

RAPID GRASSMANNIAN AVERAGING WITH CHEBYSHEV POLYNOMIALS

Anonymous authors

Paper under double-blind review

ABSTRACT

We propose new algorithms to efficiently average a collection of points on a Grassmannian manifold in both the centralized and decentralized settings. Grassmannian points are used ubiquitously in machine learning, computer vision, and signal processing to represent data through (often low-dimensional) subspaces. While averaging these points is crucial to many tasks (especially in the decentralized setting), existing methods unfortunately remain computationally expensive due to the non-Euclidean geometry of the manifold. Our proposed algorithms, Rapid Grassmannian Averaging (RGrAv) and Decentralized Rapid Grassmannian Averaging (DRGrAv), overcome this challenge by leveraging the spectral structure of the problem to rapidly compute an average using only small matrix multiplications and QR factorizations. We provide a theoretical guarantee of optimality and present numerical experiments which demonstrate that our algorithms outperform state-of-the-art methods in providing high accuracy solutions in minimal time. Additional experiments showcase the versatility of our algorithms to tasks such as K -means clustering on video motion data, establishing RGrAv and DRGrAv as powerful tools for generic Grassmannian averaging.

1 INTRODUCTION

Grassmannian manifolds, which represent sets of K -dimensional linear subspaces of N -dimensional spaces (Edelman et al., 1998), have been used extensively in machine learning (Huang et al., 2018; Zhang et al., 2018; Slama et al., 2015), computer vision (Harandi et al., 2013; Lui & Beveridge, 2008; Turaga et al., 2011), and signal processing (Gallivan et al., 2003; Mondal et al., 2007; Xu & Hassibi, 2008). Applications include Principal Component Analysis (PCA) (Jolliffe & Cadima, 2016), low-rank matrix completion (Keshavan et al., 2010), multi-task feature learning (Mishra et al., 2019), clustering (Gruber & Theis, 2006), array processing (Love et al., 2003; DeLude et al., 2022), and distance metric learning (Meyer et al., 2009).

An essential primitive operation is finding an average of a collection of points on the manifold. There are several distinct yet reasonable definitions for an average of points on a Grassmannian (Marrinan et al., 2014). Arguably the most natural analog of the Euclidean mean for points on a Riemannian manifold (such as a Grassmannian) is the Fréchet (or Karcher) mean, defined as the point which minimizes the sum of squared distances to all sample points (Fréchet, 1948). Unfortunately, the Fréchet mean rarely admits a closed form solution, instead necessitating approximation via iterative algorithms (Jeuris et al., 2012). Such algorithms are often computationally expensive, scale poorly with dimension, and are not easily decentralized.

The induced arithmetic mean (IAM) is an alternative manifold average computed by first determining the Euclidean mean of the manifold sample points once embedded “naturally” in some Euclidean space and subsequently projecting this Euclidean mean “naturally” back onto the manifold. For Grassmannian manifolds, the standard embedding is the set of projection matrices and the standard projection operation is simply the closest matrix by Frobenius distance (Sarlette & Sepulchre, 2009). This manifold average may be computed much more efficiently in practice and lends itself well to decentralization as the Euclidean mean may be computed by average consensus (Nedic & Ozdaglar, 2009). In such a decentralized setting, computing this Euclidean mean would be the only operation that requires communication in order to compute the IAM.

054 As the dimensionality of data grows, it becomes increasingly important to consider decentralized
055 algorithms (Nedić et al., 2018) as data might be spread across many machines and only be accessible
056 for processing via distributed algorithms (Beltrán et al., 2023). Distributed computation might be
057 required as well for situations where the data associated to each agent must be treated with privacy
058 protections, where aggregation of all data onto a single node may be prohibited (Han et al., 2017).
059 While a central server is sometimes employed in this regime, it is similarly common for the use of
060 such a server to be infeasible or simply inefficient when compared to fully decentralized approaches
061 (Sun et al., 2021; Feller et al., 2012).

062 We propose a novel method to efficiently compute the IAM of a collection of points on a Grass-
063 mannian manifold. Our method is highly amenable to decentralization, meaning it can be readily
064 deployed to multi-agent systems or used in data centers operating on big data. Our algorithms op-
065 erate similarly to the famous power method, with the distinction that Chebyshev polynomials are
066 employed to leverage a “dual-banded” property of the problem in order to achieve never-before-
067 seen efficiency in computation and communication. We demonstrate merit through a theoretical
068 guarantee on the optimality of our approach among a class of polynomial-based algorithms, syn-
069 thetic numerical experiments comparing our algorithms against state-of-the-art, and experiments on
070 real-world problems showcasing the versatility of our algorithms.

071 072 073 074 075 2 RELATED WORK

076
077
078
079
080 The problem of computing an appropriate average on specific manifolds has been investigated for
081 many different manifolds, e.g., spheres S^N , special orthogonal matrices $SO(N)$, Stiefel matrices
082 $St(N, K)$, even Grassmannian points $Gr(N, K)$ (Downs, 1972; Buss & Fillmore, 2001; Galperin,
083 1993; Hueper & Manton, 2004; Absil et al., 2004; Moakher, 2002; Fiori et al., 2014; Yun, 2018;
084 Hauberg et al., 2014). Focus is often given to the Fréchet mean (Chakraborty et al., 2020; Cheng
085 et al., 2016; Le, 2001), however alternatives are becoming increasingly more popular (Fletcher et al.,
086 2008; 2009; Arnaudon et al., 2012; Marrinan et al., 2014; Chakraborty & Vemuri, 2015; Lee & Jung,
087 2024). Similarly, the problem of consensus on a manifold in a multi-agent setting has been explored
088 in works such as Sepulchre (2011); Tron et al. (2012).

089 There have been several algorithms proposed for decentralized optimization on manifolds such as
090 Grassmannians. Sarlette & Sepulchre (2009) proposes a decentralized gradient-based algorithm
091 to solve the problem of computing the IAM for connected compact homogeneous manifolds, e.g.
092 $SO(N)$ and $Gr(N, K)$. Deng & Hu (2023) proposes two decentralized gradient-based algorithms
093 for general optimization problems on Riemannian manifolds. Mishra et al. (2019) proposes a de-
094 centralized gradient-based gossip algorithm for general optimization problems on a Grassmannian
095 manifold. Similar works include Chen et al. (2021; 2023); Zhang & Sun (2017).

096 A problem which is closely related to Grassmannian averaging is that of PCA. Ye & Zhang (2021)
097 proposes the DeEPCA algorithm to solve the decentralized PCA problem. While computing a Grass-
098 mannian average is not the intended application of DeEPCA, it may be adapted to this task fairly
099 naturally. Gang et al. (2021); Gang & Bajwa (2022); Froelicher et al. (2023) similarly propose
100 distributed algorithms for PCA; for a more comprehensive review of this field, see Wu et al. (2018).

101 As we will see later in this paper, the problem of Grassmannian averaging is related to the problem
102 of spectral estimation (see Section 3.1). There exist many tailored algorithms in this field, varying
103 based on factors such as eigenvalue vs. eigenvector estimation, matrix rank, symmetry, estimation
104 of leading vs trailing quantities, size of eigengap, etc (Liesen & Strakos, 2013; Lanczos, 1950;
105 Knyazev, 2001; Sleijsen & Van der Vorst, 2000; Zhou & Saad, 2007; Ghanem & Ghosh, 2007;
106 Martinsson & Tropp, 2020). Many such algorithms optimize for centralized computation, using
107 iterative and sometimes stochastic approaches. One of the most elegant solutions in this space is the
power method, from which we take inspiration (see Section 3.3).

3 BACKGROUND

3.1 AVERAGING SUBSPACES

Given a collection of M subspaces, our goal is to determine the average subspace as efficiently as possible. We choose to use the standard IAM definition of “average” (Sarlette & Sepulchre, 2009) as it leads to what we believe is the most efficient algorithm. Formally, let $\text{St}(N, K) := \left\{ \mathbf{U} \in \mathbb{R}^{N \times K} \mid \mathbf{U}^\top \mathbf{U} = \mathbf{I}_K \right\}$ be the set of $N \times K$ Stiefel matrices where $N \geq K$ and let $\text{Gr}(N, K) := \{[\mathbf{U}] \mid \mathbf{U} \in \text{St}(N, K)\}$ (where the equivalence is defined as $[\mathbf{U}] := \{\mathbf{U}\mathbf{Q} \mid \mathbf{Q} \in \text{St}(K, K)\}$) be the Grassmannian representing the set of all K -dimensional subspaces of \mathbb{R}^N . The average of our collection $\{\mathbf{U}_m\}_{m=1}^M$ is then denoted $[\bar{\mathbf{U}}]$ and defined by the following optimization problem

$$[\bar{\mathbf{U}}] := \underset{[\mathbf{U}] \in \text{Gr}(N, K)}{\text{argmin}} \left\| \left(\frac{1}{M} \sum_{m=1}^M \mathbf{U}_m \mathbf{U}_m^\top \right) - \mathbf{U} \mathbf{U}^\top \right\|_{\text{F}}^2 \quad (1)$$

Equation (1) may be manipulated algebraically to be interpreted equivalently in terms of the eigenvectors of $\bar{\mathbf{P}} := \frac{1}{M} \sum_{m=1}^M \mathbf{U}_m \mathbf{U}_m^\top$. Let

$$\bar{\mathbf{P}} = \tilde{\mathbf{V}} \tilde{\mathbf{\Lambda}} \tilde{\mathbf{V}}^\top = [\mathbf{V} \quad \mathbf{V}_\perp] \begin{bmatrix} \mathbf{\Lambda} & \mathbf{0} \\ \mathbf{0} & \mathbf{\Lambda}_\perp \end{bmatrix} \begin{bmatrix} \mathbf{V}^\top \\ \mathbf{V}_\perp^\top \end{bmatrix}$$

denote an eigendecomposition where $\mathbf{V} \in \text{St}(N, K)$ and the entries of $\tilde{\mathbf{\Lambda}}$ are non-increasing. Assuming $\lambda_K > \lambda_{K+1}$ (where λ_k denotes the k th largest eigenvalue of $\bar{\mathbf{P}}$), it can be shown that the solution to eq. (1) is precisely $[\bar{\mathbf{U}}] = [\mathbf{V}]$ (see Appendix B.2 for a proof). Consequently, determining the span of the leading K eigenvectors of $\bar{\mathbf{P}}$ is tantamount to solving eq. (1), which is the perspective we will later use to motivate our algorithms.

As we continue to discuss this problem, it is informative to keep in mind the following properties. The eigenvalues of $\bar{\mathbf{P}}$ are conveniently bounded by $\mathbf{0} \preceq \tilde{\mathbf{\Lambda}} \preceq \mathbf{I}_N$ and satisfy $\text{tr}(\tilde{\mathbf{\Lambda}}) = K$, which may be determined by inspection. As a result, λ_K, λ_{K+1} are bounded as $\frac{1}{N-K+1} \leq \lambda_K \leq 1$ and $0 \leq \lambda_{K+1} \leq \frac{K}{K+1}$. For convenience, we occasionally abuse notation to let $[\mathbf{X}]$ denote the Grassmannian equivalence class for the span of the columns of arbitrary (not necessarily Stiefel) matrix $\mathbf{X} \in \mathbb{R}^{N \times K}$.

3.2 AVERAGING SUBSPACES IN A DECENTRALIZED NETWORK

Decentralized or distributed optimization problems arise in numerous real-world scenarios where centralized approaches are impractical or undesirable. These problems are characterized by using information spread across multiple agents or nodes in a network. Motivating reasons include privacy, communication constraints, data storage limitations, scalability, etc.

In the context of this paper, consider the setting where there are M agents, each holding a subspace $[\mathbf{U}_m]$, connected by a some undirected communication graph \mathcal{G} . We then want each agent to learn the solution $[\bar{\mathbf{U}}]$ to eq. (1) under the restriction that each agent only communicate with their neighbors in \mathcal{G} .

Average consensus (AC) is a useful primitive in decentralized optimization to quickly approximate the average of real numbers in a decentralized manner (Nedic & Ozdaglar, 2009). If the m th agent in a decentralized network holds a matrix \mathbf{A}_m , AC allows each agent to approximate $\frac{1}{M} \sum_{m=1}^M \mathbf{A}_m$ with minimal rounds of neighbor-only communication. Unfortunately, the non-convex manifold structure of $\text{Gr}(N, K)$ precludes us from efficiently applying AC directly to solve eq. (1), as the Euclidean mean of elements from a non-convex set in general does not lay in said set. While we could have each agent compute the matrices $\mathbf{U}_m \mathbf{U}_m^\top$ and then use AC to approximate $\bar{\mathbf{P}}$, this would incur a communication cost of $\mathcal{O}(N^2)$ (the size of $\bar{\mathbf{P}}$) which may be much larger than $\mathcal{O}(NK)$. We could instead use AC to average the matrices \mathbf{U}_m with preferable communication cost $\mathcal{O}(NK)$, however

the arbitrary choice of representative Stiefel matrix U_m from the Grassmannian equivalence class $[U_m]$ makes this approach ill-posed.

In order to achieve the $\mathcal{O}(NK)$ communication cost without being ill-posed, one can have all agents compute $U_m U_m^\top X$ and then use AC to approximate $\bar{P}X$, where $X \in \mathbb{R}^{N \times K}$ is some matrix agreed upon by all agents a priori. Section 3.3 elaborates on how quantities of the form $\bar{P}X$ can be used to estimate the desired leading eigenvectors. In practice, the requirement that all agents agree upon X might be overly strict; in many scenarios, it often suffices to have each agent m have instance X_m which are all approximately equal (i.e. there exists some X for which $X_m \approx X$ for all $m \in [M]$).

After one iteration of an algorithm, each agent will have some local approximation of the quantity $\bar{P}X$. If the matrix X is retained in memory for each agent, then AC may be applied to a linear combination of the $\bar{P}X$ approximations and X to approximate some quantity $\bar{P}(a\bar{P}X + bX) = (a\bar{P}^2 + b\bar{P})X$. Applying this logic recursively reveals that such an algorithm can approximate $f_t(\bar{P})X$ after t iterations, where $f_t \in \mathcal{P}_t$ is a t th order polynomial; in Section 4 we will consider more thoroughly these polynomials and how they can translate to desirable algorithms.

Gradient tracking is a famous technique in decentralized optimization that improves upon the convergence rate of AC-based methods (Shi et al., 2015; Xu et al., 2015; Qu & Li, 2017; Deng & Hu, 2023; Ye & Zhang, 2021). In essence, it sets up a recursion using standard AC whose fixed point satisfies both a consensus condition (meaning all agents agree) and a stationarity condition (meaning the solution is locally optimal). We will later employ a form of gradient tracking over quantities of the form $f_t(\bar{P})X$ for our decentralized algorithms.

3.3 THE POWER METHOD

The power method is a classical algorithm to estimate the leading eigenspace of a positive semidefinite matrix $A \in \mathbb{R}^{N \times N}$. A single power iteration applies the matrix A and orthonormalizes (for numerical stability) the result, e.g.

$$U^{(t)} = \text{QR}\left(AU^{(t-1)}\right),$$

where $\text{QR}(\cdot)$ computes the $Q \in \mathbb{R}^{N \times K}$ matrix from a QR factorization of the argument. The power method loop may be unrolled (thanks to the property $\text{QR}(X\text{QR}(Y)) = \text{QR}(XY)$) to reveal the following form

$$U^{(t)} = \text{QR}\left(\underbrace{AA \cdots AA}_{t \text{ times}} U^{(0)}\right) = \text{QR}\left(A^t U^{(0)}\right)$$

Let $V_* \in \text{St}(N, K)$ be a basis for the leading K eigenvectors of A . For random initialization $U^{(0)} \in \mathbb{R}^{N \times K}$, the span of $U^{(t)}$ converges to the span of the V_* provided $\text{rank}\left(V_*^\top U^{(0)}\right) = K$ (Golub & Van Loan, 2013, Chapter 8.2).

At iteration t , the power method effectively applies the function $f_t(\lambda) = \lambda^t$ to each eigenvalue of A . Consider for example the case where $\lambda_K(A) = 1$ and $\lambda_K(A) - \lambda_{K+1}(A) > 0$. As t increases, the ratio between trailing and leading eigenvalues of A^t shrinks exponentially; formally, for $1 \leq k \leq K$ and $K + 1 \leq \ell \leq N$ we have the following

$$\frac{\lambda_\ell(A^t)}{\lambda_k(A^t)} \leq \frac{\lambda_{K+1}(A^t)}{\lambda_K(A^t)} = \left(\frac{\lambda_{K+1}(A)}{\lambda_K(A)}\right)^t < 1$$

While the power method works well, large values of λ_{K+1} can slow convergence. We will later show how polynomials other than λ^t can overcome this shortcoming and use them in our algorithms.

4 METHODS

4.1 MOTIVATION

The goal of the RGrAv algorithms is to solve eq. (1) as efficiently as possible. Recall from Section 3 that determining the span of the K leading eigenvectors of \bar{P} is tantamount to solving eq. (1). Mo-

216
217
218
219
220
221
222
223
224
225
226
227
228
229
230
231
232
233
234
235
236
237
238
239
240
241
242
243
244
245
246
247
248
249
250
251
252
253
254
255
256
257
258
259
260
261
262
263
264
265
266
267
268
269

tivated by the decentralized setting and the power method, we restrict our consideration to iterative algorithms which after t iterations compute $f_t(\bar{\mathbf{P}})\bar{\mathbf{U}}^{(0)}$ for some t th-order polynomial f_t and initial estimate $\bar{\mathbf{U}}^{(0)}$ (see Section 3.2). If we can choose f_t such that $f_t(\bar{\mathbf{P}})$ has its trailing $N - K$ eigenvalues significantly reduced compared to its leading K eigenvalues (relative to initial $\bar{\mathbf{P}}$), then the span of the matrix product $f_t(\bar{\mathbf{P}})\bar{\mathbf{U}}^{(0)}$ will be approximately our desired solution $[\mathbf{V}]$ for arbitrary initial $\bar{\mathbf{U}}^{(0)}$. We then focus on choosing such a so-called “noise-canceling” polynomial f_t , in the sense that the trailing eigenvalues get “canceled”.

We consider the case where the spectrum of $\bar{\mathbf{P}}$ is dual-banded, i.e. there exists some $0 < \alpha < \beta < 1$ such that $\mathbf{\Lambda}_\perp \preceq \alpha \mathbf{I}_{N-K}$, $\beta \mathbf{I}_K \preceq \mathbf{\Lambda}$, and $\beta - \alpha \gg 0$ (e.g. $\beta - \alpha = \frac{1}{3}$). These values α, β are unknown, but may be estimated heuristically from domain knowledge. This situation can arise, for instance, when points are normally distributed on the manifold (see Section 5.1), where the heuristic estimation of α, β comes from estimation of the variance of the dataset. For simplicity, we refer to the intervals $[0, \alpha]$ and $[\beta, 1]$ as the “stop-band” and “pass-band”, respectively. Similar to the power method, we would like our polynomial f_t to decrease the ratio between eigenvalues in the stop-band relative to eigenvalues in the pass-band. However, unlike the power method, knowledge of this dual-banded structure (even heuristically) allows us to choose polynomials which optimize the worst case value of this ratio criteria. Leveraging this spectral structure is how we will choose our optimal polynomials f_t^* (see Theorem 1).

There is, however, an important consideration for high-dimensional data. In the case where $N \gg MK$ we are guaranteed that the nullspace of $\bar{\mathbf{P}}$ is (at least $N - MK$) high-dimensional by rank subadditivity, meaning there will be a cluster of eigenvalues at 0. For this reason, we will constrain our polynomials f_t to always satisfy $f_t(0) = 0$. Since the ratio criteria of Theorem 1 is invariant to scaling of f_t , we provide one final constraint of $f_t(1) = 1$ simply for uniqueness of solution and numerical stability.

Theorem 1. For $t \geq 1$, the minimization problem

$$\underset{f_t \in \mathcal{P}'_t}{\text{minimize}} \frac{\max_{\lambda \in [0, \alpha]} |f_t(\lambda)|}{\min_{\lambda \in [\beta, 1]} |f_t(\lambda)|}, \quad (2)$$

where \mathcal{P}'_t is the set of t th order polynomials such that $f_t(0) = 0$ and $f_t(1) = 1$, is solved by

$$f_t^*(\lambda) = \prod_{s=0}^{t-1} \frac{\lambda - r_{s,t}}{1 - r_{s,t}}, \quad r_{s,t} := \alpha \frac{\cos\left(\frac{\pi(s+1/2)}{t}\right) + \cos\left(\frac{\pi}{2t}\right)}{1 + \cos\left(\frac{\pi}{2t}\right)}$$

which is a modification of a Chebyshev polynomial of the first kind.¹

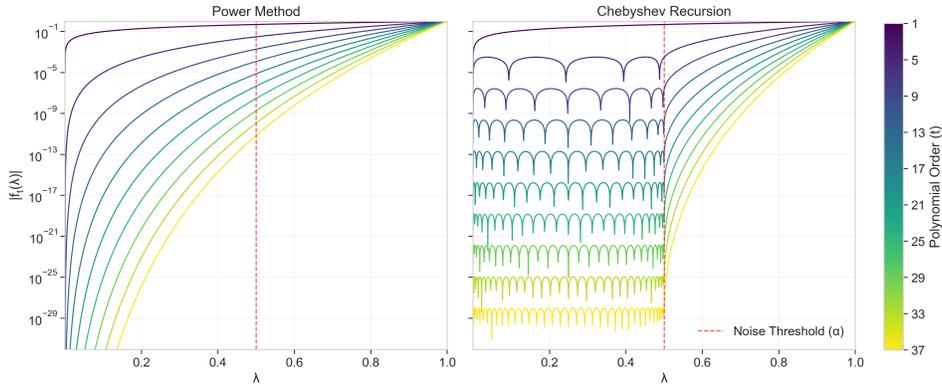


Figure 1: A visual comparison between the power method and our method. Each method’s corresponding t th-order polynomial is applied to the eigenvalues λ in the domain $[0, 1]$. The Chebyshev recursion with threshold parameter $\alpha = 0.5$ results in the polynomial oscillations being reduced and flattened in the range $[0, \alpha]$.

¹A proof may be found in Appendix B.1

4.2 RGRAV ALGORITHMS

The polynomial f_t^* may be exactly implemented by iteratively multiplying each factor in the product given in Theorem 1; the “finite” variants of the RGrAv algorithms (Algorithms 5 and 6) do precisely this. While this approach may be acceptable when the number of iterations t is known in advance, intermediate solutions can be quite sub-optimal.

Ideally, one would be able to describe f_t^* in terms of only a constant number of previous f_s^* , e.g. f_{t-1}^*, f_{t-2}^* . This would yield an efficient algorithm with optimal intermediate solutions whose memory/compute costs do not grow as $t \rightarrow \infty$. Unfortunately, this is not the case; fortunately, empirically, f_t^* is *well-approximated* in terms of f_{t-1}^*, f_{t-2}^* (see Figure 1). For $t \geq 2$, coefficients a_t, b_t, c_t are chosen such that

$$\tilde{f}_t^*(\lambda) = a_t((\lambda + b_t)f_{t-1}(\lambda) + c_t f_{t-2}(\lambda))$$

matches $f_t^*(\lambda)$ in its leading three terms, i.e. $f_t^*(\lambda) - \tilde{f}_t^*(\lambda) \in \mathcal{P}_{t-3}$ (see Algorithm 7); the “asymptotic” variants of the RGrAv algorithms (Algorithms 1 and 4) use this \tilde{f}_t^* .

As discussed in Section 3.3, orthonormalization must occur periodically for numerical stability. To minimize the frequency of the orthonormalization schedule, our algorithms effectively cache the operation of the most recent exact orthonormalization to a matrix \mathbf{S} and then efficiently approximate the orthonormalization procedure for iterations between the schedule by left-application of \mathbf{S} . We choose the QR factorization for our orthonormalization method, however alternative methods would be acceptable.

In the centralized setting, the RGrAv algorithms need not worry about inaccuracy from average consensus and so the order of operations focuses on minimizing the number of computations performed. In the decentralized setting, the order of operations focuses on minimizing the error in the gradient tracking procedure. Additionally, in the decentralized setting one must take care to use a “stable” orthonormalization method which will not change drastically with small perturbations in the input. We omit the nuances of numerical linear algebra that lead to this problem (see Golub & Van Loan (2013, Chapter 5) for more information) and simply present Algorithm 2, which is a stable wrapper for any implementation of a possibly unstable QR factorization.

Algorithm 1 Asymptotic DRGrAv (Decentralized Rapid Grassmannian Averaging)

Input: $\alpha \in [0, 1)$, $\{\bar{\mathbf{U}}_m^{(0)}\}_{m=1}^M$, $\{\mathbf{U}_m\}_{m=1}^M$

Output: $\bar{\mathbf{U}}^{(T)}$

$\mathbf{S}_m \leftarrow \mathbf{I}_K$

for $t = 1, 2, 3, \dots$ **do**

$\mathbf{A}_m^{(t)} \leftarrow \mathbf{U}_m \mathbf{U}_m^\top \bar{\mathbf{U}}^{(t-1)}$

▷ Local Power Iteration

if $t = 1$ **then**

$\mathbf{Y}_m^{(1)} \leftarrow \mathbf{A}_m^{(1)}$

▷ $[\mathbf{Y}^{(1)}] \approx [\bar{\mathbf{P}}\bar{\mathbf{U}}^{(0)}]$

$\mathbf{Z}_m^{(1)} \leftarrow \mathbf{Y}_m^{(1)}$

▷ Gradient Tracking

else

$a_t, b_t, c_t \leftarrow \text{ChebyshevCoefficients}(t, \alpha)$

▷ See Algorithm 7

$\mathbf{Y}_m^{(t)} \leftarrow a_t(\mathbf{A}_m^{(t)} + b_t \bar{\mathbf{U}}_m^{(t-1)} + c_t \bar{\mathbf{U}}_m^{(t-2)})$

▷ $[\mathbf{Y}^{(t)}] \approx [\tilde{f}_t^*(\bar{\mathbf{P}})\bar{\mathbf{U}}^{(0)}]$

$\mathbf{Z}_m^{(t)} \leftarrow \hat{\mathbf{Z}}_m^{(t-1)} + \mathbf{Y}_m^{(t)} - \mathbf{Y}_m^{(t-1)}$

▷ Gradient Tracking

end if

$\hat{\mathbf{Z}}_m^{(t)} = \text{AverageConsensus}(\mathbf{Z}_m^{(t)})$

if t is on the orthonormalization schedule **then**

$\bar{\mathbf{U}}_m^{(t)}, \mathbf{S}_m \leftarrow \text{StableQR}(\hat{\mathbf{Z}}_m^{(t)})$

▷ (Numerical Stability)

else

$\bar{\mathbf{U}}_m^{(t)} \leftarrow \hat{\mathbf{Z}}_m^{(t)} \mathbf{S}_m$

▷ (Numerical Stability)

end if

$\bar{\mathbf{U}}_m^{(t-1)} \leftarrow \bar{\mathbf{U}}_m^{(t-1)} \mathbf{S}_m$

▷ (Numerical Stability)

end for

Algorithm 2 *StableQR***Input:** $\hat{Z} \in \mathbb{R}^{N \times K}$ **Output:** $U \in \text{St}(N, K)$, $S \in \mathbb{R}^{K \times K}$ $Q, R \leftarrow \text{QR}(\hat{Z})$ ▷ Arbitrary QR implementation $D \leftarrow \text{sgn}(\text{Diag}(R))$ $U \leftarrow QD$ $S \leftarrow R^{-1}D$ ▷ Upper triangular inverse

5 EXPERIMENTS

5.1 DECENTRALIZED GRASSMANNIAN AVERAGING

In these experiments, we consider the problem where a network of M connected agents each has a local instance of a Grassmannian basis $U_m \in \text{St}(N, K)$ and we would like for all agents to learn an average Grassmannian basis of all U_m in a strictly decentralized manner (i.e. there is no central server, all communication is neighbor-to-neighbor). Our experiments had parameters $M = 64$, $N = 150$, $K = 30$. To demonstrate the practicality of DRGrAv in both well-connected and sparse communication graphs, we performed experiments for two communication graphs: the hypercube graph and the cycle graph.

We compared DRGrAv to several alternative methods for Grassmannian averaging. Given below are the algorithms, their sources, and considerations for their tuning such that the comparison would be fair.

DRGrAv (this paper): Contrary to the following algorithms for which hyperparameters were chosen over large ranges to be empirically optimal, the hyperparameter α is chosen here heuristically as 0.15. We also choose to use the approximate asymptotic variant of DRGrAv, and set the orthonormalization schedule to orthonormalize at every iteration (to match DeEPCA). We believe it is unrealistic in practice to exactly know the optimal choice of α , so by comparing our heuristically-tuned algorithm against optimally-tuned competitors (detailed below) we hope to demonstrate that our algorithm is competitive against alternatives, regardless of however optimal their hyperparameter tuning.

DeEPCA (Ye & Zhang, 2021): While this algorithm is not intended for decentralized Grassmannian averaging, we found that it could be easily adapted to this setting and gave competitive results: one simply substitutes $U_m U_m^\top$ for the paper’s A_j . Also, for numerical stability, the paper’s QR + SignAdjust procedure is replaced with the StableQR procedure. At a high level, these are all that’s needed to adapt the method; see the `deepca.py` for a comprehensive algorithm description.

DPRGD/DPRGT (Deng & Hu, 2023): These algorithms are adapted to the decentralized Grassmannian averaging problem by using $-\frac{1}{2} \left\| U_m^\top \bar{U}_m^{(t)} \right\|_F^2$ for the paper’s f_i . These algorithms each have a single hyperparameter for step size (referred to as α in Deng & Hu (2023)), which was chosen for each algorithm (up to precision of 2 significant figures) by searching for the value in $[10^{-4}, 10^3]$ which approximately minimized the MSE; the solutions were on the order of 10^0 .

COM (Consensus Optimization on Manifolds) (Sarlette & Sepulchre, 2009): This algorithm is the discrete-time variant of the continuous-time dynamics presented in Equation 20 of Sarlette & Sepulchre (2009). There is a single hyperparameter for step size (referred to as α in Sarlette & Sepulchre (2009)), which was chosen (up to precision of 2 significant figures) by searching for the value in $[10^{-4}, 10^3]$ which approximately minimized the MSE; the solutions were on the order of 10^{-1} .

Gossip (Mishra et al., 2019): This is Algorithm 1 of Mishra et al. (2019) where $\frac{1}{4} \sum_{m=1}^M \sum_{n \in \mathcal{N}(m)} d^2 \left(\left[\bar{U}_m^{(t)} \right], \left[\bar{U}_n^{(t)} \right] \right)$ is used as the paper’s g (where $\mathcal{N}(m)$ is the set of neighbors of m). This algorithm is a gossip algorithm, *not* an average-consensus-based algorithm. As a result, to keep the comparison fair we let each agent perform a gradient step in parallel for each round of consensus the other algorithms perform. In precise terms, during each round of con-

sensus this algorithm will select edges from the graph uniformly at random until there no longer remain any 2 neighboring agents who both have not yet been selected; each of these agents then performs a gradient step and the process repeats. There are 2 hyperparameters² for step size a and b , which were chosen (up to precision of 2 significant figures) by searching for values in $a \in [10^{-4}, 10^3]$, $b \in [10^{-8}, 10^0]$ which approximately minimized the MSD; the solutions were a on the order of 10^0 and b on the order of 10^{-4} .

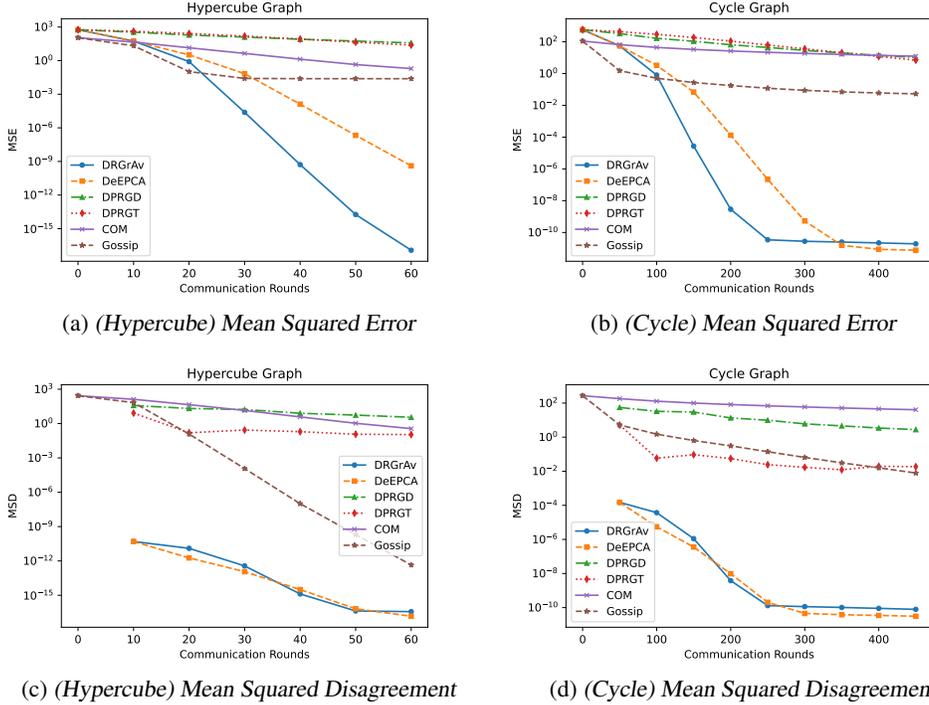


Figure 2: Plots of Mean Squared Error/Disagreement for the example decentralized Grassmannian averaging problem. DRGrAv is our proposed algorithm, DeEPCA is from Ye & Zhang (2021), DPRGD and DPRGT are from Deng & Hu (2023), COM is from Sarlette & Sepulchre (2009), and Gossip is from Mishra et al. (2019). The units of the x axes are communication rounds, not algorithm iterations.

Table 1: Comparison of runtimes for various algorithms (using the hypercube graph). The first five data columns display time (in milliseconds) until the MSE across agents goes below the given tolerance. COM and Gossip do not in general converge to any specific point, so their metric for tolerance is instead MSD. The minimal quantity in each column is bolded. The final column represents time per algorithm iteration. These decentralized algorithms are not truly run on separate devices, only simulated as such, so these runtimes should be interpreted broadly as general evidence that DRGrAv would perform well in a true decentralized setting.

Time (ms) until tolerance...	1e-3	1e-6	1e-9	1e-12	1e-15	Per Iter.
DRGrAv (This Paper)	35.4	47.4	47.4	56.7	66.1	11.8
DeEPCA Ye & Zhang (2021)	35.8	45.7	53.0	61.1	80.6	9.13
DPRGD Deng & Hu (2023)	1860	61200	>100000	>100000	>100000	9.24
DPRGT Deng & Hu (2023)	2270	2910	3470	4200	4780	14.5
COM* Sarlette & Sepulchre (2009)	5050	7290	9260	11000	13200	16.3
Gossip* Mishra et al. (2019)	1280	1730	2150	2590	3390	427

In order to have a fair comparison, all average-consensus-based algorithms mentioned above used the same consensus protocol. Both graphs used the optimal Laplacian-based communication matrix (i.e. $\mathbf{W} = \mathbf{I} - \frac{1}{\gamma}\mathbf{L}$ for the hypercube graph, $\mathbf{W} \approx \mathbf{I} - \frac{1}{2}\mathbf{L}$ for the cycle graph, where \mathbf{L} is the

²Technically, Mishra et al. (2019) has a third hyperparameter, denoted ρ . However, their algorithm only ever uses the quantity ρa , so w.l.o.g. we let $\rho = 1$ and control only a .

432 corresponding graph Laplacian matrix) for 10 rounds of communication in the hypercube graph
 433 case and 50 rounds of communication in the cycle graph case.

434
 435 A single synthetic dataset $\{U_m\}_{m=1}^M$ of “normally distributed points with standard deviation $\frac{\pi}{4}$ ” was
 436 used for all experiments. In precise terms, said dataset was generated by sampling a center point U_C
 437 uniformly at random on $\text{Gr}(N, K)$ and then computing $U_m := \exp_{U_C}(T_m)$ for all $m \in [M]$, where
 438 $T_m := \tilde{U}_m \tilde{\Sigma}_m \tilde{V}_m^T$ is a random tangent vector at U_C such that \tilde{U}_m, \tilde{V}_m are sampled uniformly at
 439 random from sets $\{\tilde{U} \mid \tilde{U} \in \text{St}(N, K), U_C^T \tilde{U} = \mathbf{0}\}, \text{St}(K, K)$ respectively and $\tilde{\Sigma} := \text{Diag}(\frac{\pi}{4} z)$
 440 where $z \sim \mathcal{N}(\mathbf{0}_K, \mathbf{I}_K)$ is a vector draw of K i.i.d. standard normal random variables.
 441

442 The *Mean Squared Error* quantity at time t was computed as $\frac{1}{M} \sum_{m=1}^M d^2\left([\bar{U}_m^{(t)}], [\bar{U}]\right)$ where
 443 d is the extrinsic (or chordal) distance on the Grassmannian defined as $d([\bar{U}_1], [\bar{U}_2]) :=$
 444 $2^{-1/2} \left\| U_1 U_1^T - U_2 U_2^T \right\|_F$ and \bar{U} is the true IAM average (see Section 3), computed using the
 445 `torch.linalg.eigh` function on \bar{P} directly (runtime of 161 milliseconds). Similarly, the *Mean*
 446 *Squared Disagreement* quantity represents the amount to which the agents’ estimates vary at time t
 447 and was computed as $\frac{2}{M(M-1)} \sum_{m=1}^M \sum_{n=m+1}^M d^2\left([\bar{U}_m^{(t)}], [\bar{U}_n^{(t)}]\right)$.
 448
 449

450 The results of these experiments are shown in Figure 2 and Table 1. Six iterations were chosen
 451 for the hypercube graph case because at this point DRGrAv reaches floating point tolerance (FP-
 452 tol), demonstrating that such precision is achievable by an algorithm in such time; Nine iterations
 453 were chosen for the cycle graph in order to demonstrate the effective consensus-permitted tolerance
 454 (ECptol) of around 10^{-10} . DRGrAv performs the best out of all algorithms, converging in the hy-
 455 percube graph case to FPtol in only 6 iterations and converging in the cycle graph case to ECptol
 456 in only 5 iterations. The adapted DeEPCA method performs most closely to DRGrAv, however still
 457 lags behind several orders of magnitude in MSE. In the cycle graph, DeEPCA manages to barely
 458 beat DRGrAv in at the end, however given both are more or less at ECptol we do not think this
 459 provides strong evidence to prefer DeEPCA to DRGrAv. Since both COM and Gossip begin with
 460 $\bar{U}_m^{(0)} \leftarrow U_m$ instead of some pre-agreed upon starting $U^{(0)}$, they are able to have superior perfor-
 461 mance to DRGrAv in the short term; however after only 2 iterations this short term behavior ends.
 462 DPRGD and DPRGT, being generic algorithms for use on any compact submanifold, do not lever-
 463 age any of the structure specific to the Grassmannian problem and consequently are not empirically
 464 competitive to algorithms which do, e.g. DRGrAv. All algorithms presented have their specific ideal
 465 use cases, and we claim that the problem of decentralized Grassmannian averaging is the ideal use
 466 case of DRGrAv.

467 5.2 K-MEANS FOR VIDEO MOTION CLUSTERING

468
 469 We consider the application of the RGrAv algorithm to the problem of video motion analysis, extend-
 470 ing the work of Marrinan et al. (2014). Their study applied centralized subspace averaging methods
 471 to multiple tasks on the DARPA Mind’s Eye video dataset. In our work, we focus specifically on the
 472 (centralized) task of K -means clustering and compare against the best algorithm presented in their
 473 work.

474 The Mind’s Eye dataset consists of a set of “tracklets” — short grayscale videos sequences of moving
 475 objects, primarily people. Each tracklet consists of 48 frames of size 32×32 pixels. To prepare
 476 these tracklets for subspace analysis, they are flattened into matrices $X_t \in \mathbb{R}^{1024 \times 48}$, where each
 477 column represents a vectorized frame. The subspaces $U_t = \text{span}(X_t)$ spanned by these columns
 478 are treated as points on the Grassmannian, effectively encoding the essential motion patterns in the
 479 video.

480 Each tracklet is annotated with a label describing the type of motion it contains, such as “walk” or
 481 “ride-bike.” These labels provide ground truth for evaluating the effectiveness of clustering algo-
 482 rithms; clusters are considered high-quality if most tracklets in the cluster share the same label. In
 483 Marrinan et al. (2014), the authors find that the choice of averaging algorithm does not significantly
 484 affect cluster quality as the number of clusters increases, whereas runtime can differ significantly.
 485 As a result, we compare to the fastest averaging algorithm from their work – the flag mean – and
 show that applying RGrAv can reduce runtime for the clustering task.

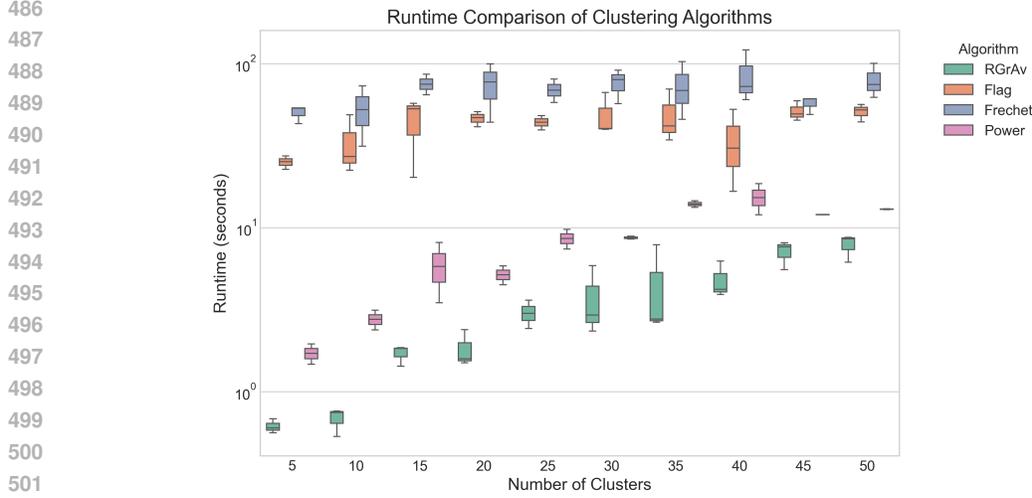


Figure 3: A comparison of runtime for K -means with various averaging algorithms and numbers of clusters K . The four colors represent the averaging algorithm as RGrAv (green), flag mean (orange), Fréchet mean (blue), and power method (pink). The four algorithms produce clusters with similar quality (excluded for brevity), but the RGrAv algorithm is significantly faster, showing $2 \times 10 \times$ speedup over the other averaging algorithms.

The standard K -means algorithm can be extended to cluster points on the Grassmannian by defining two primitives: a distance metric and an averaging operation. The standard K -means algorithm with these operations is shown in Algorithm 3. The centers \bar{U}_c are initialized randomly. At each iteration, the points U_t in the dataset are each assigned to their closest mean using the metric to form clusters. The means are then updated to the average of their respective clusters, and the steps are repeated until the means converge.

Algorithm 3 Grassmannian K -Means

Input: Subspaces $\{U_t\}_{t=1}^T$; Averaging algorithm: $\text{Ave}(U_1, \dots, U_n)$; Metric: $d(U_1, U_2)$

Output: Means $\{U_c\}_{c=1}^C$

```

 $\{\bar{U}_c^{(0)}\}_{c=1}^C \leftarrow \{\text{rand}(\text{St}(N, K))\}_{c=1}^C$  ▷ Initialize
while not converged do
   $\{i_t\}_{t=1}^T = \{i : d(U_t, \bar{U}_i^{(k)}) \leq d(U_t, \bar{U}_j^{(k)}) \quad \forall j\}_{t=1}^T$  ▷ Assign clusters
   $\bar{U}_c^{(k+1)} = \text{Ave}(\{U_t : i_t = c\})$  ▷ Compute new means
  converged =  $\max_c d(\bar{U}_c^{(k)}, \bar{U}_c^{(k+1)}) < \text{tol}$  ▷ Check termination
  k = k + 1
end while

```

For our clustering experiments, we test on the first 200 tracklets in the Mind’s eye dataset, which have a total of 24 unique labels. Given the success of DeEPCA (Ye & Zhang, 2021) in our benchmarks as a runner-up, we use the centralized version (the block power method) in these experiments as well. We test the performance of K -means with four different averaging algorithms, namely RGrAv, the power method, the Fréchet mean, and the flag mean. The distance operation for cluster assignment is chosen to be the chordal distance for computational efficiency. The Fréchet mean is computed via iterated gradient descent on the sum of squared distances cost function. Similar to Marrinan et al. (2014), we find that the four averaging algorithms produce clusters with similar quality across various values for K (these results are not shown for brevity). However as can be seen in Figure 3, the runtime varies significantly between algorithms. The Fréchet mean has the slowest runtime regardless of number of clusters, while RGrAv offers a $2 \times 10 \times$ speedup over the other averaging algorithms for all K values.

6 REPRODUCIBILITY STATEMENT

Section 5.1 is intended to be a comprehensive description sufficient for reproducibility; however, in addition all experiments from this section may be reproduced by running the `scripts/decentralized_grassmannian_averaging.py` script in the supplemental code; precise instructions and data formatting are described in the header of this file. The K -means experiment can be reproduced in three steps. First is by downloading the SUMMET dataset and putting it in the right subdirectory. Simply follow instructions in `data_sources/video_separation.py`. Next to run the experiment, return to the base directory and run the command `python -m scripts.tracklet_clustering` to run it as a module. Finally once this completes, there should be a `.pkl` file with the results. To create the visualization seen in this paper, run `python -m vis.visualize_tracklet` and look in the new `plots/` subdirectory for the results.

REFERENCES

- P-A Absil, Robert Mahony, and Rodolphe Sepulchre. Riemannian geometry of grassmann manifolds with a view on algorithmic computation. *Acta Applicandae Mathematica*, 80:199–220, 2004.
- Marc Arnaudon, Frédéric Barbaresco, and Le Yang. Medians and means in riemannian geometry: existence, uniqueness and computation. In *Matrix Information Geometry*, pp. 169–197. Springer, 2012.
- Enrique Tomás Martínez Beltrán, Mario Quiles Pérez, Pedro Miguel Sánchez Sánchez, Sergio López Bernal, Gérôme Bovet, Manuel Gil Pérez, Gregorio Martínez Pérez, and Alberto Huertas Celdrán. Decentralized federated learning: Fundamentals, state of the art, frameworks, trends, and challenges. *IEEE Communications Surveys & Tutorials*, 2023.
- Samuel R Buss and Jay P Fillmore. Spherical averages and applications to spherical splines and interpolation. *ACM Transactions on Graphics (TOG)*, 20(2):95–126, 2001.
- Rudrasis Chakraborty and Baba C. Vemuri. Recursive Fréchet Mean Computation on the Grassmannian and Its Applications to Computer Vision. In *2015 IEEE International Conference on Computer Vision (ICCV)*, pp. 4229–4237, Santiago, Chile, December 2015. IEEE. ISBN 978-1-4673-8391-2. doi: 10.1109/ICCV.2015.481. URL <http://ieeexplore.ieee.org/document/7410838/>.
- Rudrasis Chakraborty, Liu Yang, Søren Hauberg, and Baba C Vemuri. Intrinsic grassmann averages for online linear, robust and nonlinear subspace learning. *IEEE transactions on pattern analysis and machine intelligence*, 43(11):3904–3917, 2020.
- Jun Chen, Haishan Ye, Mengmeng Wang, Tianxin Huang, Guang Dai, Ivor W Tsang, and Yong Liu. Decentralized riemannian conjugate gradient method on the stiefel manifold. *arXiv preprint arXiv:2308.10547*, 2023.
- Shixiang Chen, Alfredo Garcia, Mingyi Hong, and Shahin Shahrampour. Decentralized riemannian gradient descent on the stiefel manifold. In *International Conference on Machine Learning*, pp. 1594–1605. PMLR, 2021.
- Guang Cheng, Jeffrey Ho, Hesamoddin Salehian, and Baba C Vemuri. Recursive computation of the fréchet mean on non-positively curved riemannian manifolds with applications. In *Riemannian Computing in Computer Vision*, pp. 21–43. Springer, 2016.
- Coleman DeLude, Santhosh Karnik, Mark Davenport, and Justin Romberg. Broadband beamforming via linear embedding, 2022. URL <https://arxiv.org/abs/2206.07143>.
- Kangkang Deng and Jiang Hu. Decentralized projected riemannian gradient method for smooth optimization on compact submanifolds. *arXiv preprint arXiv:2304.08241*, 2023.
- Thomas D Downs. Orientation statistics. *Biometrika*, 59(3):665–676, 1972.
- Alan Edelman, Tomás A Arias, and Steven T Smith. The geometry of algorithms with orthogonality constraints. *SIAM journal on Matrix Analysis and Applications*, 20(2):303–353, 1998.

- 594 Eugen Feller, C. Morin, and A. Esnault. A case for fully decentralized dynamic vm consolidation
595 in clouds. *4th IEEE International Conference on Cloud Computing Technology and Science*
596 *Proceedings*, pp. 26–33, 2012. doi: 10.1109/CloudCom.2012.6427585.
- 597
598 Simone Fiori, Tetsuya Kaneko, and Toshihisa Tanaka. Tangent-bundle maps on the grassmann man-
599 ifold: Application to empirical arithmetic averaging. *IEEE Transactions on Signal Processing*,
600 63(1):155–168, 2014.
- 601 P Thomas Fletcher, Suresh Venkatasubramanian, and Sarang Joshi. Robust statistics on riemannian
602 manifolds via the geometric median. In *2008 IEEE Conference on Computer Vision and Pattern*
603 *Recognition*, pp. 1–8. IEEE, 2008.
- 604
605 P Thomas Fletcher, Suresh Venkatasubramanian, and Sarang Joshi. The geometric median on rie-
606 mannian manifolds with application to robust atlas estimation. *NeuroImage*, 45(1):S143–S152,
607 2009.
- 608 Maurice Fréchet. Les éléments aléatoires de nature quelconque dans un espace distancié. In *Annales*
609 *de l’institut Henri Poincaré*, volume 10, pp. 215–310, 1948.
- 610
611 David Froelicher, Hyunghoon Cho, Manaswitha Edupalli, Joao Sa Sousa, Jean-Philippe Bossuat,
612 Apostolos Pyrgelis, Juan R Troncoso-Pastoriza, Bonnie Berger, and Jean-Pierre Hubaux. Scalable
613 and privacy-preserving federated principal component analysis. In *2023 IEEE Symposium on*
614 *Security and Privacy (SP)*, pp. 1908–1925. IEEE, 2023.
- 615 Kyle A Gallivan, Anuj Srivastava, Xiuwen Liu, and Paul Van Dooren. Efficient algorithms for
616 inferences on grassmann manifolds. In *IEEE Workshop on Statistical Signal Processing, 2003*,
617 pp. 315–318. IEEE, 2003.
- 618
619 GA Galperin. A concept of the mass center of a system of material points in the constant curvature
620 spaces. *Communications in Mathematical Physics*, 154:63–84, 1993.
- 621 Arpita Gang and Waheed U Bajwa. A linearly convergent algorithm for distributed principal com-
622 ponent analysis. *Signal Processing*, 193:108408, 2022.
- 623
624 Arpita Gang, Bingqing Xiang, and Waheed U Bajwa. Distributed principal subspace analysis for
625 partitioned big data: Algorithms, analysis, and implementation. *IEEE Transactions on Signal and*
626 *Information Processing over Networks*, 7:699–715, 2021.
- 627 Roger Ghanem and Debraj Ghosh. Efficient characterization of the random eigenvalue problem in a
628 polynomial chaos decomposition. *International Journal for Numerical Methods in Engineering*,
629 72(4):486–504, 2007.
- 630
631 Gene H. Golub and Charles F. Van Loan. *Matrix computations*. Johns Hopkins studies in the
632 mathematical sciences. The Johns Hopkins University Press, Baltimore, fourth edition edition,
633 2013. ISBN 978-1-4214-0794-4.
- 634 Peter Gruber and Fabian J Theis. Grassmann clustering. In *2006 14th European Signal Processing*
635 *Conference*, pp. 1–5. IEEE, 2006.
- 636
637 Shuo Han, Ufuk Topcu, and George J. Pappas. Differentially private distributed constrained opti-
638 mization. *IEEE Transactions on Automatic Control*, 62(1):50–64, 2017. doi: 10.1109/TAC.2016.
639 2541298.
- 640 Mehrtash Harandi, Conrad Sanderson, Chunhua Shen, and Brian C Lovell. Dictionary learning
641 and sparse coding on grassmann manifolds: An extrinsic solution. In *Proceedings of the IEEE*
642 *international conference on computer vision*, pp. 3120–3127, 2013.
- 643
644 Soren Hauberg, Aasa Feragen, and Michael J. Black. Grassmann averages for scalable robust pca.
645 In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*,
646 June 2014.
- 647 Zhiwu Huang, Jiqing Wu, and Luc Van Gool. Building deep networks on grassmann manifolds. In
Proceedings of the AAAI Conference on Artificial Intelligence, volume 32, 2018.

- 648 K Hueper and J Manton. The karcher mean of points on $so(n)$. *Talk at Cesame (UCL, Belgium)*,
649 2004.
- 650
- 651 Ben Jeuris, Raf Vandebril, and Bart Vandereycken. A survey and comparison of contemporary
652 algorithms for computing the matrix geometric mean. *Electronic Transactions on Numerical*
653 *Analysis*, 39:379–402, 2012.
- 654
- 655 Ian T Jolliffe and Jorge Cadima. Principal component analysis: a review and recent developments.
656 *Philosophical transactions of the royal society A: Mathematical, Physical and Engineering Sci-*
657 *ences*, 374(2016):20150202, 2016.
- 658 Raghunandan H Keshavan, Andrea Montanari, and Sewoong Oh. Matrix completion from a few
659 entries. *IEEE transactions on information theory*, 56(6):2980–2998, 2010.
- 660
- 661 Andrew V. Knyazev. Toward the optimal preconditioned eigensolver: Locally optimal block
662 preconditioned conjugate gradient method. *SIAM Journal on Scientific Computing*, 23(2):
663 517–541, 2001. doi: 10.1137/S1064827500366124. URL [https://doi.org/10.1137/
664 S1064827500366124](https://doi.org/10.1137/S1064827500366124).
- 665 Cornelius Lanczos. An iteration method for the solution of the eigenvalue problem of linear differ-
666 ential and integral operators. 1950.
- 667
- 668 Huiling Le. Locating fréchet means with application to shape spaces. *Advances in Applied Probab-*
669 *ility*, 33(2):324–338, 2001.
- 670
- 671 Jongmin Lee and Sungkyu Jung. Huber means on riemannian manifolds. *arXiv preprint*
672 *arXiv:2407.15764*, 2024.
- 673 Jörg Liesen and Zdenek Strakos. *Krylov subspace methods: principles and analysis*. Numerical
674 Mathematics and Scie, 2013.
- 675
- 676 David J Love, Robert W Heath, and Thomas Strohmer. Grassmannian beamforming for multiple-
677 input multiple-output wireless systems. *IEEE transactions on information theory*, 49(10):2735–
678 2747, 2003.
- 679
- 680 Yui Man Lui and J Ross Beveridge. Grassmann registration manifolds for face recognition. In
681 *European conference on computer vision*, pp. 44–57. Springer, 2008.
- 682
- 683 Tim Marrinan, Bruce Draper, J. Ross Beveridge, Michael Kirby, and Chris Peterson. Finding the
684 Subspace Mean or Median to Fit Your Need. In *2014 IEEE Conference on Computer Vision*
685 *and Pattern Recognition*, pp. 1082–1089, Columbus, OH, USA, June 2014. IEEE. ISBN 978-
686 1-4799-5118-5. doi: 10.1109/CVPR.2014.142. URL [https://ieeexplore.ieee.org/
687 document/6909538](https://ieeexplore.ieee.org/document/6909538).
- 688
- 689 Per-Gunnar Martinsson and Joel A Tropp. Randomized numerical linear algebra: Foundations and
690 algorithms. *Acta Numerica*, 29:403–572, 2020.
- 691
- 692 Gilles Meyer, Michel Journée, Silvere Bonnabel, and Rodolphe Sepulchre. From subspace learning
693 to distance learning: a geometrical optimization approach. In *2009 IEEE/SP 15th Workshop on*
694 *Statistical Signal Processing*, pp. 385–388. IEEE, 2009.
- 695
- 696 Bamdev Mishra, Hiroyuki Kasai, Pratik Jawanpuria, and Atul Saroop. A riemannian gossip ap-
697 proach to subspace learning on grassmann manifold. *Machine Learning*, 108:1783–1803, 2019.
- 698
- 699 Maher Moakher. Means and averaging in the group of rotations. *SIAM journal on matrix analysis*
700 *and applications*, 24(1):1–16, 2002.
- 701
- 702 Bishwarup Mondal, Satyaki Dutta, and Robert W Heath. Quantization on the grassmann manifold.
703 *IEEE Transactions on Signal Processing*, 55(8):4208–4216, 2007.
- 704
- 705 Angelia Nedic and Asuman Ozdaglar. Distributed subgradient methods for multi-agent optimiza-
706 tion. *IEEE Transactions on Automatic Control*, 54(1):48–61, 2009.

- 702 Angelia Nedić, Alex Olshevsky, and Michael G. Rabbat. Network topology and communication-
703 computation tradeoffs in decentralized optimization. *Proceedings of the IEEE*, 106(5):953–976,
704 2018. doi: 10.1109/JPROC.2018.2817461.
- 705 Guannan Qu and Na Li. Harnessing smoothness to accelerate distributed optimization. *IEEE Trans-*
706 *actions on Control of Network Systems*, 5(3):1245–1260, 2017.
- 708 Alain Sarlette and Rodolphe Sepulchre. Consensus optimization on manifolds. *SIAM journal on*
709 *Control and Optimization*, 48(1):56–76, 2009.
- 710 Rodolphe Sepulchre. Consensus on nonlinear spaces. *Annual reviews in control*, 35(1):56–64, 2011.
- 712 Wei Shi, Qing Ling, Gang Wu, and Wotao Yin. Extra: An exact first-order algorithm for decentral-
713 ized consensus optimization. *SIAM Journal on Optimization*, 25(2):944–966, 2015.
- 714 Rim Slama, Hazem Wannous, Mohamed Daoudi, and Anuj Srivastava. Accurate 3d action recogni-
715 tion using learning on the grassmann manifold. *Pattern Recognition*, 48(2):556–567, 2015.
- 717 Gerard LG Sleijpen and Henk A Van der Vorst. A jacobi–davidson iteration method for linear
718 eigenvalue problems. *SIAM review*, 42(2):267–293, 2000.
- 719 Tao Sun, Dongsheng Li, and Bao Wang. Decentralized federated averaging. *IEEE Transactions*
720 *on Pattern Analysis and Machine Intelligence*, 45:4289–4301, 2021. doi: 10.1109/TPAMI.2022.
721 3196503.
- 723 Roberto Tron, Bijan Afsari, and René Vidal. Riemannian consensus for manifolds with bounded
724 curvature. *IEEE Transactions on Automatic Control*, 58(4):921–934, 2012.
- 725 Pavan Turaga, Ashok Veeraraghavan, Anuj Srivastava, and Rama Chellappa. Statistical computa-
726 tions on grassmann and stiefel manifolds for image and video-based recognition. *IEEE Transac-*
727 *tions on Pattern Analysis and Machine Intelligence*, 33(11):2273–2286, 2011.
- 729 Sissi Xiaoxiao Wu, Hoi-To Wai, Lin Li, and Anna Scaglione. A review of distributed algorithms for
730 principal component analysis. *Proceedings of the IEEE*, 106(8):1321–1340, 2018.
- 731 Jinming Xu, Shanying Zhu, Yeng Chai Soh, and Lihua Xie. Augmented distributed gradient meth-
732 ods for multi-agent optimization under uncoordinated constant stepsizes. In *2015 54th IEEE*
733 *Conference on Decision and Control (CDC)*, pp. 2055–2060. IEEE, 2015.
- 734 Weiyu Xu and Babak Hassibi. Compressed sensing over the grassmann manifold: A unified ana-
735 lytical framework. In *2008 46th Annual Allerton Conference on Communication, Control, and*
736 *Computing*, pp. 562–567. IEEE, 2008.
- 738 Haishan Ye and Tong Zhang. Deepca: Decentralized exact pca with linear convergence rate. *Journal*
739 *of Machine Learning Research*, 22(238):1–27, 2021.
- 740 Se-Young Yun. Noisy Power Method with Grassmann Average. In *2018 IEEE International*
741 *Conference on Big Data and Smart Computing (BigComp)*, pp. 709–712, Shanghai, January
742 2018. IEEE. ISBN 978-1-5386-3649-7. doi: 10.1109/BigComp.2018.00132. URL <https://ieeexplore.ieee.org/document/8367212/>.
- 743
744
745 Jiayao Zhang, Guangxu Zhu, Robert W Heath Jr, and Kaibin Huang. Grassmannian learning: Em-
746 bedding geometry awareness in shallow and deep learning. *arXiv preprint arXiv:1808.02229*,
747 2018.
- 748 Peng Zhang and Sumei Sun. Decentralized network anomaly detection via a riemannian cluster
749 approach. In *GLOBECOM 2017-2017 IEEE Global Communications Conference*, pp. 1–6. IEEE,
750 2017.
- 751 Yunkai Zhou and Yousef Saad. A chebyshev–davidson algorithm for large symmetric eigenprob-
752 lems. *SIAM Journal on Matrix Analysis and Applications*, 29(3):954–971, 2007.
- 753
754
755

A ADDITIONAL ALGORITHMS

Algorithm 4 Asymptotic RGrAv (Rapid Grassmannian Averaging)

Input: $\alpha \in [0, 1)$, $\bar{U}^{(0)}$, $\{U_m\}_{m=1}^M$
Output: $\bar{U}^{(T)}$

```

S =  $\mathbf{I}_K$ 
for  $t = 1, 2, 3, \dots$  do
   $A_m^{(t)} \leftarrow U_m U_m^T \bar{U}^{(t-1)}$ 
   $\hat{A}^{(t)} \leftarrow \frac{1}{M} \sum_{m=1}^M A_m^{(t)}$  ▷ Power Iteration
  if  $t = 1$  then
     $\hat{Z}^{(1)} \leftarrow \hat{A}^{(1)}$  ▷  $[\hat{Z}^{(1)}] = [\bar{P}\bar{U}^{(0)}]$ 
  else
     $a_t, b_t, c_t \leftarrow \text{ChebyshevCoefficients}(t, \alpha)$  ▷ See Algorithm 7
     $\hat{Z}^{(t)} \leftarrow a_t (\hat{A}^{(t)} + b_t \bar{U}^{(t-1)} + c_t \bar{U}^{(t-2)})$  ▷  $[\hat{Z}^{(t)}] = [\tilde{f}_t^*(\bar{P})\bar{U}^{(0)}]$ 
  end if
  if  $t$  is on the orthonormalization schedule then
     $\bar{U}^{(t)}, S \leftarrow \text{StableQR}(\hat{Z}^{(t)})$  ▷ (Numerical Stability)
  else
     $\bar{U}^{(t)} \leftarrow \hat{Z}^{(t)} S$  ▷ (Numerical Stability)
  end if
   $\bar{U}^{(t-1)} \leftarrow \bar{U}^{(t-1)} S$  ▷ (Numerical Stability)
end for

```

Algorithm 5 Finite RGrAv (Rapid Grassmannian Averaging)

Inputs: $\alpha \in [0, 1)$, $T \in \mathbb{N}$, $\bar{U}^{(0)}$, $\{U_m\}_{m=1}^M$
Output: $\bar{U}^{(T)}$

```

S  $\leftarrow \mathbf{I}_K$ 
for  $t = 1, 2, 3, \dots, T$  do
   $A_m^{(t)} \leftarrow U_m U_m^T \bar{U}^{(t-1)}$ 
   $\hat{A}^{(t)} \leftarrow \frac{1}{M} \sum_{m=1}^M A_m^{(t)}$  ▷ Power Iteration
   $r_t \leftarrow \text{ChebyshevRoot}(t, T, \alpha)$  ▷ See Algorithm 8
   $\hat{Z}^{(t)} \leftarrow \frac{1}{1-r_t} (\hat{A}^{(t)} - r_t \bar{U}^{(t-1)})$ 
  if  $t$  is on the orthonormalization schedule then
     $\bar{U}^{(t)}, S \leftarrow \text{StableQR}(\hat{Z}^{(t)})$  ▷ (Numerical Stability)
  else
     $\bar{U}^{(t)} \leftarrow \hat{Z}^{(t)} S$  ▷ (Numerical Stability)
  end if
end for

```

B AUXILIARY THEOREMS

Lemma 2. Suppose $f^* \in \mathcal{P}'_t$ is a solution to eq. (2). Then the f^* has no roots in $[\alpha, 1]$.

Proof. If f^* had a root in $[\beta, 1]$, then $\min_{\lambda \in [\beta, 1]} |f^*(\lambda)| = 0$ and the objective function is unbounded, so f^* cannot have any roots in $[\beta, 1]$. This also allows us to conclude f^* is not the 0 polynomial, which is used in what follows.

Now, we show that if f^* has a root in (α, β) , moving this root to α strictly decreases the

Algorithm 6 *Finite DRGrAv (Decentralized Rapid Grassmannian Averaging)*

Input: $\alpha \in [0, 1), T \in \mathbb{N}, \{\bar{\mathbf{U}}_m^{(0)}\}_{m=1}^M, \{\mathbf{U}_m\}_{m=1}^M$

Output: $\bar{\mathbf{U}}^{(T)}$

$\mathbf{S}_m \leftarrow \mathbf{I}_K$

for $t = 1, 2, 3, \dots, T$ **do**

$\mathbf{A}_m^{(t)} \leftarrow \mathbf{U}_m \mathbf{U}_m^\top \bar{\mathbf{U}}^{(t-1)}$ ▷ Local Power Iteration

$r_t \leftarrow \text{ChebyshevRoot}(t, T, \alpha)$ ▷ See Algorithm 8

$\mathbf{Y}_m^{(t)} \leftarrow \frac{1}{1-r_t} (\mathbf{A}_m^{(t)} - r_t \bar{\mathbf{U}}_m^{(t-1)})$

if $t = 1$ **then**

$\mathbf{Z}_m^{(1)} \leftarrow \mathbf{Y}_m^{(1)}$ ▷ Gradient Tracking

else

$\mathbf{Z}_m^{(t)} \leftarrow \hat{\mathbf{Z}}_m^{(t-1)} + \mathbf{Y}_m^{(t)} - \mathbf{Y}_m^{(t-1)}$ ▷ Gradient Tracking

end if

$\hat{\mathbf{Z}}_m^{(t)} = \text{AverageConsensus}(\mathbf{Z}_m^{(t)})$

if t is on the orthonormalization schedule **then**

$\bar{\mathbf{U}}_m^{(t)}, \mathbf{S}_m \leftarrow \text{StableQR}(\hat{\mathbf{Z}}_m^{(t)})$ ▷ (Numerical Stability)

else

$\bar{\mathbf{U}}_m^{(t)} \leftarrow \hat{\mathbf{Z}}_m^{(t)} \mathbf{S}_m$ ▷ (Numerical Stability)

end if

end for

Algorithm 7 *ChebyshevCoefficients*

Input: $t \geq 2, \alpha \in [0, 1)$

Output: a_t, b_t, c_t

for $s = t - 2, t - 1, t$ **do**

if $s = 0$ **then**

$g_0 \leftarrow 1$

else

$r_s \leftarrow \cos\left(\frac{\pi}{2s}\right)$

$z_s \leftarrow \frac{1+r_s}{2}$

$\tau_s \leftarrow T_s^\alpha(z_s - r_s)$ ▷ T_s is the s th-order Chebyshev polynomial of the first kind

$g_s \leftarrow \frac{\tau_s}{z_s^2}$

end if

end for

for $s = t - 1, t$ **do**

$a_s \leftarrow 2 \frac{g_s - 1}{z_s}$

$q_s \leftarrow -\frac{r_s}{z_s}$

end for

$b_t \leftarrow tq_t - (t - 1)q_{t-1}$

if $t = 2$ **then**

$c_t \leftarrow 0$

else

$c_t \leftarrow \frac{1}{4}a_{t-1} \left(2t(t-1)(q_t - q_{t-1})^2 - \frac{t}{z_t^2} + \frac{t-1}{z_{t-1}^2} \right)$

end if

Algorithm 8 *ChebyshevRoot*

Input: $t \in \mathbb{N}, T \in \mathbb{N}, \alpha \in [0, 1)$

Output: $r_{t,T}$

$$r_{t,T} \leftarrow \alpha \frac{\cos\left(\frac{\pi(t+1/2)}{T}\right) + \cos\left(\frac{\pi}{2T}\right)}{1 + \cos\left(\frac{\pi}{2T}\right)}$$

objective function's value showing that f^* could not have been a solution to the minization problem.

If f^* had a root in (α, β) , then we can write $f^*(\lambda) = (\lambda - \lambda_1)g(\lambda)$ for some $\lambda_1 \in (\alpha, \beta)$ and $g \in \mathcal{P}'_{t-1}$ where g is not the 0 polynomial. It follows that

$$\begin{aligned} \frac{\max_{\lambda \in [0, \alpha]} |f^*(\lambda)|}{\min_{\lambda \in [\beta, 1]} |f^*(\lambda)|} &= \frac{\max_{\lambda \in [0, \alpha]} |(\lambda - \lambda_1)g(\lambda)|}{\min_{\lambda \in [\beta, 1]} |(\lambda - \lambda_1)g(\lambda)|} \\ &= \frac{\sup_{\lambda \in [0, \alpha]} |(\lambda - \alpha)g(\lambda)| \left| \frac{\lambda - \lambda_1}{\lambda - \alpha} \right|}{\min_{\lambda \in [\beta, 1]} |(\lambda - \alpha)g(\lambda)| \left| \frac{\lambda - \lambda_1}{\lambda - \alpha} \right|} \\ &> \frac{\max_{\lambda \in [0, \alpha]} |(\lambda - \alpha)g(\lambda)|}{\min_{\lambda \in [\beta, 1]} |(\lambda - \alpha)g(\lambda)|} \end{aligned}$$

where the final line results from the fact that $\left| \frac{\lambda - \lambda_1}{\lambda - \alpha} \right| < 1$ for $\lambda \in [\beta, 1]$ and $\left| \frac{\lambda - \lambda_1}{\lambda - \alpha} \right| > 1$ for $\lambda \in [0, \alpha)$. Now we observe that the function $h(\lambda) := (\lambda - \alpha)g(\lambda) \in \mathcal{P}'_t$ achieves a strictly lower value for the objective function via moving the root $\lambda_1 \in (\alpha, \beta)$ to $\lambda = \alpha$. Thus, f^* would not be a solution to eq. (2) and therefore f^* cannot have roots in (α, β) .

Now we show that if f^* had a root at $\lambda = \alpha$, we could slightly move the root to some point to the left of α and decrease the objective function's value.

Suppose f^* has a root at $\lambda = \alpha$. We can write $f^*(\lambda) = (\lambda - \alpha)g(\lambda)$ for some $g \in \mathcal{P}'_{t-1}$, where g is not the 0 polynomial. Let $\delta > 0$ be such that

$$\max_{\lambda \in [\alpha - \frac{\delta}{2}, \alpha]} |f^*(\lambda)| < \min_{\epsilon \in [0, \alpha]} \max_{\lambda \in [0, \alpha]} |(\lambda - \epsilon)g(\lambda)|$$

which exists since $\min_{\epsilon \in [0, \alpha]} \max_{\lambda \in [0, \alpha]} |(\lambda - \epsilon)g(\lambda)| > 0$ (as $g \neq 0$ from the beginning of the proof) and f^* is continuous. It follows that

$$\begin{aligned} \frac{\max_{\lambda \in [0, \alpha]} |f^*(\lambda)|}{\min_{\lambda \in [\beta, 1]} |f^*(\lambda)|} &= \frac{\max_{\lambda \in [0, \alpha]} |(\lambda - \alpha)g(\lambda)|}{\min_{\lambda \in [\beta, 1]} |(\lambda - \alpha)g(\lambda)|} \\ &= \frac{\sup_{\lambda \in [0, \alpha] \setminus \{\alpha - \delta\}} |(\lambda - (\alpha - \delta))g(\lambda)| \left| \frac{\lambda - \alpha}{\lambda - (\alpha - \delta)} \right|}{\min_{\lambda \in [\beta, 1]} |(\lambda - (\alpha - \delta))g(\lambda)| \left| \frac{\lambda - \alpha}{\lambda - (\alpha - \delta)} \right|} \\ &= \frac{\max\{\sup_{\lambda \in [0, \alpha - \frac{\delta}{2}] \setminus \{\alpha - \delta\}} |(\lambda - (\alpha - \delta))g(\lambda)| \left| \frac{\lambda - \alpha}{\lambda - (\alpha - \delta)} \right|, \Phi\}}{\min_{\lambda \in [\beta, 1]} |(\lambda - (\alpha - \delta))g(\lambda)| \left| \frac{\lambda - \alpha}{\lambda - (\alpha - \delta)} \right|} \\ &> \frac{\max_{\lambda \in [0, \alpha]} |(\lambda - (\alpha - \delta))g(\lambda)|}{\min_{\lambda \in [\beta, 1]} |(\lambda - (\alpha - \delta))g(\lambda)|} \end{aligned}$$

with $\Phi = \max_{\lambda \in [\alpha - \frac{\delta}{2}, \alpha]} |(\lambda - (\alpha - \delta))g(\lambda)| \left| \frac{\lambda - \alpha}{\lambda - (\alpha - \delta)} \right|$. The final line results from the fact that $\left| \frac{\lambda - \alpha}{\lambda - (\alpha - \delta)} \right| < 1$ for $\lambda \in [\beta, 1]$, $\left| \frac{\lambda - \alpha}{\lambda - (\alpha - \delta)} \right| \geq 1$ for $\lambda \in [0, \alpha - \frac{\delta}{2})$, and since $\left| \frac{\lambda - \alpha}{\lambda - (\alpha - \delta)} \right| \leq 1$ for $\lambda \in [\alpha - \frac{\delta}{2}, \alpha]$,

$$\begin{aligned} \max_{\lambda \in [\alpha - \frac{\delta}{2}, \alpha]} |(\lambda - (\alpha - \delta))g(\lambda)| \left| \frac{\lambda - \alpha}{\lambda - (\alpha - \delta)} \right| &< \min_{\epsilon \in [0, \alpha]} \max_{\lambda \in [0, \alpha]} |(\lambda - \epsilon)g(\lambda)| * 1 \\ &\leq \max_{\lambda \in [0, \alpha]} |(\lambda - (\alpha - \delta))g(\lambda)| \end{aligned}$$

Now we observe that the function $h(\lambda) := (\lambda - (\alpha - \delta))g(\lambda) \in \mathcal{P}'_t$ achieves a strictly lower value for the objective function via moving the root from α to $\lambda = \alpha - \delta$. Thus, f^* would not be a solution to eq. (2) and therefore f^* cannot have a root at $\lambda = \alpha$. \square

Definition 1. A polynomial f is t -**equioscillatory** on an interval $[a, b]$ if there exists some $a \leq \gamma_0 < \gamma_1 < \dots < \gamma_{t-1} \leq b$ such that $|f(\gamma_s)| = \max_{\lambda \in [a, b]} |f(\lambda)|$ for all $0 \leq s \leq t-1$ and $f(\gamma_0) = -f(\gamma_1) = f(\gamma_2) = -f(\gamma_3) = \dots$

Lemma 3. A solution to the minimization problem

$$\text{minimize}_{f_t \in \mathcal{P}_t} \frac{\max_{\lambda \in [0, \alpha]} |f_t(\lambda)|}{\min_{\lambda \in [\beta, 1]} |f_t(\lambda)|}$$

must be t -equioscillatory.

Proof. Suppose f^* is a minimizer in to eq. (2) that is not t -equioscillatory. First, we assume that without loss of generality, $f^*(\alpha) > 0$ since if $f^*(\alpha) < 0$, we could replace f^* with $-f^*$ and it would still be a minimizer to the eq. (2), and we cannot have $f^*(\alpha) = 0$ by Lemma 2.

Let $\{\lambda_1 = 0, \lambda_2, \dots, \lambda_m\} \in [0, \alpha]$ for some $m \leq t$ denote the distinct roots of f^* in increasing order that lie in $[0, \alpha]$. We have m -intervals, $\mathcal{I} = \{I_i\}_{i=1}^m$ where $I_i = [\lambda_i, \lambda_{i+1}]$ and $\lambda_{m+1} = \alpha$. We have that for any $i \in [m]$, $\text{sgn}(f^*(\lambda)) = c$ for every $\lambda \in \text{int}(I_i)$ where $c \in \{-1, 1\}$. In other words, in the interior of any intervals from \mathcal{I} , f^* takes strictly all positive values or strictly all negative values. Define, for any interval $I \subseteq \mathbb{R}$ and function g

$$\text{sgn}(g|_I) = c$$

where $c \in \{-1, 1\}$ is the value such that $\text{sgn}(f^*(\lambda)) = c$ for every $\lambda \in \text{int}(I)$. If such a value does not exist, $\text{sgn}(g|_I)$ is left undefined.

We also define for any compact $A \subseteq \mathbb{R}$ and $g \in C(\mathbb{R})$,

$$\|g\|_A := \max_{x \in A} |g(x)|$$

Let

$$L = \{I \in \mathcal{I} : \max_{\lambda \in I} |f^*(\lambda)| < \|f^*\|_{[0, \alpha]}\}$$

$$M = \{I \in \mathcal{I} : \max_{\lambda \in I} |f^*(\lambda)| = \|f^*\|_{[0, \alpha]}\} = \mathcal{I} \setminus L$$

$$J = \bigcup_{I \in L} I$$

Define $g(\lambda) = |f^*(\lambda)| - \|f^*\|_{[0, \alpha]}$ and denote $\epsilon = \min_{\lambda \in J} |g(\lambda)|$, i.e the minimum distance by which any point in J misses one of $\pm \|f^*\|_{[0, \alpha]}$.

Let $k = |\{\lambda \in [0, \alpha] : f^*(\lambda) = \|f^*\|_{[0, \alpha]}\}|$. Note that $k < t$ by assumption. Let $M = \{M_1, \dots, M_k\}$ be intervals listed from left to right where $i_j \in [m]$ are indices such that $M_j = I_{i_j} = [\lambda_{i_j}, \lambda_{i_j+1}]$ for $j \in [k]$ and $\lambda_{i_j} < \lambda_{i_m}$ for $1 \leq j < m \leq k$.

Now define

$$R = \{\lambda_{i_j} : \text{sgn}(f^*|_{M_j}) \neq \text{sgn}(f^*|_{M_{j-1}}) \text{ for some } 2 \leq j \leq k\}$$

We have that since $k \leq t-1$ by assumption, then $|R| \leq k-1 \leq t-2$.

With $R = \{r_1, \dots, r_q\}$ for some $q \leq t-2$, we define a polynomial $r : \mathbb{R} \rightarrow \mathbb{R}$ based on the sign of f^* on M_k .

Case A: If $\text{sgn}(f^*|_{M_k}) = 1$, we define

$$r(\lambda) := c_r \lambda (\lambda - \beta) \prod_{i=1}^q (\lambda - r_i)$$

Case B: Otherwise (when $\text{sgn}(f^*|_{M_k}) = -1$), define

$$r(\lambda) := c_r \lambda (\lambda - 1) \prod_{i=1}^q (\lambda - r_i)$$

In either case, r be a polynomial of degree at most t . We now select $c_r \in \mathbb{R}$ so that $\|r(\lambda)\|_{[0, \alpha]} < \epsilon$ and set the $\text{sgn}(c_r)$ so that $\text{sgn}(r|_{M_1}) = -\text{sgn}(f^*|_{M_1})$. We now show $\text{sgn}(r|_{M_j}) = -\text{sgn}(f^*|_{M_j})$ for every $j \in [k]$ via induction. Our base case holds via how we set c_r above.

Suppose inductively that $\text{sgn}(r|_{M_{j-1}}) = -\text{sgn}(f^*|_{M_{j-1}})$ for some $j \leq k$. Then, if $\text{sgn}(f^*|_{M_j}) = \text{sgn}(f^*|_{M_{j-1}})$, we have that as no root was added to r between these intervals and so

$$\text{sgn}(r|_{M_j}) = \text{sgn}(r|_{M_{j-1}}) = -\text{sgn}(f^*|_{M_{j-1}})$$

with the last equality by the induction hypothesis. Otherwise, if $\text{sgn}(f^*|_{M_j}) \neq \text{sgn}(f^*|_{M_{j-1}})$, then since a root is added between them

$$\text{sgn}(r|_{M_j}) = -\text{sgn}(r|_{M_{j-1}}) = \text{sgn}(f^*|_{M_{j-1}}) = -\text{sgn}(f^*|_{M_j})$$

So in both cases, we have $\text{sgn}(r|_{M_j}) = -\text{sgn}(f^*|_{M_j})$ which closes the induction.

Now we obtain that for any $j \in [k]$,

$$\|f^* + r\|_{M_j} < \|f^*\|_{[0,\alpha]}$$

as, for $M_j \neq [\lambda_m, \alpha]$, the maximum $\|f^*\|_{[0,\alpha]}$ is achieved by $|f^*|$ in the interior of M_j where $|r| > 0$ as r has no roots in this interior and $\text{sgn}(r|_{M_j}) = -\text{sgn}(f^*|_{M_j})$. Additionally, one recognizes that if $M_j = [\lambda_m, \alpha]$, then $r(\lambda) < 0$ for $x \in (\lambda_m, \alpha]$ and so the above bound still holds even when the maximum is attained on the boundary of the interval, i.e when $f^*(\alpha) = \|f^*\|_{[0,\alpha]}$. Furthermore, we have that

$$\begin{aligned} \|f^* + r\|_J &\leq \max_{\lambda \in J} |f^*(\lambda)| + \|r\|_{[0,\alpha]} \\ &< \max_{\lambda \in J} |f^*(\lambda)| + \epsilon = \|f^*\|_{[0,\alpha]} \end{aligned}$$

by ϵ 's definition. Now we have obtained

$$\|f^* + r\|_{[0,\alpha]} < \|f^*\|_{[0,\alpha]} \quad (3)$$

Now we turn to bounding the denominator of the objective function using f^* from above.

Note that $\text{sgn}(f^*|_{[\alpha,1]}) = 1$ as $f^*(\alpha) > 0$ and f^* has no roots in $[\alpha, 1]$ by Lemma 2. We will now show that in either case,

$$\min_{\lambda \in [\beta,1]} |(f^* + r)(\lambda)| \geq \min_{\lambda \in [\beta,1]} |f^*(\lambda)| \quad (4)$$

Case A: Since $\text{sgn}(f^*|_{M_k}) = 1$, $\text{sgn}(r|_{M_k}) = -1$ which gives $\text{sgn}(r|_{[\lambda_{i_k}, \beta]}) = -1$ as r has no roots in (λ_{i_k}, β) . Since r has a simple root at $\lambda = \beta$ with no other roots greater than $\lambda = \beta$, we have that $\text{sgn}(r|_{[\beta,1]}) = 1$. It immediately follows that

$$\min_{\lambda \in [\beta,1]} |(f^* + r)(\lambda)| \geq \min_{\lambda \in [\beta,1]} |f^*(\lambda)|$$

Case B: Since $\text{sgn}(f^*|_{M_k}) = -1$, $\text{sgn}(r|_{M_k}) = 1$ which gives $\text{sgn}(r|_{[\lambda_{i_k}, 1]}) = 1$ as r has no roots in $(\lambda_{i_k}, 1)$. So $\text{sgn}(r|_{[\beta,1]}) = 1$ and it follows that

$$\min_{\lambda \in [\beta,1]} |(f^* + r)(\lambda)| \geq \min_{\lambda \in [\beta,1]} |f^*(\lambda)|$$

Combining eq. (3) and eq. (4) gives

$$\frac{\|f^*\|_{[0,\alpha]}}{\min_{\lambda \in [\beta,1]} |f^*(\lambda)|} \geq \frac{\|f^* + r\|_{[0,\alpha]}}{\min_{\lambda \in [\beta,1]} |(f^* + r)(\lambda)|}$$

We note that $(f^* + r)(0) = 0$ and $f^* + r$ is an at most t -degree polynomial and therefore $f^* + r$ is a feasible function for the minimization problem eq. (2). Thus, f^* is not a solution to eq. (2) as $f^* + r$ is feasible and achieves a strictly smaller value for the objective function, which proves the lemma. \square

Lemma 4. Define $\mathcal{P}_t'' := \{f_t \in \mathcal{P}_t \mid f_t(0) = 0, f_t \text{ is } t\text{-equioscillatory on } [0, \alpha], f_t(\beta) = 1\}$ for $\beta > \alpha$. The problem

$$\text{minimize}_{f_t \in \mathcal{P}_t''} \max_{\lambda \in [0,\alpha]} |f_t(\lambda)|$$

is solved by

$$\begin{aligned} f_t^*(\lambda) &= \prod_{s=0}^{t-1} \frac{\lambda - r_{s,t}}{\beta - r_{s,t}} \\ r_{s,t} &:= \alpha \frac{\cos\left(\frac{\pi(s+1/2)}{t}\right) + \cos\left(\frac{\pi}{2t}\right)}{1 + \cos\left(\frac{\pi}{2t}\right)} \end{aligned}$$

1026 *Proof.* First note that f_t^* is indeed feasible; from inspection one realizes $f_t^*(\beta) = 1$ and $r_{t-1,t} = 0$
 1027 (consequently $f_t^*(0) = 0$) and as a variant of a Chebyshev polynomial (of the first kind), f_t^* is
 1028 t -equioscillatory on $[0, \alpha]$ with extremal points

$$1029 \gamma_s = \alpha \frac{\cos\left(\frac{\pi s}{t}\right) + \cos\left(\frac{\pi}{2t}\right)}{1 + \cos\left(\frac{\pi}{2t}\right)}, \quad 0 \leq s \leq t-1$$

1032 Assume for the sake of contradiction that there exists some $g \in \mathcal{P}_t''$ such that $\max_{\lambda \in [0, \alpha]} |g(\lambda)| <$
 1033 $\max_{\lambda \in [0, \alpha]} |f_t^*(\lambda)|$. This would then imply that the difference polynomial $h(\lambda) := f_t^*(\lambda) - g(\lambda)$
 1034 satisfies the following

$$1035 \forall s = 0, 2, \dots : \quad h(\gamma_s) = f_t^*(\gamma_s) - g(\gamma_s)$$

$$1036 \quad \quad \quad > f_t^*(\gamma_s) - \max_{\lambda \in [0, \alpha]} |f_t^*(\lambda)|$$

$$1037 \quad \quad \quad = 0$$

$$1038$$

$$1039 \forall s = 1, 3, \dots : \quad h(\gamma_s) = f_t^*(\gamma_s) - g(\gamma_s)$$

$$1040 \quad \quad \quad < f_t^*(\gamma_s) + \max_{\lambda \in [0, \alpha]} |f_t^*(\lambda)|$$

$$1041 \quad \quad \quad = 0$$

$$1042$$

$$1043 \quad \quad \quad h(0) = f_t^*(0) - g(0)$$

$$1044 \quad \quad \quad = 0$$

$$1045 \quad \quad \quad h(\beta) = f_t^*(\beta) - g(\beta)$$

$$1046 \quad \quad \quad = 0$$

1048 Since h changes sign on every γ_s , it must have at least one root in each of the $t-1$ intervals between
 1049 consecutive γ_s . By construction, h also has a 2 more roots at 0 and β , meaning h has at minimum
 1050 $t+1$ distinct roots. However, h is a t th order polynomial, leading to a contradiction. \square
 1051

1052 B.1 PROOF OF THEOREM 1

1053 *Proof.* First, by Lemma 3, since any minimizer f^* of eq. (2) is t -equioscillatory, all t of f^* 's roots
 1054 lie in $[0, \alpha]$. This implies that f^* is increasing in $[\beta, 1]$ and therefore we have

$$1056 \min_{\lambda \in [\beta, 1]} |f^*(\lambda)| = f^*(\beta)$$

1058 which is assumed to be positive without loss of generality also as in Lemma 3. Now we can scale
 1059 f^* so that $f(\beta) = 1$ without changing the value of the objective function, i.e

$$1060 \frac{\max_{\lambda \in [0, \alpha]} |f^*(\lambda)|}{f^*(\beta)} = \frac{\max_{\lambda \in [0, \alpha]} |C f^*(\lambda)|}{C f^*(\beta)} = \max_{\lambda \in [0, \alpha]} |C f^*(\lambda)|$$

1063 where $C = \frac{1}{f^*(\beta)}$. We have now transformed the problem into

$$1064 \text{minimize } \max_{f_t \in \mathcal{P}_t''} \max_{\lambda \in [0, \alpha]} |f_t(\lambda)|$$

1066 as in Lemma 4 which is solved by

$$1068 f_t^*(\lambda) = \prod_{s=0}^{t-1} \frac{\lambda - r_{s,t}}{\beta - r_{s,t}}$$

$$1069 \quad \quad \quad r_{s,t} := \alpha \frac{\cos\left(\frac{\pi(s+1/2)}{t}\right) + \cos\left(\frac{\pi}{2t}\right)}{1 + \cos\left(\frac{\pi}{2t}\right)}$$

1073 as desired.

1075 B.2 SIMPLIFICATION OF EQUATION (1)

1076 We begin by simplifying the original optimization problem

$$1077 [\bar{U}] = \underset{[U] \in \text{Gr}(N, K)}{\text{argmin}} \left\| \bar{P} - UU^T \right\|_F^2$$

$$\begin{aligned}
1080 & = \operatorname{argmin}_{[\mathbf{U}] \in \operatorname{Gr}(N, K)} -2 \langle \bar{\mathbf{P}}, \mathbf{U}\mathbf{U}^\top \rangle + \|\bar{\mathbf{P}}\|_{\text{F}}^2 + \|\mathbf{U}\mathbf{U}^\top\|_{\text{F}}^2 \\
1081 & = \operatorname{argmin}_{[\mathbf{U}] \in \operatorname{Gr}(N, K)} -\langle \bar{\mathbf{P}}, \mathbf{U}\mathbf{U}^\top \rangle \\
1082 & = \operatorname{argmin}_{[\mathbf{U}] \in \operatorname{Gr}(N, K)} -\langle \tilde{\mathbf{V}}\tilde{\mathbf{\Lambda}}\tilde{\mathbf{V}}^\top, \mathbf{U}\mathbf{U}^\top \rangle \\
1083 & = \operatorname{argmin}_{[\mathbf{U}] \in \operatorname{Gr}(N, K)} -\langle \tilde{\mathbf{\Lambda}}, \tilde{\mathbf{V}}^\top \mathbf{U}\mathbf{U}^\top \tilde{\mathbf{V}} \rangle \\
1084 & = \left[\tilde{\mathbf{V}} \operatorname{argmin}_{\mathbf{W} \in \operatorname{St}(N, K)} -\langle \tilde{\mathbf{\Lambda}}, \mathbf{W}\mathbf{W}^\top \rangle \right] \\
1085 & = \left[\begin{bmatrix} \mathbf{V} & \mathbf{V}_\perp \end{bmatrix} \operatorname{argmin}_{\mathbf{W} \in \operatorname{St}(N, K)} -\langle \begin{bmatrix} \mathbf{\Lambda} & \mathbf{0} \\ \mathbf{0} & \mathbf{\Lambda}_\perp \end{bmatrix}, \mathbf{W}\mathbf{W}^\top \rangle \right] \\
1086 & \\
1087 & \\
1088 & \\
1089 & \\
1090 & \\
1091 & \\
1092 & \\
1093 & \\
1094 & \\
1095 & \\
1096 & \\
1097 & \\
1098 & \\
1099 & \\
1100 & \\
1101 & \\
1102 & \\
1103 & \\
1104 & \\
1105 & \\
1106 & \\
1107 & \\
1108 & \\
1109 & \\
1110 & \\
1111 & \\
1112 & \\
1113 & \\
1114 & \\
1115 & \\
1116 & \\
1117 & \\
1118 & \\
1119 & \\
1120 & \\
1121 & \\
1122 & \\
1123 & \\
1124 & \\
1125 & \\
1126 & \\
1127 & \\
1128 & \\
1129 & \\
1130 & \\
1131 & \\
1132 & \\
1133 &
\end{aligned}$$

Since $\mathbf{W} \in \operatorname{St}(N, K)$, it follows that $\operatorname{tr}(\mathbf{W}\mathbf{W}^\top) = K$ and all diagonal elements of $\mathbf{W}\mathbf{W}^\top$ lay in the range $[-1, 1]$. Consequently, $\min_{\mathbf{W} \in \operatorname{St}(N, K)} -\langle \tilde{\mathbf{\Lambda}}, \mathbf{W}\mathbf{W}^\top \rangle \geq -\operatorname{tr}(\tilde{\mathbf{\Lambda}})$. Since $\mathbf{W} = \begin{bmatrix} \mathbf{Q} \\ \mathbf{0} \end{bmatrix}$ for arbitrary $\mathbf{Q} \in \operatorname{St}(K, K)$ satisfies this inequality with equality, it may be substituted as the solution and a final simplification completes the proof.

$$\begin{aligned}
1101 & \\
1102 & \\
1103 & \\
1104 & \\
1105 & \\
1106 & \\
1107 & \\
1108 & \\
1109 & \\
1110 & \\
1111 & \\
1112 & \\
1113 & \\
1114 & \\
1115 & \\
1116 & \\
1117 & \\
1118 & \\
1119 & \\
1120 & \\
1121 & \\
1122 & \\
1123 & \\
1124 & \\
1125 & \\
1126 & \\
1127 & \\
1128 & \\
1129 & \\
1130 & \\
1131 & \\
1132 & \\
1133 &
\end{aligned}$$

$$\begin{aligned}
[\bar{\mathbf{U}}] & = \left[\begin{bmatrix} \mathbf{V} & \mathbf{V}_\perp \end{bmatrix} \operatorname{argmin}_{\mathbf{W} \in \operatorname{St}(N, K)} -\langle \begin{bmatrix} \mathbf{\Lambda} & \mathbf{0} \\ \mathbf{0} & \mathbf{\Lambda}_\perp \end{bmatrix}, \mathbf{W}\mathbf{W}^\top \rangle \right] \\
& = \left[\begin{bmatrix} \mathbf{V} & \mathbf{V}_\perp \end{bmatrix} \begin{bmatrix} \mathbf{Q} \\ \mathbf{0} \end{bmatrix} \right] \\
& = [\mathbf{V}\mathbf{Q}] \\
& = [\mathbf{V}]
\end{aligned}$$