in \mathbb{R}^{d_1 \times k}, R \in \mathbb{R}^{d_2 \times k}$. Consider any
 1529 measurement vector $b \in \mathbb{R}^m$. Then in the ball $\|\theta\|_2 \leq B$, the objective in Eq. E.6 satisfies

$$1530 \quad \|\nabla^2 f(\theta_1) - \nabla^2 f(\theta_2)\|_{\text{op}} \leq 12B \left(\sum_{i=1}^m \|A_i\|_{\text{F}}^2 \right) \|\theta_1 - \theta_2\|_2.$$

1531 To prove this result, we first characterize the quadratic form of the Hessian of Eq. E.6.

1532 **Lemma E.4.** For the objective $f(\cdot)$ as defined in Eq. E.6, for any $\theta, v \in \mathbb{R}^{(d_1+d_2)k}$ with
 1533 entries denoted $\theta = (L, R)$ and $v = (U, V)$, the Hessian quadratic form can be expressed as

$$1534 \quad D^2 f(L, R)[U, V] = 2\|\mathcal{A}(UR^\top + LV^\top)\|_2^2 + 4\langle \mathcal{A}(LR^\top) - b, \mathcal{A}(UV^\top) \rangle.$$

1535 *Proof.* We begin by deriving the gradient form using the limit definition:

$$1536 \quad \nabla f(L, R)[U, V] \stackrel{\text{def}}{=} \lim_{t \rightarrow 0} \frac{1}{t} (f(L + tU, R + tV) - f(L, R))$$

$$1537 \quad = \lim_{t \rightarrow 0} \frac{1}{t} \left(\|\mathcal{A}((L + tU)(R + tV)^\top) - b\|_2^2 - \|\mathcal{A}(LR^\top) - b\|_2^2 \right).$$

1538 By linearity of the map $\mathcal{A}(\cdot)$,

$$1539 \quad \mathcal{A}((L + tU)(R + tV)^\top) = \mathcal{A}(LR^\top + tUR^\top + tLV^\top + t^2UV^\top)$$

$$1540 \quad = \mathcal{A}(LR^\top) + t\mathcal{A}(UR^\top + LV^\top) + t^2\mathcal{A}(UV^\top).$$

1541 Substituting this into the above limit yields

$$1542 \quad \nabla f(L, R)[U, V] = \lim_{t \rightarrow 0} \frac{1}{t} \left(\|\Delta + t\mathcal{A}(UR^\top + LV^\top) + t^2\mathcal{A}(UV^\top)\|_2^2 - \|\Delta\|_2^2 \right)$$

1543 where $\Delta \stackrel{\text{def}}{=} \mathcal{A}(LR^\top) - b$. Expanding the squared norm, this yields

$$1544 \quad \nabla f(L, R)[U, V] = \lim_{t \rightarrow 0} \frac{1}{t} \left(\|\Delta\|_2^2 + \|t\mathcal{A}(UR^\top + LV^\top) + t^2\mathcal{A}(UV^\top)\|_2^2 \right.$$

$$1545 \quad \left. + 2\langle \Delta, t\mathcal{A}(UR^\top + LV^\top) + t^2\mathcal{A}(UV^\top) \rangle - \|\Delta\|_2^2 \right)$$

$$1546 \quad = \lim_{t \rightarrow 0} \frac{1}{t} \left(\|t\mathcal{A}(UR^\top + LV^\top) + t^2\mathcal{A}(UV^\top)\|_2^2 \right.$$

$$1547 \quad \left. + 2\langle \Delta, t\mathcal{A}(UR^\top + LV^\top) + t^2\mathcal{A}(UV^\top) \rangle \right)$$

$$1548 \quad = \lim_{t \rightarrow 0} \left(t\|\mathcal{A}(UR^\top + LV^\top) + t\mathcal{A}(UV^\top)\|_2^2 \right.$$

$$1549 \quad \left. + 2\langle \Delta, \mathcal{A}(UR^\top + LV^\top) + t\mathcal{A}(UV^\top) \rangle \right)$$

$$1550 \quad = 0 + 2\langle \Delta, \mathcal{A}(UR^\top + LV^\top) + 0 \rangle$$

$$1551 \quad = 2\langle \mathcal{A}(LR^\top) - b, \mathcal{A}(UR^\top + LV^\top) \rangle$$

1552 where the last line follows from $\Delta \stackrel{\text{def}}{=} \mathcal{A}(LR^\top) - b$. Given this expression for the gradient, we
 1553 can then define the Hessian quadratic form on input $v = (U, V)$ as

$$1554 \quad \nabla^2 f(L, R)[U, V] \stackrel{\text{def}}{=} \lim_{t \rightarrow 0} \frac{1}{t} (\nabla f(L + tU, R + tV)[U, V] - \nabla f(L, R)[U, V]).$$

1555 Given the above expression for the gradient, this yields

$$1556 \quad \nabla^2 f(L, R)[U, V] = \lim_{t \rightarrow 0} \frac{1}{t} \left(2\langle \mathcal{A}((L + tU)(R + tV)^\top) - b, \mathcal{A}(U(R + tV)^\top + (L + tU)V^\top) \rangle \right. \\ 1557 \quad \left. - 2\langle \Delta, \mathcal{A}(UR^\top + LV^\top) \rangle \right).$$

1558 As noted above,

$$1559 \quad \mathcal{A}((L + tU)(R + tV)^\top) = \mathcal{A}(LR^\top) + t\mathcal{A}(UR^\top + LV^\top) + t^2\mathcal{A}(UV^\top).$$

1560 Substituting this in to the first term in the expression for the Hessian, and using linearity of
 1561 $\mathcal{A}(\cdot)$ to simplify terms, yields

$$1562 \quad \nabla^2 f(L, R)[U, V] = \lim_{t \rightarrow 0} \frac{2}{t} \left(\langle \Delta + t\mathcal{A}(UR^\top + LV^\top) + t^2\mathcal{A}(UV^\top), \mathcal{A}(UR^\top + LV^\top) + 2t\mathcal{A}(UV^\top) \rangle \right. \\ 1563 \quad \left. - \langle \Delta, \mathcal{A}(UR^\top + LV^\top) \rangle \right) \\ 1564 \quad = \lim_{t \rightarrow 0} \frac{2}{t} \left(\langle \Delta, \mathcal{A}(UR^\top + LV^\top) \rangle + \langle t\mathcal{A}(UR^\top + LV^\top) + t^2\mathcal{A}(UV^\top), \mathcal{A}(UR^\top + LV^\top) \rangle \right. \\ 1565 \quad \left. + \langle \Delta, 2t\mathcal{A}(UV^\top) \rangle + \langle t\mathcal{A}(UR^\top + LV^\top) + t^2\mathcal{A}(UV^\top), 2t\mathcal{A}(UV^\top) \rangle \right. \\ 1566 \quad \left. - \langle \Delta, \mathcal{A}(UR^\top + LV^\top) \rangle \right) \\ 1567 \quad = \lim_{t \rightarrow 0} \frac{2}{t} \left(\langle t\mathcal{A}(UR^\top + LV^\top) + t^2\mathcal{A}(UV^\top), \mathcal{A}(UR^\top + LV^\top) \rangle \right. \\ 1568 \quad \left. + \langle \Delta, 2t\mathcal{A}(UV^\top) \rangle + \langle t\mathcal{A}(UR^\top + LV^\top) + t^2\mathcal{A}(UV^\top), 2t\mathcal{A}(UV^\top) \rangle \right) \\ 1569 \quad = \lim_{t \rightarrow 0} \frac{2}{t} \left(t\langle \mathcal{A}(UR^\top + LV^\top) + t\mathcal{A}(UV^\top), \mathcal{A}(UR^\top + LV^\top) \rangle \right. \\ 1570 \quad \left. + 2t\langle \Delta, \mathcal{A}(UV^\top) \rangle + 2t^2\langle \mathcal{A}(UR^\top + LV^\top) + t\mathcal{A}(UV^\top), \mathcal{A}(UV^\top) \rangle \right).$$

1571 Taking the limit as $t \rightarrow 0$,

$$1572 \quad \nabla^2 f(L, R)[U, V] = 2 \left(\langle \mathcal{A}(UR^\top + LV^\top), \mathcal{A}(UR^\top + LV^\top) \rangle + 2\langle \Delta, \mathcal{A}(UV^\top) \rangle \right) \\ 1573 \quad = 2 \left(\|\mathcal{A}(UR^\top + LV^\top)\|_2^2 + 2\langle \mathcal{A}(LR^\top) - b, \mathcal{A}(UV^\top) \rangle \right) \\ 1574 \quad = 2\|\mathcal{A}(UR^\top + LV^\top)\|_2^2 + 4\langle \mathcal{A}(LR^\top) - b, \mathcal{A}(UV^\top) \rangle$$

1575 which yields the result. ■

1576 We can prove Lemma E.3, which bounds the Lipschitz constant of the Hessian of $f(\cdot)$ as
 1577 defined in Eq. E.6.

1578 *Proof of Lemma E.3.* Given some vector $v \in \mathbb{R}^{(d_1+d_2)k}$, group the entries of the vector v
 1579 into two matrices: let $v = (U, V)$ for matrices $U \in \mathbb{R}^{d_1 \times k}$, $V \in \mathbb{R}^{d_2 \times k}$. Then the quadratic
 1580 form of the Hessian is

$$1581 \quad \langle v, \nabla^2 f(\theta)v \rangle = D^2 f(L, R)[U, V].$$

1582 The objective $f(\cdot)$ satisfies Assumption 2 with respect to Lipschitz constant H if $\forall \theta_1, \theta_2, v \in$
 1583 $\mathbb{R}^{(d_1+d_2)k}$,

$$1584 \quad |\langle v, (\nabla^2 f(\theta_1) - \nabla^2 f(\theta_2))v \rangle| \leq H \|v\|_2^2 \|\theta_1 - \theta_2\|_2.$$

1585 For any pair $\theta_1, \theta_2 \in \mathbb{R}^{(d_1+d_2)k}$, denote the entries by $\theta_1 = (L_1, R_1)$ and $\theta_2 = (L_2, R_2)$, and for
 1586 any $v \in \mathbb{R}^{(d_1+d_2)k}$, denote the entries by $v = (U, V)$ as above. Then the above inequality is
 1587 equivalent to

$$1588 \quad |D^2 f(L_1, R_1)[U, V] - D^2 f(L_2, R_2)[U, V]| \leq H (\|U\|_F^2 + \|V\|_F^2) \|(L_1, R_1) - (L_2, R_2)\|_2.$$

1589 where

$$1590 \quad \|(L, R)\|_2 \stackrel{\text{def}}{=} \sqrt{\|L\|_F^2 + \|R\|_F^2} = \|\theta\|_2.$$

1591 By Lemma E.4, for any $\theta = (L, R)$ and any $v = (U, V)$, the Hessian satisfies

$$1592 \quad D^2 f(L, R)[U, V] = 2\|\mathcal{A}(UR^\top + LV^\top)\|_2^2 + 4\langle \mathcal{A}(LR^\top) - b, \mathcal{A}(UV^\top) \rangle.$$

1593 By linearity of $\mathcal{A}(\cdot)$, we can expand the squared norm:

$$1594 \quad D^2 f(L, R)[U, V] = 2\|\mathcal{A}(UR^\top) + \mathcal{A}(LV^\top)\|_2^2 + 4\langle \mathcal{A}(LR^\top) - b, \mathcal{A}(UV^\top) \rangle$$

$$1595 \quad = 2\left(\|\mathcal{A}(LV^\top)\|_2^2 + \|\mathcal{A}(UR^\top)\|_2^2\right) + 4\langle \mathcal{A}(UR^\top), \mathcal{A}(LV^\top) \rangle + 4\langle \mathcal{A}(LR^\top) - b, \mathcal{A}(UV^\top) \rangle.$$

1596 Thus

(E.7)

$$1597 \quad D^2 f(L_1, R_1)[U, V] - D^2 f(L_2, R_2)[U, V] = 2\left(\|\mathcal{A}(L_1V^\top)\|_2^2 - \|\mathcal{A}(L_2V^\top)\|_2^2\right)$$

$$1598 \quad (E.8) \quad + 2\left(\|\mathcal{A}(UR_1^\top)\|_2^2 - \|\mathcal{A}(UR_2^\top)\|_2^2\right)$$

$$1599 \quad (E.9) \quad + 4\left(\langle \mathcal{A}(UR_1^\top), \mathcal{A}(L_1V^\top) \rangle - \langle \mathcal{A}(UR_2^\top), \mathcal{A}(L_2V^\top) \rangle\right)$$

$$1600 \quad (E.10) \quad + 4\langle \mathcal{A}(L_1R_1^\top) - \mathcal{A}(L_2R_2^\top), \mathcal{A}(UV^\top) \rangle.$$

1601 We bound the magnitude of each term in sequence. For the first term in Line E.7, we observe

$$1602 \quad \left| \|\mathcal{A}(L_1V^\top)\|_2^2 - \|\mathcal{A}(L_2V^\top)\|_2^2 \right| = |\langle \mathcal{A}(L_1V^\top) + \mathcal{A}(L_2V^\top), \mathcal{A}(L_1V^\top) - \mathcal{A}(L_2V^\top) \rangle|$$

$$1603 \quad = |\langle \mathcal{A}\left((L_1 + L_2)V^\top\right), \mathcal{A}\left((L_1 - L_2)V^\top\right) \rangle|$$

$$1604 \quad \leq \|\mathcal{A}\left((L_1 + L_2)V^\top\right)\|_2 \|\mathcal{A}\left((L_1 - L_2)V^\top\right)\|_2.$$

1605 The operator norm of $\mathcal{A}(\cdot)$ can be bounded as

$$1606 \quad \|\mathcal{A}(X)\|_2^2 = \sum_{i=1}^m \langle A_i, X \rangle^2 \leq \|X\|_F^2 \sum_{i=1}^m \|A_i\|_F^2.$$

1607 Thus

$$1608 \quad (\text{E.11}) \quad \|\mathcal{A}\|_2 \leq \left(\sum_{i=1}^m \|A_i\|_{\mathbb{F}}^2 \right)^{1/2}.$$

1609 We can thus bound the term in Line E.7 as

$$\begin{aligned} 1610 \quad \left| \|\mathcal{A}(L_1 V^\top)\|_2^2 - \|\mathcal{A}(L_2 V^\top)\|_2^2 \right| &\leq \|\mathcal{A}((L_1 + L_2)V^\top)\|_2 \|\mathcal{A}((L_1 - L_2)V^\top)\|_2 \\ 1611 &\leq \|\mathcal{A}\|_2^2 \|(L_1 + L_2)V^\top\|_{\mathbb{F}} \|(L_1 - L_2)V^\top\|_{\mathbb{F}} \\ 1612 &\leq \|\mathcal{A}\|_2^2 \|L_1 + L_2\|_{\mathbb{F}} \|V\|_{\mathbb{F}} \|L_1 - L_2\|_{\mathbb{F}} \|V\|_{\mathbb{F}}. \end{aligned}$$

1613 In the ball $\|\theta\|_2 \leq B$, we have that for $\theta = (L, R)$

$$1614 \quad \|L\|_{\mathbb{F}} \leq \|\theta\|_2 \leq B.$$

1615 Thus $\|L_1 + L_2\|_{\mathbb{F}} \leq 2B$, so we can bound

$$\begin{aligned} 1616 \quad \left| \|\mathcal{A}(L_1 V^\top)\|_2^2 - \|\mathcal{A}(L_2 V^\top)\|_2^2 \right| &\leq 2B \|\mathcal{A}\|_2^2 \|V\|_{\mathbb{F}}^2 \|L_1 - L_2\|_{\mathbb{F}} \\ 1617 &\leq 2B \|\mathcal{A}\|_2^2 \|V\|_{\mathbb{F}}^2 \|(L_1, R_1) - (L_2, R_2)\|_2 \end{aligned}$$

1618 where the last line follows by

$$1619 \quad \|(L_1, R_1) - (L_2, R_2)\|_2 \stackrel{\text{def}}{=} \sqrt{\|L_1 - L_2\|_{\mathbb{F}}^2 + \|R_1 - R_2\|_{\mathbb{F}}^2} \geq \|L_1 - L_2\|_{\mathbb{F}}.$$

1620 An analogous argument bounds the term in Line E.8 as

$$1621 \quad \left| \|\mathcal{A}(UR_1^\top)\|_2^2 - \|\mathcal{A}(UR_2^\top)\|_2^2 \right| \leq 2B \|\mathcal{A}\|_2^2 \|U\|_{\mathbb{F}}^2 \|(L_1, R_1) - (L_2, R_2)\|_2.$$

1622 To bound the term in Line E.9, we first observe that

$$\begin{aligned} 1623 \quad &\langle \mathcal{A}(UR_1^\top), \mathcal{A}(L_1 V^\top) \rangle - \langle \mathcal{A}(UR_2^\top), \mathcal{A}(L_2 V^\top) \rangle \\ 1624 &= \langle \mathcal{A}(UR_1^\top), \mathcal{A}(L_1 V^\top) \rangle - \langle \mathcal{A}(L_1 V^\top), \mathcal{A}(UR_2^\top) \rangle \\ 1625 &\quad + \langle \mathcal{A}(L_1 V^\top), \mathcal{A}(UR_2^\top) \rangle - \langle \mathcal{A}(UR_2^\top), \mathcal{A}(L_2 V^\top) \rangle \\ 1626 &= \langle \mathcal{A}(UR_1^\top) - \mathcal{A}(UR_2^\top), \mathcal{A}(L_1 V^\top) \rangle + \langle \mathcal{A}(L_1 V^\top) - \mathcal{A}(L_2 V^\top), \mathcal{A}(UR_2^\top) \rangle \\ 1627 &= \langle \mathcal{A}(U(R_1 - R_2)^\top), \mathcal{A}(L_1 V^\top) \rangle + \langle \mathcal{A}((L_1 - L_2)V^\top), \mathcal{A}(UR_2^\top) \rangle \end{aligned}$$

1628 Thus

$$\begin{aligned} 1629 \quad &|\langle \mathcal{A}(UR_1^\top), \mathcal{A}(L_1 V^\top) \rangle - \langle \mathcal{A}(UR_2^\top), \mathcal{A}(L_2 V^\top) \rangle| \\ 1630 &= |\langle \mathcal{A}(U(R_1 - R_2)^\top), \mathcal{A}(L_1 V^\top) \rangle + \langle \mathcal{A}((L_1 - L_2)V^\top), \mathcal{A}(UR_2^\top) \rangle| \\ 1631 &\leq \|\mathcal{A}(U(R_1 - R_2)^\top)\|_2 \|\mathcal{A}(L_1 V^\top)\|_2 + \|\mathcal{A}((L_1 - L_2)V^\top)\|_2 \|\mathcal{A}(UR_2^\top)\|_2 \\ 1632 &\leq \|\mathcal{A}\|_2^2 \|U(R_1 - R_2)^\top\|_{\mathbb{F}} \|L_1 V^\top\|_{\mathbb{F}} + \|\mathcal{A}\|_2^2 \|(L_1 - L_2)V^\top\|_{\mathbb{F}} \|UR_2^\top\|_{\mathbb{F}} \\ 1633 &\leq \|\mathcal{A}\|_2^2 \|U\|_{\mathbb{F}} \|R_1 - R_2\|_{\mathbb{F}} \|L_1\|_{\mathbb{F}} \|V\|_{\mathbb{F}} + \|\mathcal{A}\|_2^2 \|L_1 - L_2\|_{\mathbb{F}} \|V\|_{\mathbb{F}} \|U\|_{\mathbb{F}} \|R_2\|_{\mathbb{F}} \\ 1634 &= \|\mathcal{A}\|_2^2 \|U\|_{\mathbb{F}} \|V\|_{\mathbb{F}} (\|R_1 - R_2\|_{\mathbb{F}} \|L_1\|_{\mathbb{F}} + \|L_1 - L_2\|_{\mathbb{F}} \|R_2\|_{\mathbb{F}}) \end{aligned}$$

1635 In the ball $\|\theta\|_2 \leq B$, we have that for $\|L_1\|_F \leq B$ and $\|R_2\|_F \leq B$. Thus

$$1636 \quad |\langle \mathcal{A}(UR_1^\top), \mathcal{A}(L_1V^\top) \rangle - \langle \mathcal{A}(UR_2^\top), \mathcal{A}(L_2V^\top) \rangle|$$

$$1637 \quad \leq B\|\mathcal{A}\|_2^2\|U\|_F\|V\|_F(\|R_1 - R_2\|_F + \|L_1 - L_2\|_F).$$

1638 Recall that for any $a, b \in R$, it holds that $a + b \leq 2\sqrt{a^2 + b^2}$ and $2ab \leq a^2 + b^2$. Thus

$$1639 \quad \|R_1 - R_2\|_F + \|L_1 - L_2\|_F \leq 2\sqrt{\|L_1 - L_2\|_F^2 + \|R_1 - R_2\|_F^2} = 2\|(L_1, R_1) - (L_2, R_2)\|_2$$

1640 and

$$1641 \quad \|U\|_F\|V\|_F \leq \frac{1}{2}(\|U\|_F^2 + \|V\|_F^2).$$

1642 Combining these bounds yields

$$1643 \quad |\langle \mathcal{A}(UR_1^\top), \mathcal{A}(L_1V^\top) \rangle - \langle \mathcal{A}(UR_2^\top), \mathcal{A}(L_2V^\top) \rangle|$$

$$1644 \quad \leq B\|\mathcal{A}\|_2^2(\|U\|_F^2 + \|V\|_F^2)\|(L_1, R_1) - (L_2, R_2)\|_2.$$

1645 Lastly, we bound the term in Line E.10. We begin by observing that

$$1646 \quad \mathcal{A}(L_1R_1^\top) - \mathcal{A}(L_2R_2^\top) = \mathcal{A}(L_1R_1^\top) - \mathcal{A}(L_1R_2^\top) + \mathcal{A}(L_1R_2^\top) - \mathcal{A}(L_2R_2^\top)$$

$$1647 \quad = \mathcal{A}\left(L_1(R_1 - R_2)^\top\right) + \mathcal{A}\left((L_1 - L_2)R_2^\top\right)$$

1648 Thus

$$1649 \quad |\langle \mathcal{A}(L_1R_1^\top) - \mathcal{A}(L_2R_2^\top), \mathcal{A}(UV^\top) \rangle| = |\langle \mathcal{A}\left(L_1(R_1 - R_2)^\top\right) + \mathcal{A}\left((L_1 - L_2)R_2^\top\right), \mathcal{A}(UV^\top) \rangle|$$

$$1650 \quad = |\langle \mathcal{A}\left(L_1(R_1 - R_2)^\top\right), \mathcal{A}(UV^\top) \rangle + \langle \mathcal{A}\left((L_1 - L_2)R_2^\top\right), \mathcal{A}(UV^\top) \rangle|$$

$$1651 \quad \leq \|\mathcal{A}\left(L_1(R_1 - R_2)^\top\right)\|_2\|\mathcal{A}(UV^\top)\|_2 + \|\mathcal{A}\left((L_1 - L_2)R_2^\top\right)\|_2\|\mathcal{A}(UV^\top)\|_2$$

$$1652 \quad \leq \|\mathcal{A}\|_2^2\|L_1(R_1 - R_2)^\top\|_F\|UV^\top\|_F + \|\mathcal{A}\|_2^2\|(L_1 - L_2)R_2^\top\|_F\|UV^\top\|_F$$

$$1653 \quad \leq \|\mathcal{A}\|_2^2(\|L_1\|_F\|R_1 - R_2\|_F\|U\|_F\|V\|_F + \|L_1 - L_2\|_F\|R_2\|_F\|U\|_F\|V\|_F)$$

$$1654 \quad = \|\mathcal{A}\|_2^2\|U\|_F\|V\|_F(\|L_1\|_F\|R_1 - R_2\|_F + \|L_1 - L_2\|_F\|R_2\|_F).$$

1655 In the ball $\|\theta\|_2 \leq B$, we have that for $\|L_1\|_F \leq B$ and $\|R_2\|_F \leq B$. Thus

$$1656 \quad |\langle \mathcal{A}(L_1R_1^\top) - \mathcal{A}(L_2R_2^\top), \mathcal{A}(UV^\top) \rangle| \leq B\|\mathcal{A}\|_2^2\|U\|_F\|V\|_F(\|R_1 - R_2\|_F + \|L_1 - L_2\|_F).$$

1657 Recalling the bounds $\|V\|_F\|U\|_F$ and $(\|L_1 - L_2\|_F + \|R_1 - R_2\|_F)$ established above, we have

$$1658 \quad |\langle \mathcal{A}(L_1R_1^\top) - \mathcal{A}(L_2R_2^\top), \mathcal{A}(UV^\top) \rangle| \leq B\|\mathcal{A}\|_2^2(\|U\|_F^2 + \|V\|_F^2)\|(L_1, R_1) - (L_2, R_2)\|_2.$$

1659 Employing triangle inequality and each of the above bounds implies

$$1660 \quad |D^2f(L_1, R_1)[U, V] - D^2f(L_2, R_2)[U, V]| \leq 2 \cdot 2B\|\mathcal{A}\|_2^2\|V\|_F^2\|(L_1, R_1) - (L_2, R_2)\|_2$$

$$1661 \quad \quad \quad + 2 \cdot 2B\|\mathcal{A}\|_2^2\|U\|_F^2\|(L_1, R_1) - (L_2, R_2)\|_2$$

$$1662 \quad \quad \quad + 4B\|\mathcal{A}\|_2^2(\|U\|_F^2 + \|V\|_F^2)\|(L_1, R_1) - (L_2, R_2)\|_2$$

$$1663 \quad \quad \quad + 4B\|\mathcal{A}\|_2^2(\|U\|_F^2 + \|V\|_F^2)\|(L_1, R_1) - (L_2, R_2)\|_2$$

$$1664 \quad \leq 12B\|\mathcal{A}\|_2^2(\|U\|_F^2 + \|V\|_F^2)\|(L_1, R_1) - (L_2, R_2)\|_2$$

1665 which implies the result. ■

1666 E.2. Deferred Proofs from Section A.2.

1667 *Proof of Lemma A.11.* Observe that for any $\theta \in \mathbb{R}^d$,

$$1668 \quad \nabla f(\theta) = A^\top(A\theta - y) = A^\top(A\theta - A\theta^* - A\eta) = A^\top A(\theta - \theta^*) + A\eta.$$

1669 Thus for ρ a standard Gaussian,

$$\begin{aligned} 1670 \quad \text{EGOP}(f) &= \mathbb{E}_{\theta \sim \rho} \left[\left(A^\top(A\theta - y) \right) \left(A^\top(A\theta - y) \right)^\top \right] \\ 1671 \quad &= A^\top \mathbb{E}_{\theta \sim \rho} \left[A\theta\theta^\top A + yy^\top \right] A \\ 1672 \quad &= \left(A^\top A \right)^2 + (A^\top y)(A^\top y)^\top \end{aligned}$$

1673 using the fact that ρ is mean-zero and isotropic. Thus $\text{EGOP}(f)$ is a rank-1 perturbation
1674 of the matrix $(A^\top A)^2$, which is itself PSD and has eigenvalues $\sigma_k(A)^4$. Moreover if we let
1675 $A = Q_1 \Sigma Q_2^\top$ denote the SVD of A , we can rewrite

$$1676 \quad \text{EGOP}(f) = Q_2 \left(\Sigma^4 + \Sigma Q_1^\top yy^\top Q_1 \Sigma \right) V^\top$$

1677 where Σ^4 is diagonal and has entries $\{\sigma_i^4(A)\}_{i=1}^n$. By Golub (1973) [16], the eigenvalues of
1678 $\text{EGOP}(f)$ thus satisfy

$$1679 \quad \sigma_i^4(A) \geq \lambda_i(\text{EGOP}(f)) \geq \sigma_{i+1}^4(A) \quad \forall i \in [2, \dots, n]$$

1680 and

$$1681 \quad \sigma_1^4(A) + \|\Sigma Q_1^\top y\|_2^2 \geq \lambda_1(\text{EGOP}) \geq \sigma_2^4(A)$$

1682 which implies the result. ■

1683 Appendix F. Extended Discussion of Related Works.

1684 In this section, we expand on the overview of related work presented in Section 1.1.

1685 *Geometric sensitivity of adaptive methods.* A large body of research has been devoted to
1686 understanding the settings in which the per-coordinate learning rates of adaptive algorithms
1687 confer a benefit over more basic methods, such as (S)GD. Recently there has been renewed
1688 interest in distinguishing the properties of adaptive algorithms versus SGD because several
1689 empirical studies suggest that adaptive methods outperform SGD when training transformer
1690 models [23, 42]. Traditional analyses of Adagrad in the context of online convex optimization
1691 establish regret bounds which can be either better or worse than those enjoyed by SGD by
1692 up to a factor of \sqrt{d} , for d the number of problem parameters [4, 14]. More recent research
1693 has studied the rates at which adaptive methods converge to stationary points. Several works
1694 study convergence guarantees, measured in terms of the ℓ_2 norm of the gradient, on smooth
1695 non-convex objectives [11, 38]. These results show that Adagrad, Adam and related variants
1696 achieve rates matching those enjoyed by SGD, while having the key distinction that Adagrad
1697 and its variants do not require a-priori knowledge of the objective function's Lipschitz constant
1698 [11, 38].

1699 In order to shed light on the question of when adaptive algorithms enjoy provably stronger
 1700 guarantees than SGD, a recent line of work studies convergence under more refined geometric
 1701 assumptions, with particular emphasis on assumptions that are *not* rotationally invariant
 1702 [21, 25, 40]. Xie et al. [40] establish convergence guarantees in terms of the L_∞ smoothness
 1703 constant of the objective and also in terms of the related Hessian 1-norm. They show that
 1704 rotationally invariant geometric assumptions do not suffice to capture settings when Adam
 1705 out-performs SGD through experiments examining the sensitivity of Adam to orthonormal
 1706 rotations [40].

1707 Jiang et al. [21] and Liu et al. [25] study convergence of Adagrad on objectives that satisfy
 1708 coordinate-wise smoothness. Both works prove similar convergence guarantees for Adagrad
 1709 in terms of ℓ_1 gradient norm, rather than the ℓ_2 gradient norm measure more widely studied;
 1710 Jiang et al. [21] establish convergence guarantees showing that SGD remains worst-case
 1711 optimal even in the setting of coordinate-wise smoothness when measuring stationarity with
 1712 the ℓ_2 gradient norm, motivating the need to measure stationarity with the ℓ_1 norm in order to
 1713 prove separation between Adagrad and SGD. Using this ℓ_1 stationarity measure, Jiang et al.
 1714 [21] establish provable separation between SGD’s and Adagrad’s convergence to stationary
 1715 points of non-convex objectives. They show that when objective functions exhibit certain
 1716 geometric properties, measured by their coordinate-wise smoothness constants, Adagrad’s
 1717 upper bounds are lower than SGD’s lower bounds by a factor of d [21]. Our analysis builds
 1718 on that of Jiang et al. [21] and Liu et al. [25], as a consequence of our results is that the
 1719 EGOP reparameterization proposed in this work acts to transform objectives into the setting
 1720 identified by Jiang et al. [21] where Adagrad’s convergence guarantees compare most favorably
 1721 with SGD’s.

1722 Ling et al. [24] develop a rotationally equivariant extension of Adam, termed VectorAdam,
 1723 in an effort to reduce axis-aligned artifacts that arise when using adaptive methods to solve
 1724 geometric optimization problems such as differentiable rendering. This method targets ap-
 1725 plications when the problem parameters $\theta \in \mathbb{R}^{r \cdot n}$ represent a collection of r different vectors
 1726 in \mathbb{R}^n . VectorAdam uses the squared gradient norm of each n -dimensional vector to rescale
 1727 the learning rates of all entries in each of the r vectors comprising θ , making the algorithm
 1728 equivariant to transformations of the form $Q \text{mat } \theta$, where $\text{mat } \theta \in \mathbb{R}^{n \times r}$ is the reshaping of
 1729 $\theta \in \mathbb{R}^{r \cdot n}$ and $Q \in \mathbb{R}^{n \times n}$ is orthonormal [24].

1730 *Change-of-basis for Optimization.* Several recent works propose that when using Adam and
 1731 its variants to train neural networks, different choices of orthonormal transformation can re-
 1732 duce the computational cost associated with these algorithms and improve their performance
 1733 [27, 37, 43]. Gupta et al. [17] introduced Shampoo, an efficient preconditioning method for
 1734 optimizing tensor spaces. Vyas et al. [37] formalized a connection between Shampoo and a
 1735 variant of Adagrad, leading to a new proposed method called SOAP. Designed for training
 1736 neural networks, SOAP computes an orthonormal reparameterization based on the singular
 1737 vectors of the matrix-valued network gradients and performs optimization in this basis. Vyas
 1738 et al. [37] empirically examine the performance of SOAP and find that it outperforms both
 1739 Adam and Shampoo in LLM pre-training. Zhao et al. [43] propose GaLore, a method that
 1740 simultaneously performs reparameterization and dimensionality reduction. GaLore computes
 1741 a similar orthogonal basis to that used in SOAP, but instead of a full-dimensional change-
 1742 of-basis GaLore retains only leading basis vectors in order to reduce dimension [43]. Maes

1743 et al. [27] empirically study Adam’s rotational sensitivity and examine the power of exist-
 1744 ing geometric assumptions (such as those leveraged by Xie et al. [40]) to explain Adam’s
 1745 performance when training transformer architectures. They propose an orthonormal repara-
 1746 meterization, similar to those used by SOAP and GaLore, and show empirically that this
 1747 can improve Adam’s performance [27].

1748 SOAP and GaLore both call for periodically re-computing the change-of-basis matrices.
 1749 In their experiments, the fewest re-computations performed by Vyas et al. [37] was once every
 1750 100 batches, and they show that the performance gap between SOAP and Adam narrows as
 1751 the number of batches between re-computations increases. In contrast, our we show that with
 1752 our proposed choice of basis, a single up-front computation of the change-of-basis suffices to
 1753 improve the performance of both Adam and Adagrad in a variety of settings, as shown in
 1754 Section 6.

1755 Outside of the domain neural network training, several works have considered data-driven
 1756 methods for performing dimensionality reduction when optimizing objectives with low-rank
 1757 EGOP matrices, including the works of Cartis et al. [3] and Cosson et al. [7]. For functions
 1758 with low-dimensional active subspaces, Cartis et al. [3] study using the EGOP eigenbasis
 1759 to reparameterize and reduce the dimension of optimization problems. They demonstrate
 1760 that this approach yields computational speedups for loss functions whose EGOP matrix is
 1761 low-rank. Cosson et al. [7] develop gradient descent algorithms which leverage the EGOP
 1762 eigenbasis. Their method first estimates the EGOP and computes the corresponding leading
 1763 r -eigenvectors, and then performs gradient descent steps along the directions of these r vectors.
 1764 In the setting of exactly low-rank functions, they prove convergence results illustrating that
 1765 this approach improves the dimensional dependency of gradient descent.

1766 **Appendix G. Kronecker Product Constraints in SOAP/Shampoo.**

1767 In this section, we provide a detailed comparison of the change of basis employed by EGOP
 1768 reparameterization and those performed by SOAP and Shampoo [37, 17]. We instantiate an
 1769 example of EGOP reparameterization for an objective whose parameters can be viewed as
 1770 matrices, as this highlights some of the conceptual differences between the methods.

1771 In many key applications of adaptive algorithms, including training neural networks, the
 1772 parameters over which optimization occurs are matrix-valued. As a simple illustration, con-
 1773 sider optimizing a single-layer linear fully-connected network. Let n_{in} denote the number of
 1774 input features to the layer, and n_{out} denote the number of output features. This network
 1775 then has $n_{\text{in}} \cdot n_{\text{out}}$ parameters, which can be expressed either as a vector, $\theta \in \mathbb{R}^{n_{\text{in}}n_{\text{out}}}$, or as
 1776 a matrix, denoted $\text{mat}(\theta) \in \mathbb{R}^{n_{\text{in}} \times n_{\text{out}}}$.

1777 Denote the training data and labels by $A \in \mathbb{R}^{n_{\text{in}} \times n_{\text{samples}}}$ and $Y \in \mathbb{R}^{n_{\text{out}} \times n_{\text{samples}}}$, and
 1778 consider minimizing loss function

$$1779 \quad f(\theta) = \frac{1}{2} \|\text{mat}(\theta)^\top A - Y\|_F^2.$$

1780 Similarly, the vector-valued gradients $\nabla f(\theta) \in \mathbb{R}^{n_{\text{in}}n_{\text{out}}}$ can also be viewed as matrices,
 1781 $\text{mat}(\nabla f(\theta)) \in \mathbb{R}^{n_{\text{in}} \times n_{\text{out}}}$.

1782 In the method proposed in this work, we consider gradients to be vector-valued when form-
 1783 ing the EGOP. Thus for this single-layer objective, $\text{EGOP}(f) \in \mathbb{R}^{n_{\text{in}}n_{\text{out}} \times n_{\text{in}}n_{\text{out}}}$. We emphasize
 1784 this vector-view of gradients for clarity, because the EGOP matrix has distinct eigenvectors

1785 from the expectation of the matrix product, $\text{mat}(\nabla f(\theta)) \text{mat}(\nabla f(\theta))^\top$. The related methods
 1786 discussed in the introduction, including SOAP and GaLore, consider transformations by the
 1787 matrices

$$\begin{aligned} 1788 \quad Q_L &= \text{eigenvectors}(\mathbb{E}[\text{mat}(\nabla f(\theta)) \text{mat}(\nabla f(\theta))^\top]) \\ 1789 \quad Q_R &= \text{eigenvectors}(\mathbb{E}[\text{mat}(\nabla f(\theta))^\top \text{mat}(\nabla f(\theta))]) \end{aligned}$$

1790 and closely related transformations [27, 37, 43]. In general the eigenvectors in Q_L, Q_R corre-
 1791 spond to different transformations than the eigenvectors of the EGOP. In particular, letting
 1792 V denote the eigenbasis of the EGOP formed from vector-valued gradients, $V \neq Q_L \otimes Q_R$.
 1793 Moreover, the class of orthonormal matrices obtainable from the EGOP eigenbasis is strictly
 1794 more general than the class $Q_L \otimes Q_R$ for pairs of orthogonal matrices Q_L, Q_R , as formalized
 1795 below:

1796 **Lemma G.1.** *For any orthogonal matrices $Q_L \in \mathbb{R}^{n_{in} \times n_{in}}, Q_R \in \mathbb{R}^{n_{out} \times n_{out}}$ and any objective*
 1797 *function $f: \mathbb{R}^{n_{in}n_{out}} \rightarrow \mathbb{R}$, there exists orthogonal $V \in \mathbb{R}^{n_{in}n_{out} \times n_{in}n_{out}}$ such that*

$$1798 \quad \forall \theta \in \mathbb{R}^{n_{in}n_{out}}, \quad Q_L^\top \text{mat}(\nabla f(\theta)) Q_R = \text{mat}(V \nabla f(\theta)).$$

1799 *However, there exist values of n_{in}, n_{out} , objective functions $f: \mathbb{R}^{n_{in}n_{out}} \rightarrow \mathbb{R}$, and orthog-*
 1800 *onal matrices $V \in \mathbb{R}^{n_{in}n_{out} \times n_{in}n_{out}}$ such that no orthogonal matrices $Q_L \in \mathbb{R}^{n_{in} \times n_{in}}, Q_R \in$
 1801 $\mathbb{R}^{n_{out} \times n_{out}}$ satisfy*

$$1802 \quad \forall \theta \in \mathbb{R}^{n_{in}n_{out}}, \quad Q_L^\top \text{mat}(\nabla f(\theta)) Q_R = \text{mat}(V \nabla f(\theta)).$$

1803 *Proof.* First we show that for any Q_L, Q_R , there exists suitable V satisfying the property.
 1804 We begin by noting that for any matrices A, B, C of compatible dimension,

$$1805 \quad \text{vec}(ABC) = (C^\top \otimes A) \text{vec}(B).$$

1806 Thus for any Q_L, Q_R ,

$$1807 \quad \text{vec}(Q_L^\top \text{mat}(\nabla f(\theta)) Q_R) = (Q_R^\top \otimes Q_L^\top) \nabla f(\theta)$$

1808 where the second equality uses the fact that trivially $\text{vec}(\text{mat}(\nabla f(\theta))) = \nabla f(\theta)$. The Kron-
 1809 ecker product of orthogonal matrices is orthogonal, so choice of $V = (Q_R^\top \otimes Q_L^\top)$ satisfies the
 1810 desired property.

1811 We now show there exist orthogonal matrices $V \in \mathbb{R}^{n_{in}n_{out} \times n_{in}n_{out}}$ such that no orthogonal
 1812 matrices $Q_L \in \mathbb{R}^{n_{in} \times n_{in}}, Q_R \in \mathbb{R}^{n_{out} \times n_{out}}$ satisfy

$$1813 \quad Q_R \otimes Q_L^\top = V.$$

1814 This is a consequence of the fact that not all orthogonal matrices admit Kronecker product
 1815 factorizations. Here we give one specific construction.

Let $\vec{1}_d$ denote the all-ones vector in \mathbb{R}^d . Consider V with leading column

$$v_1 = \vec{1}_{n_{in}n_{out}} / \sqrt{n_{in}n_{out}}$$

and second column

$$v_2 = (n_{\text{in}}n_{\text{out}})^{-1/2} \cdot \text{concatenate}(\vec{1}_{n_{\text{in}}n_{\text{out}}/2}, -\vec{1}_{n_{\text{in}}n_{\text{out}}/2})$$

1816 . Given the entries of v_1 are identical, this implies all the entries in the first column of Q_R
 1817 are identical, and thus that the first n_{in} columns of V comprise concatenated copies of Q_L^T .
 1818 However this contradicts the fact that the entries of v_2 are identical in magnitude and nonzero
 1819 but have positive sign for the first $n_{\text{in}}n_{\text{out}}/2$ entries and negative sign for the rest. Thus no
 1820 orthogonal matrices V with such first and second columns can be decomposed into $Q_R \otimes Q_L^T$,
 1821 so in particular there exists some vector $z \in \mathbb{R}^{n_{\text{in}}n_{\text{out}}}$ such that

$$1822 \quad Q_R \otimes Q_L^T z \neq Vz.$$

1823 Thus for any $f(\cdot)$ such that there exists θ with $\nabla f(\theta) = z$, it holds that for this value

$$1824 \quad Q_L^T \text{mat}(\nabla f(\theta)) Q_R = \text{mat}\left((Q_R \otimes Q_L^T) \nabla f(\theta)\right) \neq \text{mat}(V \nabla f(\theta)). \quad \blacksquare$$