

42 collection of manifolds, the MEB must be computed on one individual Grassmannian and the
 43 choice of which is not obvious. Determining which Grassmannian provides the best center for
 44 a collection of subspaces is one of the tasks of this manuscript, and we provide a geometrically
 45 motivated criteria for automatically selecting this manifold.

46 With subspace data, it is natural to think of the center of the Grassmannian minimum
 47 enclosing ball (GMEB) as the common information in the data set. Common subspace ex-
 48 traction can be found in subspace clustering [1], domain adaptation, and subspace alignment.
 49 These tools can be used in a plethora of tasks in pattern recognition including subspace
 50 tracking [32], face recognition [7, 8], video action recognition [7, 22], infected patient diag-
 51 nosis [18], adaptive sorting [15], model reduction [11], and many more. Common subspace
 52 extraction is frequently done by finding the ℓ_2 - or ℓ_1 -center in cases where outliers are present
 53 in the data collection, but if the data are drawn from a uniform distribution whose support
 54 is a ball, the ℓ_∞ -center gives the maximum likelihood estimator for the center of the support
 55 and thus may be preferred when all the subspaces have been drawn from a single uniform
 56 distribution [2]. Furthermore, techniques have been developed to prune outliers from data
 57 sets using the ℓ_∞ -norm, with theoretical guarantees in some circumstances [31].

58 In this paper, we present a novel technique to accurately estimate the GMEB for a
 59 collection of linear subspaces of possibly differing dimension, and a geometrically inspired
 60 order-selection rule to identify the Grassmannian that best represents the shared information
 61 in the data. Choosing the ideal manifold on which to perform the ℓ_∞ -averaging is inherently
 62 related to finding a common subspace of optimal rank, and thus the numerical experiments
 63 explore the relationships between different rank-adaptive subspace averaging methods.

64 The main contributions of the paper are summarized as follows. We propose

- 65 • a subgradient approach to solve the dual of the GMEB problem for subspaces of
 66 differing dimensions. A duality gap of zero certifies the solution as optimal.
- 67 • an unsupervised order-selection rule for the dimension of the center of the GMEB.
- 68 • a warm-start initialization for the subgradient algorithm that reduces the number of
 69 iterations needed for the subgradient algorithm to converge.
- 70 • a hybrid method for order-selection which modifies the existing rule of [28] for use
 71 with the center of the GMEB.
- 72 • a synthetic data model that allows us to measure the accuracy of an estimate for the
 73 center of the GMEB, and demonstrate the effectiveness of the proposed technique
 74 using data generated with this model.

75 Finally, we compare the proposed order-selection rules to existing methods for automatic order
 76 selection in subspace averaging with numerical experiments.

77 **2. Mathematical background: Grassmannian minimum enclosing ball.** In this sec-
 78 tion we provide the mathematical background necessary to formulate the GMEB problem for
 79 subspaces of differing dimension. We begin by stating the relevant properties of invariant
 80 metrics, a standard reference on this topic is [33]. We recall the maps defined in [38] that
 81 associate a subset of points on a single manifold with each subspace from the collection and
 82 the point-to-set distance that measures the dissimilarity of these sets. Finally, we explicitly
 83 state the minimax optimization problem that defines this GMEB.

84 Denote by $\text{Gr}(k, n)$ the Grassmann manifold of k -dimensional subspaces in \mathbb{R}^n . If A is
 85 an $n \times k$ matrix with full column rank, the column space of A , $\text{col}(A)$, defines a subspace
 86 that can be identified with a point $\mathbf{A} \in \text{Gr}(k, n)$. To simplify notation we assume without
 87 loss of generality that the chosen representative for a point $\mathbf{A} \in \text{Gr}(k, n)$ is an orthonormal
 88 basis, $A \in \mathbb{R}^{n \times k}$ with $A^T A = I$. Let $O(k)$ denote the set of $k \times k$ orthogonal matrices. If
 89 $Q_k \in O(k)$ then $\text{col}(AQ_k) = \text{col}(A) = \mathbf{A}$, and we can see that a point on this Grassmannian
 90 can be represented by any real $n \times k$ matrix that spans the same subspace. For any two points,

91 $\mathbf{A}, \mathbf{B} \in \text{Gr}(k, n)$, there exists a set of k principal angles, $0 \leq \theta_1(\mathbf{A}, \mathbf{B}) \leq \dots \leq \theta_k(\mathbf{A}, \mathbf{B}) \leq \pi/2$,
 92 defined recursively as

$$\begin{aligned} \theta_1(\mathbf{A}, \mathbf{B}) &\doteq \min_{\mathbf{a}_1 \in \mathbf{A}, \mathbf{b}_1 \in \mathbf{B}} \cos^{-1} \left(\frac{\mathbf{a}_1^T \mathbf{b}_1}{\|\mathbf{a}_1\|_2 \|\mathbf{b}_1\|_2} \right), \text{ and for } i = 2, \dots, k \\ 93 \quad (2.1) \quad \theta_i(\mathbf{A}, \mathbf{B}) &\doteq \min_{\mathbf{a}_i \in \mathbf{A}, \mathbf{b}_i \in \mathbf{B}} \cos^{-1} \left(\frac{\mathbf{a}_i^T \mathbf{b}_i}{\|\mathbf{a}_i\|_2 \|\mathbf{b}_i\|_2} \right) \\ &\text{s.t. } \mathbf{a}_j^T \mathbf{a}_i = 0 \text{ for } j < i \\ &\quad \mathbf{b}_j^T \mathbf{b}_i = 0 \text{ for } j < i. \end{aligned}$$

94 The vectors that form these angles, $\{\mathbf{a}_1, \dots, \mathbf{a}_k\}$ and $\{\mathbf{b}_1, \dots, \mathbf{b}_k\}$, are called the left and
 95 right principal vectors, respectively, and when normalized, these vectors form orthonormal
 96 bases $A, B \in \mathbb{R}^{n \times k}$, for the spaces \mathbf{A} and \mathbf{B} . The principal angles and principal vectors can be
 97 computed via the singular value decomposition (SVD) [6]. Let $A^T B = V \Sigma W^T$ be a thin SVD
 98 with the singular values sorted in nonincreasing order, so that

$$\begin{aligned} 99 \quad (2.2) \quad V &\in \mathbb{R}^{k \times k} \text{ with } V^T V = I, \\ \Sigma &\in \mathbb{R}^{k \times k} \text{ with } \Sigma = \text{diag}(\cos(\theta(\mathbf{A}, \mathbf{B}))), \text{ and} \\ W &\in \mathbb{R}^{k \times k} \text{ with } W^T W = I. \end{aligned}$$

100 Then $\theta_i(\mathbf{A}, \mathbf{B}) = \cos^{-1}(\Sigma_{ii})$ is the i th principal angle separating \mathbf{A} and \mathbf{B} , with associated left
 101 and right principal vectors $\mathbf{a}_i = A \mathbf{v}_i$ and $\mathbf{b}_i = B \mathbf{w}_i$ for $i = 1, \dots, k$.

102 Let $d : \text{Gr}(k, n) \times \text{Gr}(k, n) \rightarrow \mathbb{R}$ be a metric. If for all $\mathbf{A}, \mathbf{B} \in \text{Gr}(k, n)$ and for all
 103 $Q_n \in O(n)$ the left action of Q_n on A and B by multiplication does not change the value
 104 of the metric, that is, $d(\mathbf{A}, \mathbf{B}) = d(Q_n \mathbf{A}, Q_n \mathbf{B})$, then d is said to be orthogonally invariant.
 105 Orthogonally invariant metrics depend only on the relative position of \mathbf{A} and \mathbf{B} , so as a result
 106 of [37, Thm. 3], d can be written as a function of the vector of principal angles separating \mathbf{A} and
 107 \mathbf{B} , $\theta(\mathbf{A}, \mathbf{B}) \in \mathbb{R}^k$. Additionally, for $\text{Gr}(k, n)$ with either $k \neq 2$ or $n \neq 2$ there is an essentially
 108 unique invariant Riemannian metric (up to scaling) which yields $d(\mathbf{A}, \mathbf{B}) = \|\theta(\mathbf{A}, \mathbf{B})\|_2$, and
 109 is frequently referred to as the geodesic distance based on arc length [37].

110 Let $\mathcal{D} = \{\mathbf{X}_i\}_{i=1}^M$ be a finite collection of subspaces of \mathbb{R}^n with possibly different dimen-
 111 sions, so that $\dim(\mathbf{X}_i) = p_i$. For the set of positive integers $\mathcal{P} = \{\dim(\mathbf{X}_i) : \mathbf{X}_i \in \mathcal{D}\}$ we
 112 can consider \mathcal{D} as a collection of points lying on the disjoint union of Grassmann manifolds,
 113 $\mathbf{X}_i \in \coprod_{p \in \mathcal{P}} \text{Gr}(p, n)$. To account for the difference in subspace dimensions, we adopt the
 114 convention of [38] by redefining $d(\mathbf{U}, \mathbf{X}_i)$ as the minimum distance between \mathbf{U} and a subset
 115 of points on $\text{Gr}(k, n)$, appropriately defined for each $\mathbf{X}_i \in \mathcal{D}$. Each subspace is associated
 116 with one of two types of subset, which are defined by

$$\begin{aligned} 117 \quad (2.3) \quad \Omega_+(\mathbf{X}_i) &\doteq \{\mathbf{Y} \in \text{Gr}(k, n) : \mathbf{X}_i \subseteq \mathbf{Y}\} \text{ for } p_i < k, \text{ and} \\ \Omega_-(\mathbf{X}_i) &\doteq \{\mathbf{Y} \in \text{Gr}(k, n) : \mathbf{Y} \subseteq \mathbf{X}_i\} \text{ for } p_i \geq k. \end{aligned}$$

118 We use $\Omega(\mathbf{X}_i)$ when referring to either type generically. For \mathbf{X}_i such that $p_i < k$, $\Omega_+(\mathbf{X}_i)$
 119 is the set of all points of $\text{Gr}(k, n)$ containing \mathbf{X}_i . Alternatively when \mathbf{X}_i is a p_i -plane with
 120 $p_i > k$, $\Omega_-(\mathbf{X}_i)$ is all k -dimensional subspaces contained in \mathbf{X}_i , and when $p_i = k$ the subset
 121 of points is just the singleton, \mathbf{X}_i .

122 Finally, we overload the notation for distance so that

$$123 \quad (2.4) \quad d_{\text{Gr}(k, n)}(\mathbf{U}, \mathbf{X}_i) \doteq d_{\text{Gr}(k, n)}(\mathbf{U}, \Omega(\mathbf{X}_i)) = \min\{d(\mathbf{U}, \mathbf{Y}_i) : \mathbf{Y}_i \in \Omega(\mathbf{X}_i)\}$$

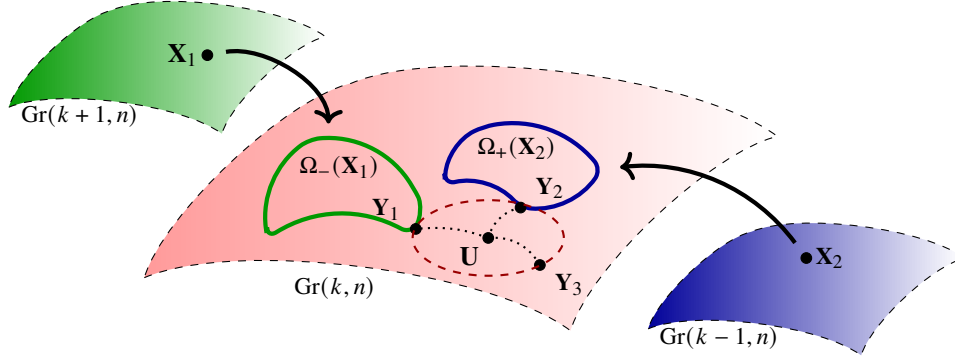


Fig. 1: Illustration of the minimum point-to-set distance on $\text{Gr}(k, n)$ between \mathbf{U} and the sets $\Omega_-(\mathbf{X}_1)$, $\Omega_+(\mathbf{X}_2)$, and \mathbf{Y}_3 , associated with points on $\text{Gr}(k+1, n)$, $\text{Gr}(k-1, n)$, and $\text{Gr}(k, n)$, respectively. The points that realize the minimum distance are $\mathbf{Y}_1 \in \Omega_-(\mathbf{X}_1)$, $\mathbf{Y}_2 \in \Omega_+(\mathbf{X}_2)$, and \mathbf{Y}_3 . The point \mathbf{U} is the center of the minimum enclosing ball of \mathbf{Y}_1 , \mathbf{Y}_2 , and \mathbf{Y}_3 .

124 when the distance is being measured on $\text{Gr}(k, n)$ and the data comes from Grassmann manifolds of possibly differing dimension. This is the proposed distance of [38], which is well-defined a fixed value of k . Figure 1 shows an illustration of this distance as the length of the shortest path between a point, \mathbf{U} , and the sets of points, $\Omega(\mathbf{X}_i)$ for $i = 1, \dots, 3$. In this particular case $\mathbf{Y}_3 \in \text{Gr}(k, n)$ so $\mathbf{Y}_3 = \mathbf{X}_3 = \Omega(\mathbf{X}_3)$.

129 The minimum in Equation (2.4) always exists because $\Omega(\mathbf{X}_i)$ is a closed subset of the Grassmannian, and the points satisfying $\mathbf{Y}_i \in \arg \min_{\mathbf{Y} \in \Omega(\mathbf{X}_i)} d(\mathbf{U}, \mathbf{Y})$ are independent of the choice of orthogonally invariant distance measure. Let $U^T \mathbf{X}_i = V \Sigma W^T$ be a thin SVD. One point that achieves the minimum distance is the columnspace of the matrix defined by

$$133 \quad (2.5) \quad \mathbf{Y}_i \doteq \begin{cases} [X_i \mathbf{w}_1, \dots, X_i \mathbf{w}_k] & \text{for } p_i \geq k; \\ [X_i \mathbf{w}_1, \dots, X_i \mathbf{w}_{p_i}, U \mathbf{v}_{p_i+1}, \dots, U \mathbf{v}_k] & \text{otherwise.} \end{cases}$$

134 This derivation can be found in, e.g. [29].

135 This formalism implies that distances can be written as a function of exactly k principal angles regardless of the dimension of \mathbf{X}_i , and conveniently the definition agrees with many pseudo-metrics commonly used in the literature that measure similarity as a function of the (possibly less than k) principal angles between subspaces of different dimension. It should be clear, however, that this is not a metric because the distance between \mathbf{A} and \mathbf{B} will be zero if \mathbf{A} is a proper subspace of \mathbf{B} , despite being non-identical.

141 This manuscript is concerned with computing the minimax center, i.e., the center of the GMEB, on $\text{Gr}(k, n)$ for the collection of subspaces, \mathcal{D} , using the point-to-set distance. However, rather than using a metric on $\text{Gr}(k, n)$ we measure dissimilarity by the squared chordal distance, $d(\mathbf{A}, \mathbf{B}) = \|\sin(\boldsymbol{\theta}(\mathbf{A}, \mathbf{B}))\|_2^2$. The minimum point-to-set distance using the squared chordal distance is

$$146 \quad (2.6) \quad \begin{aligned} d_{\text{Gr}(k, n)}(\mathbf{U}, \mathbf{X}_i) &= \|\sin(\boldsymbol{\theta}(\mathbf{U}, \mathbf{Y}_i))\|_2^2 \\ &= \frac{1}{2} \|U_k U_k^T - Y_i Y_i^T\|_F^2 \\ &= k - \text{Tr}(U^T Y_i Y_i^T U) \\ &= \min\{k, p_i\} - \text{Tr}(U^T X_i X_i^T U), \end{aligned}$$

147 where $\theta(\mathbf{U}, \mathbf{Y}_i) \in \mathbb{R}^k$ is the vector of principal angles between \mathbf{U} and a point $\mathbf{Y}_i \in \Omega(\mathbf{X}_i)$ that
 148 attains the minimum. The final equality in Equation (2.6) can be seen from the definition of
 149 \mathbf{Y}_i in Equation (2.5) and will be demonstrated in Equation (5.13). Note that it is not necessary
 150 to know \mathbf{Y}_i in order to compute $d_{\text{Gr}(k,n)}(\mathbf{U}, \mathbf{X}_i)$. With this definition and choice distance
 151 measurement, the minimax problem we wish to solve is

$$152 \quad (2.7) \quad \arg \min_{\mathbf{U} \in \text{Gr}(k,n)} \max_{i=1, \dots, M} d_{\text{Gr}(k,n)}(\mathbf{U}, \mathbf{X}_i).$$

153 Using the notion of distance from Equation (2.4), an algorithm was proposed by [25] to
 154 solve Problem (2.7) for a given value of k . Since the data is not of uniform dimension, it is
 155 one of our goals to find the solution across all possible values of k that best represents the
 156 common subspace in the data. In Section 5 we propose an order-selection rule for comparing
 157 solutions of different dimension, however we must first be able to find the solutions of different
 158 dimension efficiently. As we will see in Section 5.1, $\mathbf{U}^*(k) \in \text{Gr}(k, n)$ is not always contained
 159 in $\mathbf{U}^*(k+1) \in \text{Gr}(k+1, n)$, so it is not possible to construct the respective solutions iteratively
 160 via deflation. Instead the problem needs to be solved independently for each value of k .

161 **3. Dual formulation.** Problem (2.7) is nonconvex and challenging to optimize directly.
 162 Therefore, in this section we formulate the dual problem which can be solved efficiently. The
 163 dual variables also provide a primal-feasible solution, which can be tested for optimality.

164 Using Equation (2.6), Problem (2.7) can be written as one with matrix arguments that
 165 can be identified with the Grassmannian points they represent. That is,

$$166 \quad (3.1) \quad \arg \min_{U \in \mathbb{R}^{n \times k}} \max_{i=1, \dots, M} \left(\min\{k, p_i\} - \text{Tr}(U^T X_i X_i^T U) \right)$$

$$\text{s.t. } U^T U = I,$$

167 where U is an orthonormal basis for \mathbf{U} , X_i is an orthonormal basis for \mathbf{X}_i , and $p_i = \dim(\mathbf{X}_i)$.
 168 A solution to (2.7) is then the column space of a solution to (3.1), $\mathbf{U}^* = \text{col}(U^*)$. For ease of
 169 notation we will treat the dual problem as a minimization, so we reformulate the primal as,

$$170 \quad (3.2) \quad \arg \max_{U \in \mathbb{R}^{n \times k}} \min_{i=1, \dots, M} - \left(\min\{k, p_i\} - \text{Tr}(U^T X_i X_i^T U) \right)$$

$$\text{s.t. } U^T U = I.$$

171 Adding an auxiliary variable τ , the quadratic cost function to be minimized is replaced by a
 172 smooth linear objective that is maximized with respect to quadratic inequality constraints,

$$173 \quad (3.3) \quad \arg \max_{U \in \mathbb{R}^{n \times k}, \tau \in \mathbb{R}} \tau$$

$$\text{s.t. } -\tau - \min\{k, p_i\} + \text{Tr}(U^T X_i X_i^T U) \geq 0 \text{ for } i = 1, \dots, M,$$

$$U^T U = I.$$

174 This is essentially the same construction as in [25]. The authors of [25] go on to compute an
 175 intermediate solution to this problem via the Karush–Kuhn–Tucker conditions, and iterate to
 176 a stationary point by taking geodesic steps towards the subspace with the maximum distance
 177 to the current iterate of the primal variable. This contrasts with the proposed approach, where
 178 a solution to (2.7) is found by optimizing the dual problem.

179 Let $\lambda = [\lambda_1, \dots, \lambda_M]^T$ be a vector of Lagrange multipliers associated with the inequality
 180 constraints in (3.3). Dualizing only the inequality constraints leads to the Lagrangian

$$181 \quad (3.4) \quad \mathcal{L}(U, \tau, \lambda) = \tau + \sum_{i=1}^M \lambda_i \left(-\tau - \min\{k, p_i\} + \text{Tr}(U^T X_i X_i^T U) \right),$$

182 such that $U^T U = I$ and $\lambda_i \geq 0$ for $i = 1, \dots, M$. The dual cost function is then found by
 183 maximizing \mathcal{L} over U and τ ,

$$184 \quad (3.5) \quad f(\lambda) = \sup_{\tau} \left(\tau - \sum_{i=1}^M \lambda_i \tau \right) - \sum_{i=1}^M \lambda_i \min\{k, p_i\} + \sup_{U^T U = I} \text{Tr}(U^T \left(\sum_{i=1}^M \lambda_i X_i X_i^T \right) U).$$

185 The maximum over τ yields $f(\lambda) = \infty$ unless $\|\lambda\|_1 = 1$, in which case the first term is
 186 zero. The final term in (3.5) is a well-known problem that is maximized by the sum of
 187 the k largest eigenvalues of $\sum_{i=1}^M \lambda_i X_i X_i^T$ [23]. Let $d_1(\lambda) \geq d_2(\lambda) \geq \dots \geq d_n(\lambda)$ be the
 188 eigenvalues of $\sum_{i=1}^M \lambda_i X_i X_i^T$ and let $\mathbf{v}_1(\lambda), \mathbf{v}_2(\lambda), \dots, \mathbf{v}_n(\lambda)$ be the associated orthonormal
 189 eigenvectors. The argument λ is included to emphasize that the eigendecomposition depends
 190 on λ . The supremum is then $\sum_{j=1}^k d_j(\lambda)$, and is achieved by the matrix whose columns are
 191 the k dominant eigenvectors,

$$192 \quad (3.6) \quad U_{\lambda} \doteq [\mathbf{v}_1(\lambda), \dots, \mathbf{v}_k(\lambda)].$$

193 Thus the dual cost can be written as

$$194 \quad (3.7) \quad f(\lambda) = - \sum_{i=1}^M \lambda_i \min\{k, p_i\} + \sum_{j=1}^k d_j(\lambda),$$

195 and finally, we wish solve the problem,

$$196 \quad (3.8) \quad \arg \min_{\lambda \in \mathbb{R}^M} f(\lambda) \text{ s.t. } \|\lambda\|_1 = 1 \text{ and } \lambda_i \geq 0 \text{ for } i = 1, \dots, M.$$

197 **4. Solution via subgradient.** The dual cost in (3.7) is a locally Lipschitz convex function.
 198 However, it is not differentiable at values of λ for which $d_k(\lambda) = d_{k+1}(\lambda)$, that is, at values for
 199 which the k th and $(k+1)$ st eigenvalues of $\sum_{i=1}^M \lambda_i X_i X_i^T$ are equal [23, Corr. 3.10]. There are
 200 many efficient ways to optimize such a function. In this section we recall how the subgradient
 201 method [30] can be applied to solve this dual problem. After a subgradient has been computed,
 202 the well-developed literature of subgradient algorithms provides a variety of techniques and
 203 step sizes to optimize Problem (3.8) with associated convergence guarantees.

Recall that a vector $\mathbf{g} \in \mathbb{R}^M$ is a subgradient of $f : \mathbb{R}^M \rightarrow \mathbb{R}$ at $\mathbf{x} \in \text{dom } f$ if for all
 $\mathbf{z} \in \text{dom } f$,

$$f(\mathbf{z}) \geq f(\mathbf{x}) + \mathbf{g}^T (\mathbf{z} - \mathbf{x}).$$

204 In this case we denote that \mathbf{g} is in the subdifferential of f at \mathbf{x} by writing $\mathbf{g} \in \partial f(\mathbf{x})$. If f is
 205 differentiable at \mathbf{x} then the gradient is the only subgradient and $\mathbf{g} = \nabla f(\mathbf{x}) = \partial f(\mathbf{x})$.

206 To minimize f in Problem (3.8), the subgradient method uses the iteration

$$207 \quad (4.1) \quad \lambda^{(t+1)} = \Pi(\lambda^{(t)} - \alpha^{(t)} \mathbf{g}^{(t)}),$$

208 where $\alpha^{(t)}$ is a step size selected to guarantee that the sequence $\{\lambda^{(t)}\}_{t=1}^{\infty}$ converges (in
 209 distance) to the optimum, λ^* , and $\Pi : \mathbb{R}^M \rightarrow \{\mathbf{x} : \|\mathbf{x}\|_1 = 1, x_i \geq 0 \text{ for } i = 1, \dots, M\} \subset \mathbb{R}^M$
 210 projects the iterate into the unit simplex.

211 There is a standard trick for computing a subgradient of the dual function that can
 212 adapted to this problem from nonlinear optimization texts such as [5]. Write the Lagrangian
 213 as $\mathcal{L}(U, \tau, \lambda) = q(U, \tau) + \lambda^T \mathbf{g}(U, \tau)$, where $q(U, \tau)$ is the primal objective function and
 214 $\mathbf{g}(U, \tau) \in \mathbb{R}^M$ is the vector of constraint values. Given the dual variable, $\lambda^{(t)} \in \mathbb{R}^M$, at
 215 iteration t , let $(U_{\lambda^{(t)}}, \tau_{\lambda^{(t)}})$ be the primal variable that maximizes the Lagrangian. Then
 216 $\mathbf{g}^{(t)} = \mathbf{g}(U_{\lambda^{(t)}}, \tau_{\lambda^{(t)}})$ is a subgradient of f at $\lambda^{(t)}$.

217 In our case $U_{\lambda^{(t)}}$ is defined according by Equation (3.6) and the i th element of the
 218 constraint vector is $g_i(U_{\lambda^{(t)}}, \tau_{\lambda^{(t)}}) = -\tau_{\lambda^{(t)}} - \min\{k, p_i\} + \text{Tr}(U_{\lambda^{(t)}}^T X_i X_i^T U_{\lambda^{(t)}})$. However, the
 219 constant vector $[-\tau_{\lambda^{(t)}}, \dots, -\tau_{\lambda^{(t)}}]^T \in \mathbb{R}^M$ does not affect the direction after projection onto
 220 the unit simplex, so a subgradient of $f(\lambda^{(t)})$ is

$$221 \quad (4.2) \quad \mathbf{g}^{(t)} = \begin{pmatrix} -\min\{k, p_1\} + \text{Tr}(U_{\lambda^{(t)}}^T X_1 X_1^T U_{\lambda^{(t)}}) \\ \vdots \\ -\min\{k, p_M\} + \text{Tr}(U_{\lambda^{(t)}}^T X_M X_M^T U_{\lambda^{(t)}}) \end{pmatrix}.$$

222 We can check that $\mathbf{g}^{(t)}$ is a subgradient of f as follows. For any $\tilde{\lambda} \in \mathbb{R}^M$ such that $\|\tilde{\lambda}\|_1 = 1$
 223 and $\tilde{\lambda}_i \geq 0$ for $i = 1, \dots, M$ we have

$$224 \quad (4.3) \quad \begin{aligned} f(\lambda^{(t)}) + \mathbf{g}^{(t)T} (\tilde{\lambda} - \lambda^{(t)}) &= f(\lambda^{(t)}) + \mathbf{g}^{(t)T} \tilde{\lambda} - \mathbf{g}^{(t)T} \lambda^{(t)} \\ &= f(\lambda^{(t)}) + \mathbf{g}^{(t)T} \tilde{\lambda} - f(\lambda^{(t)}) \\ &= -\sum_{i=1}^M \tilde{\lambda}_i \min\{k, p_i\} + \text{Tr}(U_{\lambda^{(t)}}^T (\sum_{i=1}^M \tilde{\lambda}_i X_i X_i^T) U_{\lambda^{(t)}}) \\ &\leq -\sum_{i=1}^M \tilde{\lambda}_i \min\{k, p_i\} + \max_{U^T U = I} \text{Tr}(U^T (\sum_{i=1}^M \tilde{\lambda}_i X_i X_i^T) U) \\ &= f(\tilde{\lambda}), \end{aligned}$$

225 and thus $\mathbf{g}^{(t)} \in \partial f(\lambda^{(t)})$. Additionally, it can be verified that this subgradient matches the
 226 general description provided by [23, Thm. 3.9] with the associated affine shift.

227 **4.1. Convergence.** The subgradient $\mathbf{g}^{(t)}$ can be used to update $\lambda^{(t)}$ via the iteration
 228 in (4.1). The subgradient method is not a descent method, so the value of the objective
 229 function at step $t+1$ may be larger than it was at step t . Thus we keep track of the dual variable
 230 with the lowest cost at each iteration and denote it

$$231 \quad (4.4) \quad \lambda_{\text{best}}^{(t+1)} = \begin{cases} \lambda_{\text{best}}^{(t)} & f(\lambda^{(t+1)}) > f(\lambda_{\text{best}}^{(t)}); \\ \lambda^{(t+1)} & \text{otherwise.} \end{cases}$$

232 Given an upper bound on the norm of the subgradients, $\|\mathbf{g}^{(t)}\|_2 \leq G < \infty$ for all t ,
 233 classical theory makes different guarantees on the convergence of the sequence of iterates,
 234 $\{\lambda^{(t)}\}_{t=1}^\infty$, and thus on the sequence of objective function values, $\{f(\lambda_{\text{best}}^{(t)})\}_{t=1}^\infty$, depending on
 235 the choice of step size, $\alpha^{(t)}$. For example, with step sizes independent of iteration like $\alpha^{(t)} = a$
 236 or $\alpha^{(t)} = a/\|\mathbf{g}^{(t)}\|_2$ for some $a > 0$, the subgradient algorithm will converge respectively to
 237 within $G^2 a/2$ or $G a/2$ of the optimal value [5]. Alternatively, if the step size converges to zero
 238 and the sequence is nonsummable or square-summable, that is, $\lim_{t \rightarrow \infty} \alpha^{(t)} = 0$ and

$$239 \quad (4.5) \quad \sum_{t=1}^{\infty} \alpha^{(t)} = \infty \quad \text{or} \quad \sum_{t=1}^{\infty} (\alpha^{(t)})^2 < \infty,$$

240 the subgradient method converges to an optimal objective value, $\lim_{t \rightarrow \infty} f(\lambda_{\text{best}}^{(t)}) = f(\lambda^*)$.
 241 These conditions are satisfied by step sizes like, $\alpha^{(t)} = a/\sqrt{t}$ for $a > 0$, or $\alpha^{(t)} = a/(b+t)$
 242 where $a > 0$ and $b \geq 0$. Proofs of these results can be found in standard literature on convex
 243 optimization for nonsmooth problems such as [5, 13, 30].

244 Although the theory requires $\alpha^{(t)}$ to satisfy the constraints in (4.5) for convergence,
 245 the small step size leads to very slow convergence. In practice we can find an approximate
 246 solution quickly by stepping in the direction of a subgradient but requiring the dual objective to
 247 decrease at each iteration. Algorithm A.1 (in Appendix A) solves Problem (3.8) by performing
 248 a back-tracking line search in the direction of $\mathbf{g}^{(t)} \in \partial f(\lambda^{(t)})$ to ensure that the dual objective
 249 decreases at each step, however, this method is not guaranteed to converge because $\mathbf{g}^{(t)}$ is not
 250 necessarily a descent direction. The practical implementation of Algorithm A.1 is a hybrid of
 251 a back-tracking line search and a nonsummable diminishing step size and for a fixed dimension
 252 k it identifies a stationary point of the dual problem while providing a feasible solution to the
 253 primal problem. It is not intended to be a state-of-the-art subgradient algorithm, but rather
 254 just one example of an implementation that is faster than the standard $a/(b+t)$ square-summable
 255 step size. Alternatively, a well-established quasi-Newton method like the Broyden-Fletcher-
 256 Goldfarb-Shanno (BFGS) algorithm [9] can be used to solve Equation (3.8), but empirically
 257 the convergence rates are comparable to those of the algorithm presented here for this problem.

258 **4.2. Optimality.** In addition to theoretical convergence guarantees, the optimality of a
 259 solution to the dual subgradient approach can be verified in some cases. Let λ^* be a solution
 260 to Problem (3.8). There exists a matrix U_{λ^*} whose columns are the k dominant eigenvectors
 261 of $\sum_{i=1}^M \lambda_i^* X_i X_i^T$, analogous to Equation (3.6). Then U_{λ^*} satisfies $U_{\lambda^*}^T U_{\lambda^*} = I$ and is thus a
 262 feasible solution to the primal problem in (3.1). If the primal and dual objective functions are
 263 equal, strong duality holds and implies that λ^* and $\mathbf{U}^* = \text{col}(U_{\lambda^*})$ are globally optimal dual
 264 and primal variables, respectively. Empirically the duality gap approaches zero for collections
 265 of data that satisfy an implicit assumption of minimax optimization; that the data collection
 266 is free of outliers. Even when strong duality does not hold, the duality gap gives a bound on
 267 the maximum possible improvement for a solution.

268 This verification of optimality is standard for problems where the primal and dual costs
 269 are both computable, but existing techniques for finding the GMEB do not offer this feature.
 270 For instance, using a primal method like [25] does not directly provide a solution to the dual
 271 problem, and thus the duality gap is unknown. Section 7.1 contains numerical experiments
 272 that demonstrate the accuracy of the proposed subgradient method.

273 **5. Proposed order selection rule.** Given a dimension, k , and a finite collection of
 274 subspaces, $\mathcal{D} = \{\mathbf{X}_i \in \text{Gr}(p_i, n)\}_{i=1}^M$, there exist subspaces, $\mathbf{U}^*(k)$, that solve

$$275 \quad (5.1) \quad \arg \min_{\mathbf{U} \in \text{Gr}(k, n)} \max_{i=1, \dots, M} d_{\text{Gr}(k, n)}(\mathbf{U}, \mathbf{X}_i),$$

276 for $k = 1, \dots, \max_i \{\dim(\mathbf{X}_i)\}$. The argument k is now included in the notation for the GMEB
 277 center to emphasize that the subspace depends on the parameter k , and may differ significantly
 278 depending on the value of this parameter. Section 4 described a method to compute $\mathbf{U}^*(k)$
 279 from the associated dual variable, $\lambda^*(k) \in \mathbb{R}^M$. However, because \mathcal{D} contains subspaces of
 280 differing dimension, it is unclear on which Grassmannian the minimum enclosing ball should
 281 be computed. Thus, given the set \mathcal{D} , in this section we would like to determine the optimal
 282 choice for k , in addition to the associated center $\mathbf{U}^*(k)$. Please note a change in notation; the
 283 costs associated with a particular order, k , are more intuitive when the primal is formulated as a
 284 minimization problem and the dual is a maximization. Therefore, as shown in Equation (5.1),
 285 the primal minimization formulation is used for the remainder of the manuscript. The prior
 286 formulation was only used for ease of notation in the subgradient method.

287 All orthogonally invariant distances on $\text{Gr}(k, n)$ can be written as a function of the k
 288 principle angles between a pair of points. It should be clear from the definition in Equation (2.1)
 289 that each angle is bounded above by $\pi/2$, and thus that the squared chordal distance is bounded
 290 above by k . Scaling the primal objective function by $1/k$ normalizes the cost associated with

291 $\mathbf{U}^*(k)$ so that the value of

$$292 \quad (5.2) \quad c_{\text{obj}}(k) := \begin{cases} 0 & k = 0; \\ \max_{i=1, \dots, M} \frac{d_{\text{Gr}(k, n)}(\mathbf{U}^*(k), \mathbf{X}_i)}{k} & k = 1, \dots, \max_i \{\dim(\mathbf{X}_i)\}, \end{cases}$$

293 gives a fair comparison across different values of k . The normalized objective function
 294 achieves its maximum value, $c_{\text{obj}}(k) = 1$, when there exists an i such that $\mathbf{X}_i \perp \mathbf{U}^*(k)$.
 295 That is, $\mathbf{U}^*(k)$ contains no information about at least one of the points in \mathcal{D} . At the other
 296 extreme, the minimum occurs when $k = 0$, and when the point of each $\Omega_*(\mathbf{X}_i)$ closest to the
 297 center coincides with the center. That is, $c_{\text{obj}}(k) = 0$ when $\mathbf{Y}_i^*(k) = \mathbf{U}^*(k)$ for all i , where
 298 $\mathbf{Y}_i^*(k) = \arg \min_{\mathbf{Y}_i \in \Omega_*(\mathbf{X}_i)} d_{\text{Gr}(k, n)}(\mathbf{U}^*(k), \mathbf{Y}_i)$.

299 Simply minimizing $c_{\text{obj}}(k)$ with respect to k is not sufficient to identify the ideal dimension
 300 of $\mathbf{U}^*(k)$ because on average $c_{\text{obj}}(k) \leq c_{\text{obj}}(k+1)$ irrespective of the relationship between the
 301 data points, and of course $c_{\text{obj}}(0) = 0$ by definition. However, the dimension of the ideal center
 302 should represent all the common information without over-fitting, and should also indicate
 303 when no significant relationship exists between the data. Thus we propose a penalty term
 304 based on the dimensions of the data not represented by $\mathbf{U}^*(k)$ that balances the information
 305 lost by making k too small with the lack of specificity that comes from setting k too large.

306 Let $\mathbf{U}^{*\perp}(k)$ denote the orthogonal complement of $\mathbf{U}^*(k)$ and $\tilde{p}_j \doteq \min\{n - k, \dim(\mathbf{X}_j)\}$
 307 for $j = 1, \dots, M$. The expression

$$308 \quad (5.3) \quad c_{\text{pen}}(k) := \begin{cases} 1 & k = 0; \\ \min_{j=1, \dots, M} 1 - \frac{d_{\text{Gr}(\tilde{p}_j, n)}(\mathbf{U}^{*\perp}(k), \mathbf{X}_j)}{\tilde{p}_j} & k = 1, \dots, \max_j \{\dim(\mathbf{X}_j)\}, \end{cases}$$

309 represents the minimum similarity between any point in \mathcal{D} and the dimensions not contained
 310 in the center of the GMEB. A high minimum similarity between points in \mathcal{D} and $\mathbf{U}^{*\perp}(k)$
 311 implies that too much information is being left out of the central subspace, $\mathbf{U}^*(k)$. The penalty
 312 term takes a value of $c_{\text{pen}}(k) = 1$ when $\dim(\mathbf{U}^{*\perp}(k) \cap \mathbf{X}_j) = \tilde{p}_j$ for all j and $c_{\text{pen}}(k) = 0$
 313 when there exists a j for which $\mathbf{X}_j \perp \mathbf{U}^{*\perp}(k)$. The sum of the terms in (5.2) and (5.3) leads
 314 to the proposed rule for selecting the optimal order k^* ,

$$315 \quad (5.4) \quad \arg \min_{k=0, \dots, \max_i \{\dim(\mathbf{X}_i)\}} c_{\text{obj}}(k) + c_{\text{pen}}(k).$$

316 The two terms in (5.4) are computed independently so the GMEB center is not affected by
 317 the penalty term. The value of k^* that minimizes the sum of these two terms corresponds to
 318 the number of subspace dimensions needed to represent the common information present in
 319 \mathcal{D} without over-fitting. Numerical experiments in Section 7.3 demonstrate the efficacy of the
 320 order selection rule on simulated data with ground truth.

321 **5.1. Primal solutions are not nested in general for increasing values of k .** Naively,
 322 the order selection rule in Equation (5.4) can be applied by computing the costs $c_{\text{obj}}(k)$ and
 323 $c_{\text{pen}}(k)$ independently for $k = 0, \dots, \max_i \{\dim(\mathbf{X}_i)\}$ as follows,

- 324 1. Compute $\lambda^*(k)$ using the subgradient method described in Section 4.
- 325 2. Find the associated primal variable, $\mathbf{U}^*(k)$, as the k -dimensional eigenspace of the
 326 weighted sum $\sum_{i=1}^M \lambda_i^*(k) \mathbf{X}_i \mathbf{X}_i^T$.
- 327 3. Compute the orthogonal complement, $\mathbf{U}^{*\perp}(k) = \text{col}(I - \mathbf{U}^*(k) \mathbf{U}^{*T}(k))$.

328 Then k^* is selected as the value of k associated with the minimum cost, $c_{\text{obj}}(k) + c_{\text{pen}}(k)$.
 329 If $\lambda^*(k) = \lambda^*(k+1)$ for some $k < \max_i \{\dim(\mathbf{X}_i)\}$ then the solution on $\text{Gr}(k+1, n)$ can be

330 constructed in a greedy fashion as the direct sum of the solution on $\text{Gr}(k, n)$ and the $(k + 1)$ st
 331 eigenvector of $\sum_{i=1}^M \lambda_i^*(k) X_i X_i^T$. Unfortunately, the dual variables are not generally equal for
 332 increasing values of k , so a greedy approach is not appropriate.

333 Observe that the central subspaces are not nested for increasing dimensions in the follow-
 334 ing illustrative example. Let

$$335 \quad (5.5) \quad X_1 = \begin{bmatrix} \frac{\sqrt{2}}{\sqrt{3}} & 0 \\ \frac{1}{\sqrt{6}} & 0 \\ \frac{1}{\sqrt{6}} & 0 \\ 0 & \frac{\sqrt{7}}{\sqrt{8}} \\ 0 & \frac{1}{\sqrt{8}} \end{bmatrix}, \quad X_2 = \begin{bmatrix} \frac{1}{\sqrt{6}} & 0 \\ \frac{\sqrt{2}}{\sqrt{3}} & 0 \\ \frac{1}{\sqrt{6}} & 0 \\ 0 & \frac{1}{\sqrt{8}} \\ 0 & \frac{\sqrt{7}}{\sqrt{8}} \end{bmatrix}, \quad \text{and } X_3 = \begin{bmatrix} \frac{1}{\sqrt{6}} \\ \frac{1}{\sqrt{6}} \\ \frac{\sqrt{2}}{\sqrt{3}} \\ 0 \\ 0 \end{bmatrix},$$

336 be orthonormal bases for the three points $\mathbf{X}_1, \mathbf{X}_2 \in \text{Gr}(2, 5)$ and $\mathbf{X}_3 \in \text{Gr}(1, 5)$. One can check
 337 that the subspace that minimizes the maximum distance to these three points on $\text{Gr}(1, 5)$ is
 338 the mean of their first columns. That is, the optimal primal and dual variables are

$$339 \quad (5.6) \quad \mathbf{U}^*(1) = \text{col} \left(\begin{bmatrix} \frac{1}{\sqrt{3}} & \frac{1}{\sqrt{3}} & \frac{1}{\sqrt{3}} & 0 & 0 \end{bmatrix}^T \right), \quad \text{and } \lambda^*(1) = \begin{bmatrix} \frac{1}{\sqrt{3}} & \frac{1}{\sqrt{3}} & \frac{1}{\sqrt{3}} \end{bmatrix}^T,$$

340 with associated primal and dual costs of

$$341 \quad (5.7) \quad \min_{\mathbf{U} \in \text{Gr}(1,5)} \max_{i=1,2,3} d_{\text{Gr}(1,5)}(\mathbf{U}, \mathbf{X}_i) = \max_{\lambda \in \mathbb{R}^3} \min_{U^T U = I} 1 - \sum_{i=1}^3 \lambda_i \text{Tr}(U^T Y_i Y_i^T U) = \frac{1}{9}.$$

342 The duality gap in Equation (5.7) is zero, indicating that this is a global solution.

343 On $\text{Gr}(2, 5)$, however, $\Omega_+(\mathbf{X}_3)$ consists of subspaces that span X_3 and any orthogonal
 344 direction. In particular there exists $\mathbf{Y}_3 \in \Omega_+(\mathbf{X}_3)$ such that the second column of Y_3 is
 345 $[0 \ 0 \ 0 \ 1/\sqrt{2} \ 1/\sqrt{2}]^T$. This leads to a solution for the center of the minimum enclosing ball on
 346 $\text{Gr}(2, 5)$ given by primal and dual variables

$$347 \quad (5.8) \quad \mathbf{U}^*(2) = \text{col} \left(\begin{bmatrix} \frac{3}{\sqrt{22}} & \frac{3}{\sqrt{22}} & \frac{2}{\sqrt{22}} & 0 & 0 \\ 0 & 0 & 0 & \frac{1}{\sqrt{2}} & \frac{1}{\sqrt{2}} \end{bmatrix}^T \right), \quad \text{and } \lambda^*(2) = \begin{bmatrix} \frac{1}{2} & \frac{1}{2} & 0 \end{bmatrix}^T.$$

348 Notably, \mathbf{X}_3 is not in the support of the minimum enclosing ball on $\text{Gr}(2, 5)$ and thus does not
 349 influence the central subspace. Strong duality also holds for this solution with

$$350 \quad (5.9) \quad \min_{\mathbf{U} \in \text{Gr}(2,5)} \max_{i=1,2,3} d_{\text{Gr}(2,5)}(\mathbf{U}, \mathbf{X}_i) = \max_{\lambda \in \mathbb{R}^3} \min_{U^T U = I} 2 - \sum_{i=1}^3 \lambda_i \text{Tr}(U^T Y_i Y_i^T U) = \frac{14 - 3\sqrt{7}}{24}.$$

351 Since $\mathbf{U}^*(1)$ is orthogonal to the second dimension of $\mathbf{U}^*(2)$ and noncollinear with the first,
 352 and the columns of $\mathbf{U}^*(2)$ are orthogonal, we have $\mathbf{U}^*(1) \notin \mathbf{U}^*(2)$. Additionally we find that
 353 optimal order selected by applying the rule in Equation (5.4) is $k^* = 1$, because

$$354 \quad (5.10) \quad \begin{aligned} c_{\text{obj}}(0) + c_{\text{pen}}(0) &= 0 + 1 = 1, \\ c_{\text{obj}}(1) + c_{\text{pen}}(1) &= \frac{1}{1} \left(\frac{1}{9} \right) + \frac{1}{1} \left(1 - \left(\frac{\sqrt{8}}{\sqrt{9}} \right)^2 \right) \approx 0.22, \quad \text{and} \\ c_{\text{obj}}(2) + c_{\text{pen}}(2) &= \frac{1}{2} \left(\frac{14 - 3\sqrt{7}}{24} \right) + \frac{1}{2} \left(2 - \left(\frac{-1}{\sqrt{12}} \right)^2 + \left(\frac{1 - \sqrt{7}}{\sqrt{16}} \right)^2 \right) \approx 0.25. \end{aligned}$$

355 This agrees with the intuition that the center of the minimum enclosing ball represents the
 356 common information in all points without over-fitting to any subset of points, but note that the
 357 optimal order is not always the dimension of the smallest subspace. The common subspace
 358 may have dimension smaller than any of the samples or there may be no common subspace.

359 Even though the primal solutions are not always nested, a good initial guess for the dual
 360 variable will reduce computational overhead. One benefit of the subgradient approach is that
 361 $\lambda^*(k)$ is computed explicitly. Thus we can initialize the algorithm with $\lambda^{(0)}(k+1) = \lambda^*(k)$.
 362 The impact of this heuristic warm-start is discussed in the experiments in Section 7.2.

363 **5.2. Related literature on order fitting for subspace averaging.** A recent work from
 364 Santamaría *et al.* [28] also attempts to find a central subspace of ambiguous dimension. The
 365 authors minimize the mean-squared error (MSE) between a subspace and a collection of data
 366 in the space of $n \times n$ projection matrices using the squared Frobenius norm. That is,

$$367 \quad (5.11) \quad E(k) = \min_{\mathbf{U} \in \text{Gr}(k,n)} \frac{1}{M} \sum_{i=1}^M \|UU^T - X_i X_i^T\|_F^2.$$

368 Putting aside for a moment that the current work is interested in minimizing the maximum
 369 deviation rather than the mean-squared error, there remains a central difference between the
 370 technique in [28] and the proposed method. The optimization of Equation (5.11) is done
 371 in a vector space, after which the solution is mapped to the nearest point on the Grassmann
 372 manifold. This is subtly different than minimizing the MSE on the Grassmannian with respect
 373 to the squared chordal distance using the point-to-set interpretation of [38]. To see this, write
 374 half of the squared distance from [28] between the central subspace and the i th point as

$$375 \quad (5.12) \quad \begin{aligned} \frac{1}{2} \|U^*(k)U^{*T}(k) - X_i X_i^T\|_F^2 &= \frac{k+p_i}{2} - \sum_{r=1}^{\min\{k,p_i\}} \cos^2(\theta_r(\mathbf{U}^*(k), \mathbf{X}_i)) \\ &= \frac{|k-p_i|}{2} + \sum_{r=1}^{\min\{k,p_i\}} \sin^2(\theta_r(\mathbf{U}^*(k), \mathbf{X}_i)). \end{aligned}$$

376 In contrast, the point-to-set squared chordal distance on $\text{Gr}(k, n)$ is

$$377 \quad (5.13) \quad \begin{aligned} d_{\text{Gr}(k,n)}(\mathbf{U}^*(k), \mathbf{X}_i) &= \min \{d(\mathbf{U}^*(k), \mathbf{Y}_i) : \mathbf{Y}_i \in \Omega(\mathbf{X}_i)\} \\ &= \min \left\{ \frac{1}{2} \|U^*(k)U^{*T}(k) - Y_i Y_i^T\|_F^2 : \mathbf{Y}_i \in \Omega(\mathbf{X}_i) \right\} \\ &= k - \sum_{r=1}^k \cos^2(\theta_r(\mathbf{U}^*(k), \mathbf{Y}_i)) \\ &= \sum_{r=1}^{\min\{k,p_i\}} \sin^2(\theta_r(\mathbf{U}^*(k), \mathbf{X}_i)) \end{aligned}$$

378 because $0 = \theta_{p_i}(\mathbf{U}^*(k), \mathbf{Y}_i) = \theta_{p_i+1}(\mathbf{U}^*(k), \mathbf{Y}_i) = \dots = \theta_k(\mathbf{U}^*(k), \mathbf{Y}_i)$ if $p_i < k$ by the
 379 definition of \mathbf{Y}_i in Equation (2.5). Thus the distances differ by $\frac{|k-p_i|}{2}$, which is the difference
 380 in dimensions between the central subspace and the i th data point.

381 The slight difference in distance measurements lends itself to an interesting interpretation
 382 when determining the appropriate rank of the central subspace. The solution, $\mathbf{U}^*(k)$, to

$$383 \quad (5.14) \quad \arg \min_{\mathbf{U} \in \text{Gr}(k,n)} \frac{1}{M} \sum_{i=1}^M \|UU^T - X_i X_i^T\|_F^2$$

384 for a fixed k is the dominant k -dimensional eigenspace of the sum $\frac{1}{M} \sum_{i=1}^M X_i X_i^T$. That is, if

$$385 \quad (5.15) \quad \frac{1}{M} \sum_{i=1}^M X_i X_i^T = F D F^T$$

386 is an eigendecomposition with eigenvectors $F = [\mathbf{f}_1, \mathbf{f}_2, \dots, \mathbf{f}_R]$ and associated eigenvalues
 387 $d_1 \geq d_2 \geq \dots \geq d_R$, then the solution to Equation (5.14) is $\mathbf{U}^*(k) = [\mathbf{f}_1, \mathbf{f}_2, \dots, \mathbf{f}_k]$. Note
 388 that this $\mathbf{U}^*(k)$ is not the same subspace as the center of the minimum enclosing ball. The
 389 MSE in Equation (5.11) can be written as a function of all R eigenvalues,

$$390 \quad (5.16) \quad E(k) = \sum_{r=1}^k 1 - d_r + \sum_{r=k+1}^R d_r,$$

391 and the minimum of Equation (5.16) is achieved when k^* is the smallest value for which
 392 $d_{k+1} < 0.5$. This eigenvalue threshold is then fixed regardless of the dimension of the ambient
 393 space, and as we will see in Section 7.3, the selected dimension could differ drastically for
 394 noisy data depending on the ambient dimension.

395 For a different interpretation of the k^* that minimizes Equation (5.11) we can rewrite
 396 Equation (5.16) as a function of the angles between each eigenvector and the subspaces,

$$397 \quad (5.17) \quad E(k) = \sum_{r=1}^k 1 - \mathbf{f}_r^T \left(\frac{1}{M} \sum_{i=1}^M X_i X_i^T \right) \mathbf{f}_r + \sum_{r=k+1}^R \mathbf{f}_r^T \left(\frac{1}{M} \sum_{i=1}^M X_i X_i^T \right) \mathbf{f}_r$$

$$398 \quad (5.18) \quad = \sum_{r=1}^k 1 - \frac{1}{M} \sum_{i=1}^M \cos^2(\theta(\mathbf{f}_r, \mathbf{X}_i)) + \sum_{r=k+1}^R \frac{1}{M} \sum_{i=1}^M \cos^2(\theta(\mathbf{f}_r, \mathbf{X}_i))$$

$$399 \quad (5.19) \quad = \sum_{r=1}^k \frac{1}{M} \sum_{i=1}^M \sin^2(\theta(\mathbf{f}_r, \mathbf{X}_i)) + \sum_{r=k+1}^R \frac{1}{M} \sum_{i=1}^M \sin^2\left(\frac{\pi}{2} - \theta(\mathbf{f}_r, \mathbf{X}_i)\right)$$

$$400 \quad (5.20) \quad = \sum_{r=1}^k \frac{1}{M} \sum_{i=1}^M \sin^2(\theta(\mathbf{f}_r, \mathbf{X}_i)) + \sum_{r=k+1}^R \frac{1}{M} \sum_{i=1}^M \sin^2(\theta(\mathbf{f}_r, \mathbf{X}_i^\perp))$$

$$401 \quad (5.21) \quad = \sum_{r=1}^k \frac{1}{M} \sum_{i=1}^M d_{\text{Gr}(1,n)}(\mathbf{f}_r, \mathbf{X}_i) + \sum_{r=k+1}^R \frac{1}{M} \sum_{i=1}^M d_{\text{Gr}(1,n)}(\mathbf{f}_r, \mathbf{X}_i^\perp).$$

402 The equality between (5.19) and (5.20) is due to [16, Thm. 2.7] which implies that $\frac{\pi}{2} -$
 403 $\theta(\mathbf{f}_r, \mathbf{X}_i) = \theta(\mathbf{f}_r, \mathbf{X}_i^\perp)$. Note, however, that Equation (5.21) is *not* equivalent to

$$404 \quad (5.22) \quad \frac{1}{M} \sum_{i=1}^M d_{\text{Gr}(k,n)}(\mathbf{U}^*(k), \mathbf{X}_i) + \frac{1}{M} \sum_{i=1}^M d_{\text{Gr}(R-k,n)}(\mathbf{U}^{*\perp}(k), \mathbf{X}_i^\perp)$$

405 because linear combinations of the eigenvectors, \mathbf{f}_r , are not included in the expression.
 406 A new interpretation of the MSE-minimizing k becomes fairly apparent in light of Equa-
 407 tion (5.21). The optimal k^* is the one that minimizes the mean-squared chordal distance
 408 between $\{\mathbf{f}_1, \dots, \mathbf{f}_k\}$ and the data points, plus the mean-squared chordal distance between
 409 $\{\mathbf{f}_{k+1}, \dots, \mathbf{f}_R\}$ and the orthogonal complements of the data points.

410 **5.3. Hybrid rule.** It is possible to create a hybrid of the order-selection rule of [28] and
 411 the proposed method with a slight modification. In [12], a robustification of the technique
 412 in [28] is proposed that leads to a weighted eigenvalue decomposition at optimality. The

413 weights are determined using a variety of robust objective functions via a majorization-
 414 minimization scheme, which results in a down-weighting of outliers in the data. By minimizing
 415 the mean-squared error of the *weighted* average (similar to Equation (5.11)), this amounts to
 416 a hard eigenvalue threshold with the order chosen to be the number of dimensions with
 417 eigenvalues greater than 0.5.

418 For the hybrid method, weights will come from the values of the dual variable, $\lambda^*(k)$, at
 419 optimality. Since these values depend on the parameter k , the hard eigenvalue threshold is not
 420 applicable. Let $d_1(k) \geq d_2(k) \geq \dots \geq d_R(k)$ be the eigenvalues of $\sum_{i=1}^M \lambda_i^*(k) X_i X_i^T$ where
 421 $\lambda^*(k)$ is the vector of optimal dual variables computed for the GMEB on $\text{Gr}(k, n)$ using the
 422 proposed algorithm. For $k = 0$, let $\lambda_i^*(0) = \frac{1}{M}$ for $i = 1, \dots, M$. We define a modified version
 423 of the MSE from Equation (5.16) as

$$424 \quad (5.23) \quad \tilde{E}(k) = \sum_{r=1}^k 1 - d_r(k) + \sum_{r=k+1}^R d_r(k).$$

425 The order-selection rule of [28] applied to the GMEB center is then

$$426 \quad (5.24) \quad \arg \min_{k=0, \dots, \max_i \{\dim(\mathbf{X}_i)\}} \tilde{E}(k).$$

427 It should be clear that the eigenvalues $\{d_r(k)\}_{r=1}^R$ will be different for different values of
 428 $\lambda^*(k)$. In the experiments of Section 7.3, this combined method is referred to as ‘‘Hybrid’’ and
 429 performs favorably for all tests; out-performing the other techniques in 2 out of 3 scenarios.

430 **6. Synthetic data generation.** The numerical experiments in Section 7 require data
 431 for which the ground truth is known, and ideally data for which the center of the GMEB is
 432 distinct from the other generalized Grassmannian means. Thus, in this section we propose
 433 two different models for sampling points nonuniformly from a unit ball on the Grassmannian.
 434 The first is an asymmetrical nested ball structure, and the second samples more densely within
 435 a randomly selected arc of the boundary of a unit ball.

436 **6.1. Asymmetrical nested ball model.** A collection of subspaces, $\mathcal{D} = \{\mathbf{X}_i\}_{i=1}^M$, are
 437 uniformly sampled from two balls, $\mathcal{B}_{\epsilon_2}(\mathbf{Z}_2) \subset \mathcal{B}_{\epsilon_1}(\mathbf{Z}_1) \subset \text{Gr}(k_0, n)$ with centers at $\mathbf{Z}_1, \mathbf{Z}_2$
 438 and corresponding radii $\epsilon_1 > \epsilon_2$, respectively. The larger ball, $\mathcal{B}_{\epsilon_1}(\mathbf{Z}_1)$, is the minimum
 439 enclosing ball of the data so that $\mathbf{U}^*(k_0) = \mathbf{Z}_1$. The smaller ball is fully contained within
 440 the larger ball, i.e., $\mathcal{B}_{\epsilon_2}(\mathbf{Z}_2) \subset \mathcal{B}_{\epsilon_1}(\mathbf{Z}_1)$, but $\mathbf{Z}_1 \notin \mathcal{B}_{\epsilon_2}(\mathbf{Z}_2)$. Let M_1, M_2 be the number of
 441 points sampled from $\mathcal{B}_{\epsilon_1}(\mathbf{Z}_1), \mathcal{B}_{\epsilon_2}(\mathbf{Z}_2)$ respectively, with $M = M_1 + M_2$. When $M_2 = 0$, the
 442 generalized Grassmannian means are all equal to the point \mathbf{Z}_1 . When more points are sampled
 443 from $\mathcal{B}_{\epsilon_2}(\mathbf{Z}_2)$ and the fraction M_2/M_1 grows, the generalized Grassmannian means for $p < \infty$
 444 move away from \mathbf{Z}_1 in the direction of \mathbf{Z}_2 , making the averages distinct without affecting the
 445 center of the GMEB. The radius of the large ball, ϵ_1 , controls the similarity of the data points.

446
 447 As described, the data points are all sampled from a single manifold, $\text{Gr}(k_0, n)$. If ϵ_1 is
 448 small enough, then the optimal rank for the GMEB (or any of the generalized Grassmannian
 449 means) is $k^* = k_0$. This construction can be generalized in two ways.

- 450 1. For $i = 1, \dots, M$, the basis for \mathbf{X}_i can be completed to a p_i -dimensional subspace
 451 by taking the span of X_i and $p_i - k_0$ random dimensions. If the $p_i - k_0$ random
 452 dimensions are mutually orthogonal for $i = 1, \dots, M$, then the optimal rank for the
 453 GMEB is still $k^* = k_0$.
- 454 2. Points from the large ball can be sampled from one manifold, $\mathcal{B}_{\epsilon_1}(\mathbf{Z}_1) \subset \text{Gr}(k_1, n)$
 455 while points from the small ball are sampled from another, $\mathcal{B}_{\epsilon_2}(\mathbf{Z}_2) \subset \text{Gr}(k_2, n)$.
 456 If $k_1 \neq k_2$, the optimal rank of the central subspace is ambiguous. Experiments

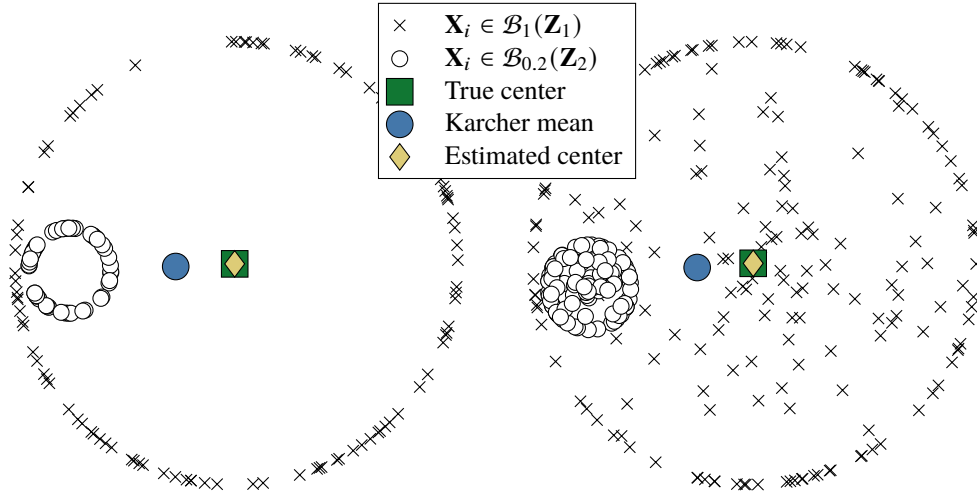


Fig. 2: Two examples of point sets from $\text{Gr}(1,3)$ generated using the nested ball model embedded into \mathbb{R}^2 by multidimensional scaling. The points from $\mathcal{B}_1(\mathbf{Z}_1)$ are indicated with x 's, points from $\mathcal{B}_{0.2}(\mathbf{Z}_2)$ are marked with white circles, the true center is the green square, the Karcher mean is the blue circle, and the estimated GMEB center is the yellow diamond.

457 show that using the proposed order selection rule, $k^* = k_1$ independent of other
 458 parameters, but using the criteria of [28], k^* depends on ϵ_1 and M_2/M_1 .
 459 As an illustrative example, Figure 2 shows 2-dimensional embeddings via multidimensional
 460 scaling of data sets on $\text{Gr}(1,3)$ that have been generated according to the asymmetrical nested
 461 ball model. The yellow diamond indicates the center of the GMEB (computed via the proposed
 462 method) and the blue circle marks the Karcher mean of each data collection.

463 **6.2. Unit ball with higher sampling density from a random arc.** Another practical
 464 scenario where the GMEB center may differ from other generalized Grassmannian means is
 465 when data has been sampled unevenly. This setting is simulated by selecting a random arc
 466 from the boundary of a unit ball and sampling additional points from that region. A collection
 467 of subspaces, $\mathcal{D} = \{\mathbf{X}_i\}_{i=1}^M$, are uniformly sampled from the ball $\mathcal{B}_{\epsilon_1}(\mathbf{Z}_1) \subset \text{Gr}(k_0, n)$ with
 468 center at \mathbf{Z}_1 and radius ϵ_1 . M_1 points are sampled from $\mathcal{B}_{\epsilon_1}(\mathbf{Z}_1)$ so that $\mathbf{U}^*(k_0) = \mathbf{Z}_1$. Two
 469 points are randomly selected from the boundary of $\mathcal{B}_{\epsilon_1}(\mathbf{Z}_1)$, and M_2 additional points are
 470 uniformly sampled from the arc connecting them on the boundary to create $M = M_1 + M_2$
 471 samples. The data points are all sampled from a single manifold, $\text{Gr}(k_0, n)$, and for sufficiently
 472 small ϵ_1 , the optimal rank for the GMEB (or any of the generalized Grassmannian means)
 473 is $k^* = k_0$. To generalize this construction, additional dimensions can be included to create
 474 points from a disjoint union of Grassmannians.

475 For $i = 1, \dots, M$, the basis for \mathbf{X}_i can be completed to a p_i dimensional subspace by taking
 476 the span of \mathbf{X}_i and $p_i - k_0$ random dimensions. If the $p_i - k_0$ random dimensions are mutually
 477 orthogonal for $i = 1, \dots, M$, then the optimal rank for the GMEB is still $k^* = k_0$. Figure 3
 478 shows 2-dimensional embeddings via multidimensional scaling of data sets on $\text{Gr}(1,3)$ that
 479 have been generated as a unit ball with higher sampling density along a random arc. The
 480 yellow diamond indicates the center of the GMEB (computed via the proposed method) and
 481 the blue circle marks the Karcher mean of each data collection.

482 It should be noted that using either data model the point at the center of $\mathcal{B}_{\epsilon_1}(\mathbf{Z}_1)$ is

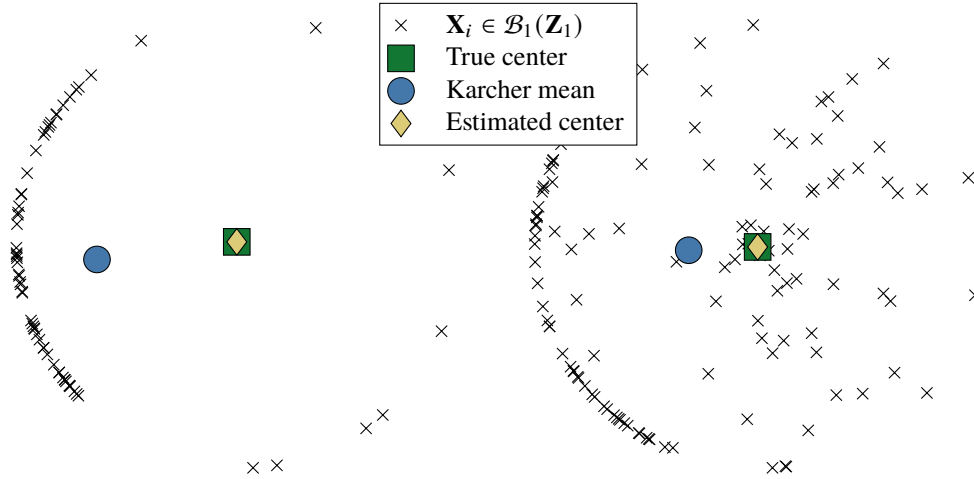


Fig. 3: Two examples of point sets from $\text{Gr}(1,3)$ on the unit ball, $\mathcal{B}_1(\mathbf{Z}_1)$, sampled with nonuniform density on the boundary, embedded into \mathbb{R}^2 by multidimensional scaling. Points from $\mathcal{B}_1(\mathbf{Z}_1)$ are indicated with x 's, the true center is the green square, the Karcher mean is the blue circle, and the estimated GMEB center is the yellow diamond.

483 only the ground-truth center of the minimum enclosing ball of the data collection, $\mathbf{U}(k^*)$, if
 484 the points have been sampled with a high enough density from the surface of the ball. The
 485 minimum number uniformly distributed points needed grows with the ambient dimension, n ,
 486 so in high dimensional spaces the number of points, M , needed to create a ground-truth center
 487 may become prohibitively large. The experimental data can be generated exclusively from the
 488 boundary of the balls or interior points can be added.¹

489 **7. Numerical experiments.** The experiments in this section are meant to illustrate three
 490 properties of the proposed GMEB algorithm and associated order-selection rule. First, we
 491 demonstrate the speed and accuracy of the proposed method for estimating the center of the
 492 GMEB. Second, we demonstrate that a warm-start on $\text{Gr}(k+1, n)$ using the optimal solution
 493 from $\text{Gr}(k, n)$ can reduce the number of iterations required for the algorithm to converge.
 494 And finally, we compare results of the proposed order-selection rule and the rule of [28] in a
 495 variety of scenarios to gain intuition about when and how they differ.

496 **7.1. Experiment 1: Accuracy of the GMEB.** To test the accuracy and efficiency of
 497 the proposed dual subgradient approach, data sets are generated according to the each of two
 498 data models from Section 6. For each data collection, the GMEB center is approximated
 499 using the proposed method and the algorithm of Renard *et al.* [25], and the residual error
 500 is measured as the between the approximate centers and the true centers. For the first data
 501 set, $M = 100$ points are sampled from $\text{Gr}(3, 10)$ using the asymmetrical nested ball model
 502 in Section 6.1 with neither of the proposed generalizations. That is, $k_0 = k_1 = k_2 = 3$ so
 503 that all points are sampled from the same Grassmann manifold. $M_1 = 70$ of the points come
 504 from the boundary of $\mathcal{B}_1(\mathbf{Z}_1)$ and $M_2 = 30$ from the boundary of $\mathcal{B}_{0.125}(\mathbf{Z}_2)$. No points are
 505 sampled from the interior of either ball. Both algorithms are initialized using the extrinsic
 506 mean of the data [20, 26], that is, $\lambda^{(0)} = [1/100, 1/100, \dots, 1/100]^T$, and $\mathbf{U}^{(0)}(3)$ is the dominant

¹Matlab code for the algorithms, data generation procedures, and numerical experiments in this manuscript is available at <https://sites.google.com/site/nicolasgillis/code>.

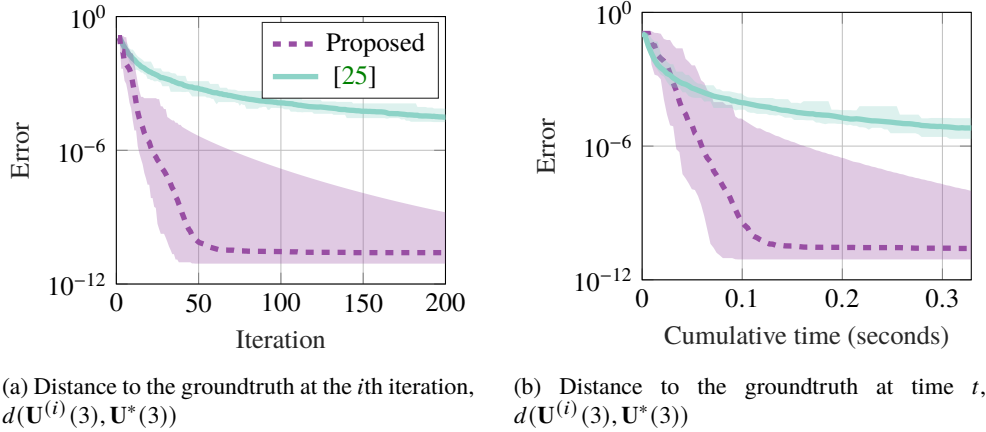


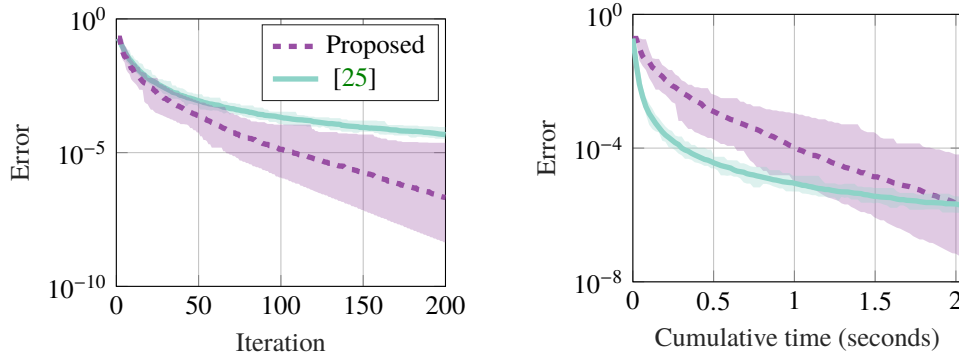
Fig. 4: Median distance to the groundtruth and cumulative time for the GMEB on $\text{Gr}(3, 10)$ of data generated with the asymmetrical nested ball model from Section 6.1 over 100 Monte Carlo trials. The data consists of 100 points in $\text{Gr}(3, 10)$. The proposed method is indicated by the dashed purple line and the method of Renard *et al.* [25] is represented by the solid turquoise line. The shaded regions span the extreme values.

507 3-dimensional eigenspace of $\sum_{i=1}^{100} \lambda_i^{(0)} X_i X_i^T$. The groundtruth center is $\mathbf{U}^*(3) = \mathbf{Z}_1$.

508 Figure 4a shows the median distance to the groundtruth over 100 Monte Carlo trials
 509 between the iterate with the lowest primal cost and the ground-truth center. Figure 4b shows
 510 the same median distance to the groundtruth relative to cumulative computation time for each
 511 algorithm. In both plots the proposed method is indicated by the dashed purple line and the
 512 method of [25] is represented by the solid turquoise line. The shaded regions denote the
 513 complete range of values across all trials. This is a setting in which all data points live on a
 514 single Grassmann manifold. Therefore the point-to-set distances reduce to the traditional
 515 Grassmannian distances and the technique of [25] is equivalent to that of [3].

516 The proposed method clearly outperforms the existing technique in terms of accuracy
 517 relative to both iterations and computation time for this collection of data. However, the
 518 cumulative computation time is affected by many of the parameters in the experimental
 519 setup. Let $P = \max_i \{\dim(\mathbf{X}_i)\}$. For the technique of [3, 25], the per iteration complexity is
 520 $\mathcal{O}(MP(nk + k^2))$ due to the M matrix products and subsequent thin SVDs. The proposed
 521 method computes these same M products and SVDs, but must additionally compute the
 522 compact SVD of a matrix of size $n \times MP$ in order to get the updated center. Assuming that
 523 $n \leq MP$ (as it is in all the experiments), the complexity of the proposed algorithm is then
 524 $\mathcal{O}(MP(nk + k^2 + n^2))$. There are an additional M SVDs for each back-tracking step taken,
 525 but those steps are infrequent and thus dominated by the other terms. From these complexities
 526 we can see that an increase in the ambient dimension, n , number of subspaces, M , or subspace
 527 dimension, P , would all lead to a relative decrease in the efficiency of the proposed method.

528 In the second example we employ the data model from Section 6.2, with the inclusion
 529 of interior points and the generalization that the data points come from a disjoint union of
 530 Grassmannians, that is, they are subspaces of differing dimensions. Initially, $M_1 = 100$ points
 531 are sampled from the boundary of $\mathcal{B}_1(\mathbf{Z}_1)$ on $\text{Gr}(3, 15)$. An additional $M_2 = 100$ points
 532 are selected from an arc on the boundary of the ball between two randomly selected points.
 533 Finally $M_3 = 100$ points are selected uniformly at random from the interior of the ball. Each



(a) Distance to the groundtruth at the i th iteration, $d(\mathbf{U}^{(i)}(3), \mathbf{U}^*(3))$ (b) Distance to the groundtruth at time t , $d(\mathbf{U}^{(t)}(3), \mathbf{U}^*(3))$

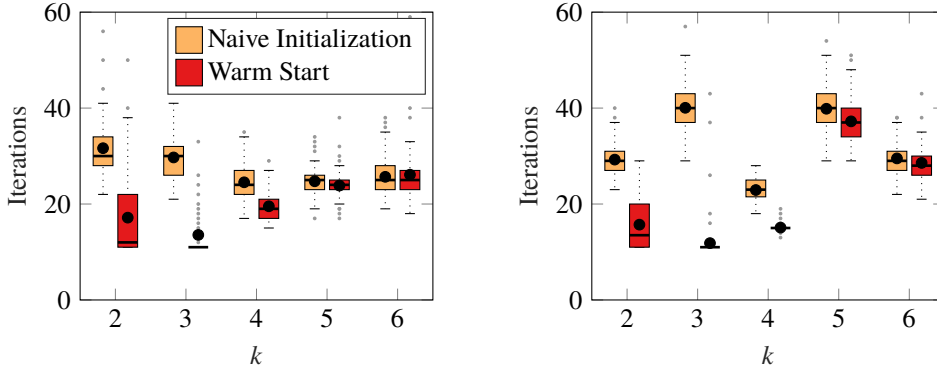
Fig. 5: Median distance to the groundtruth and cumulative computation time for the GMEB on $\text{Gr}(3, 15)$ of data generated with the nonuniform sampling model from Section 6.2 over 100 Monte Carlo trials. The data consists of 300 points in $\coprod_{p \in \mathcal{P}} \text{Gr}(p, 15)$ for $\mathcal{P} = \{3, 4, 5, 6\}$. The proposed method is indicated by the dashed purple line and the method of Renard *et al.* [25] is represented by the solid turquoise line. The shaded regions span the extreme values.

534 of the $M = 300$ points is then completed to a basis for a p_i -dimensional subspace where p_i is
 535 randomly selected from the set $\mathcal{P} = \{3, 4, 5, 6\}$. Both algorithms are again initialized using
 536 the extrinsic mean of the data on $\text{Gr}(3, 15)$ where $\lambda^{(0)} = [1/300, 1/300, \dots, 1/300]^T$, and $\mathbf{U}^{(0)}(3)$
 537 is the dominant 3-dimensional eigenspace of $\sum_{i=1}^{300} \lambda_i^{(0)} \mathbf{X}_i \mathbf{X}_i^T$. Figure 5a shows the median
 538 distance to the groundtruth over 100 Monte Carlo trials between the iterate with the lowest
 539 primal cost and the ground-truth center, while Figure 5b shows the median error relative to
 540 cumulative computation time. The proposed method is indicated by the dashed purple line
 541 and the method of Renard *et al.* [25] is represented by the solid turquoise line. The shaded
 542 regions span the extreme values. The groundtruth center is $\mathbf{U}^*(3) = \mathbf{Z}_1$.

543 As shown in Figure 5a, the proposed method achieves a higher accuracy in fewer iterations
 544 than [25]. However, the greater complexity of the proposed method means that the primal
 545 algorithm initially achieves a lower error, as shown in Figure 5b. The increased number of
 546 points in the data set and specifically in the support of the GMEB lead to a slower overall
 547 convergence for the proposed algorithm. This reduced efficiency would grow with the size of
 548 the data, however the subgradient technique is consistently achieves lower overall error given
 549 enough time. Moreover, the proposed method provides duality-gap optimality guarantees.

550 One direction for future work is to combine the two methods to get the best of both
 551 worlds; fast initial estimates of the center and high accuracy solutions over time. Using
 552 $\mathbf{U}^{(t)}(k)$ computed via t iterations of [25] as an estimate of the center, we can find dual-feasible
 553 variables that are non-zero only for points in the support set of the enclosing ball centered at
 554 $\mathbf{U}^{(t)}(k)$. For example, let $\mathcal{I} = \{i : d_{\text{Gr}(k,n)}(\mathbf{U}^{(t)}(k), \mathbf{X}_i) = \max_i d_{\text{Gr}(k,n)}(\mathbf{U}^{(t)}(k), \mathbf{X}_i)\}$. Then
 555 let $\lambda_i^{(0)} = 1/|\mathcal{I}|$ for $i \in \mathcal{I}$ and $\lambda_i^{(0)} = 0$ otherwise, and proceed with the subgradient algorithm
 556 from this warm-start. An alternative initialization strategy is proposed in Section 7.2.

557 **7.2. Experiment 2: Faster convergence by initializing with previous solutions.** To
 558 apply the order selection criteria in Section 5, the GMEB center must be computed for
 559 $k = 1, \dots, \max_i \{\dim(\mathbf{X}_i)\}$. The example in Section 5.1 demonstrates that the subspace at



(a) Results from 100 trials with the asymmetrical nested ball model where $k^* = 4$ and $M = 50$ points sampled from $\text{Gr}(p_i, 10)$ with $p_i \in \{4, 5, 6\}$.

(b) Results from 100 trials with the nonuniform sampling model where $k^* = 4$ and $M = 300$ points sampled from $\text{Gr}(p_i, 10)$ with $p_i \in \{4, 5, 6\}$.

Fig. 6: Number of iterations needed for the proposed subgradient algorithm to reach a stationary point using a naive initialization, $\lambda^{(0)}(k+1) = [1/M, 1/M, \dots, 1/M]^T$ (light orange), and a warm start, $\lambda^{(0)}(k+1) = \lambda^*(k)$ (red) for two data sets.

560 the center of the minimum enclosing ball cannot be built in a greedy fashion, because the
 561 center $\mathbf{U}^*(k-1) \in \text{Gr}(k-1, n)$ is not in general a subspace of the center $\mathbf{U}^*(k) \in \text{Gr}(k, n)$.
 562 However, the solutions are often *nearly* nested. As a result, the vector, $\lambda^*(k-1)$, that provides
 563 the optimal value of the dual objective function for the problem on $\text{Gr}(k-1, n)$ can offer
 564 a good initialization for the dual subgradient algorithm used to find the GMEB center on
 565 $\text{Gr}(k, n)$, significantly reducing the total computation time needed to identify the optimal
 566 dimension, k^* . In [36] the authors also used a warm-starting strategy on a similar problem to
 567 improve the efficiency of a rank-adaptive matrix optimization scheme. Their proposed method
 568 alternates between greedy rank increase and smooth Riemannian optimization on fixed-rank
 569 manifolds, and they show that the strategy significantly improves the number of iterations and
 570 computational time to convergence.

571 The warm-start in this experiment is via the dual variables, but leads to a more efficient
 572 solution to the primal problem as well. By way of a baseline comparison, simple initializations
 573 of $\lambda^{(0)}(k)$ would be to randomly select the dual variables or to set all of the dual variables
 574 equal so that $\lambda^{(0)}(k) = [1/M, \dots, 1/M]^T$. For these experiments the latter strategy is chosen.
 575 The initial iterate for the primal variable when the dual variables are all equal is then the
 576 uniformly weighted extrinsic mean of the data, that is, $\mathbf{U}^{(0)}(k)$ is the dominant k -dimensional
 577 eigenspace of $\sum_{i=1}^M \lambda_i^{(0)} X_i X_i^T$. On $\text{Gr}(1, n)$, no warm-start initialization is possible because
 578 $\lambda^*(0)$ is undefined, so the algorithm is run using only the naive initialization. For $k =$
 579 $2, \dots, \max_i \{\dim(\mathbf{X}_i)\}$ Figure 6 illustrates the relative speed-up due to smart initialization by
 580 comparing the number of iterations needed to find a stationary point for different choices of
 581 the initial dual variable using each of the data models. Both data models are intentionally
 582 structured so that the extrinsic mean is not the center of the GMEB on $\text{Gr}(k^*, n)$. The naive
 583 initialization is indicated by the light orange box-and-whisker plots, while the warm-start is
 584 denoted with red. The black dots mark the mean number of iterations and the solid line is the
 585 median.

586 In Figure 6a the data has been generated using the asymmetrical nested ball model with

587 $M = 50$ points sampled from $\text{Gr}(p_i, 10)$ for $p_i \in \{4, 5, 6\}$ and an optimal dimension of
 588 $k^* = 4$. The warm start converged in less iterations than the naive initialization in 359 out
 589 of 500 possible trials. An experiment using data generated by sampling more densely from
 590 a randomly selected arc of a unit ball is displayed in Figure 6b. Here, $M = 300$ points were
 591 generated on $\text{Gr}(p_i, 10)$ with $p_i \in \{4, 5, 6\}$ where $k^* = 4$. In 415 out of 500 possible trials,
 592 the warm start converged in less iterations than the naive initialization.

593 **7.3. Experiment 3: Order-selection comparison.** The previous experiments demon-
 594 strated the effectiveness of the proposed approach for computing the subspace at the center of
 595 the GMEB in a noise-free scenario. However the end-goal is to find a central subspace *and*
 596 the optimal size to best represent the common dimensions in a collection of data. Adding
 597 noise to the subspaces makes it difficult to identify how many common dimensions exist, thus
 598 the third experiment compares the ability of the proposed order-selection rule to identify the
 599 optimal dimension of the common subspace with that of the technique from Santamaria *et*
 600 *al.* [28] as the difficulty of the task varies.

601 In many machine learning applications, extracting a low-rank common subspace from data
 602 is a pre-processing task and the rank is selected with little care. Heuristic solutions often focus
 603 on different methods for locating include the elbow of the scree plot, that is, computing the
 604 SVD of the concatenated data sets, finding the the singular values that represent the significant
 605 information, and keeping the dimensions corresponding to these singular values. This can be
 606 done with a variety of techniques such as the L-method [27], which estimates the elbow as
 607 the intersection of the two lines that minimize the root mean-squared error of the projection
 608 of the points in the of the scree plot onto the lines, the method of [40], which maximizes the
 609 profile log-likelihood under an independence assumption, and even just visually inspecting
 610 the scree plot to identify the first significant change in the first derivative [34]. To justify the
 611 need for a more principled way of selecting a subspace dimension, we additionally compare
 612 to the elbow of the scree plot using the L-method, and expect it to provide bad results. In the
 613 experiments this technique is denoted ‘‘SVD.’’

614 Figure 7 shows a comparison of order-selection rules for $M = 20$ points generated using
 615 the asymmetrical nested ball model from Section 6.1 with both generalizations. The data has
 616 $M_1 = 10$ points are sampled uniformly from the boundary of $\mathcal{B}_1(\mathbf{Z}_1) \subset \text{Gr}(10, n)$ and $M_2 = 10$
 617 points are sampled from the boundary of $\mathcal{B}_5(\mathbf{Z}_2) \subset \text{Gr}(15, n)$. Each of the points is then
 618 completed to a basis for a point on $\text{Gr}(p_i, n)$ for $p_i \in \{10, 11, \dots, 20\}$ and $n = 20, 30, \dots, 200$.
 619 Zero-mean Gaussian noise is added to each basis to create noisy data sets. The signal-to-noise
 620 ratio (SNR) of the data is the total power of the signal divided by the total power of the
 621 noise. In order to have the same SNR for each subspace despite differing dimensions, the
 622 noise variance per component is scaled by the number of subspace dimensions. Since X_i is
 623 an orthonormal basis for \mathbf{X}_i , the magnitude of each basis vector is 1. Thus the total power of
 624 signal subspace is k^* , and the SNR is computed as $\text{SNR} = 10 \log_{10}(k^*/\sigma_N^2)$, where σ_N^2 is the
 625 total variance of the noise. In this example the order of the common subspace is $k^* = 10$ and
 626 $\sigma_N^2 = 1.259$ meaning that the data has an SNR of 9dB.

627 Figure 7a shows the percentage of 100 Monte Carlo trials for which the proposed order-
 628 selection rule (purple dashed line with triangle markers), the method of Santamaria *et al.* [28]
 629 (pink solid line with circle markers), the hybrid method (turquoise dotted line with square
 630 markers), and the elbow point of the SVD (orange dash-dotted line with circle markers)
 631 were able to correctly identify the optimal order of the common subspace relative to the
 632 ambient dimension. Figure 7b shows the mean selected order, averaged across all trials. We
 633 can see that when the ambient dimension is small, all methods other than the SVD tend to
 634 overestimate the order of the common subspace. This is a result of the noise dimensions
 635 being relatively close in the low-dimensional spaces. The dimension of $\text{Gr}(k, n)$ is $k(n - k)$,

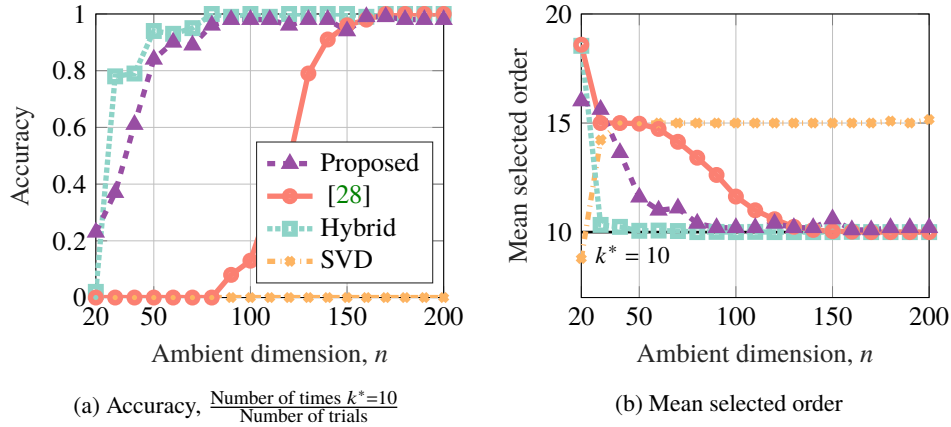


Fig. 7: Order-selection accuracy and mean selected order relative to the ambient dimension of the data from 100 Monte Carlo trials using the proposed order-selection rule (purple dashed line with triangle markers), the method of Santamaría *et al.* [28] (pink solid line with circle markers), the hybrid method (turquoise dotted line with square markers), and the elbow point of the SVD (orange dash-dotted line with circle markers). The data consists 20 points from $\bigsqcup_{p \in \mathcal{P}} \text{Gr}(p, n)$ for $\mathcal{P} = \{10, 11, \dots, 20\}$ and $n = 20, 30, \dots, 200$ with an SNR of 9 generated according to the model in Section 6.1.

636 so for $k \approx \max_i \{p_i\} \approx n$ all samples are very similar regardless of the data model. As the
 637 ambient dimension grows and the randomly selected dimensions become further apart on
 638 average, the proposed method and the hybrid method correctly select the order with a high
 639 degree of accuracy. The proposed method achieves slightly lower accuracy and has less stable
 640 performance than the hybrid method because $c_{\text{pen}}(k)$ can be significantly affected by even one
 641 subspace that is similar to $\mathbf{U}^{\perp}(k)$. However, this behavior is consistent with the assumption
 642 that every sample is valid and there are no outliers in the collection of data. As expected, [28]
 643 initially estimates the order as the dimension of the common subspace for the smaller ball and
 644 over-estimates the order as 15, while the two methods that rely on the minimum enclosing ball
 645 estimate the dimension of the common subspace for that support set. Predictably, the elbow
 646 point of the SVD has a very low accuracy regardless of the ambient dimension. In essence,
 647 this method is attempting to preserve all dimensions that are not pure noise.

648 Figure 8 shows a comparison using data from the second model, a ball that is sampled
 649 more densely from a random arc. For some $\mathbf{Z}_1 \in \text{Gr}(3, 100)$, $M_1 = 200$ points are sampled
 650 uniformly from $\mathcal{B}_{0.5}(\mathbf{Z}_1) \subset \text{Gr}(3, 100)$ and $M_2 = 25$ additional points are then sampled from
 651 a random arc on the same ball. No points were sampled from the interior of the ball. Each of
 652 these $M = 225$ subspaces is completed to basis for a point on $\text{Gr}(p_i, 100)$ for $p_i \in \{3, 4, 5\}$, and
 653 zero-mean Gaussian noise is added to each basis to create noisy data sets. In this experiment,
 654 the ambient dimension is fixed and we allow the SNR to vary from -5dB to 10dB .

655 With this data the optimal order of the common subspace is $k^* = 3$ and center of the ball is
 656 $\mathbf{U}^*(3) = \mathbf{Z}_1$. Figure 8a shows the percentage of 100 Monte Carlo trials for which the proposed
 657 order-selection rule (purple dashed line with triangle markers), the method of Santamaría
 658 *et al.* [28] (pink solid line with circle markers), the hybrid method (turquoise dotted line
 659 with square markers), and the elbow point of the SVD (orange dash-dotted line with circle
 660 markers) were able to correctly identify the optimal order of the common subspace relative

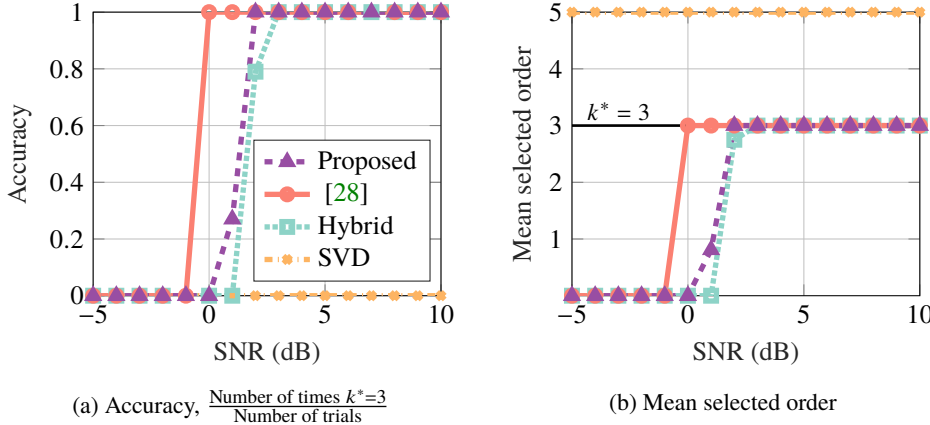


Fig. 8: Order-selection accuracy and mean selected order relative to the signal-to-noise ratio of the data (in dB) from 100 Monte Carlo trials using the proposed order-selection rule (purple dashed line with triangle markers), the method of Santamaría *et al.* [28] (pink solid line with circle markers), the hybrid method (turquoise dotted line with square markers), and the elbow point of the SVD (orange dash-dotted line with circle markers). The data consists 225 points from $\coprod_{p \in \mathcal{P}} \text{Gr}(p, 100)$ for $\mathcal{P} = \{3, 4, 5\}$ generated according to the model in Section 6.2.

661 to the signal-to-noise ratio. Figure 8b shows the mean selected order in the same trials. This
 662 experiment demonstrates the behavior of the different rules when all of the points are in the
 663 support of the minimum enclosing ball on $\text{Gr}(k^*, n)$. Each of the subspace averaging methods
 664 should theoretically select the same order in this experiment, because all of the points share
 665 the same number of dimensions and there is no ambiguity about the optimal solution. Thus
 666 even though the mean computed by [28] is not the same point as the center of the GMEB,
 667 they lead to the same estimated rank. We see that in this scenario, the behavior of the rules
 668 using ℓ_∞ -norm and the ℓ_2 -norm are similar with a sharp phase transition when the power of
 669 the signal and the power of the noise are almost equal, although the ℓ_2 -norm transitions to
 670 the correct order at a slightly higher noise power. This suggests that for situations where
 671 the data is free from outliers and the ℓ_∞ -mean is close to the ℓ_2 -mean, either technique will
 672 accurately estimate the number of common dimensions. The elbow point of the singular value
 673 decomposition fails to identify the common dimension in all trials.

674 Finally, in Figure 9 we see the ability of each method to identify when there is no subspace
 675 common to a collection of points. This is a valuable test because estimating $k^* = 0$ suggests
 676 that there is no information shared across all the data and that averaging the points is not
 677 an appropriate way to aggregate the information in the data. The data in this experiment
 678 consists of 50 subspaces chosen uniformly at random from $\text{Gr}(p_i, n)$ for $p_i \in \{3, 4, 5\}$ for
 679 $i = 1, \dots, 10$ with ambient dimensions $n = 5, 6, \dots, 15, 20, 25, \dots, 40$. The noise variance
 680 does not affect performance in this task because there is no signal so SNR undefined. In
 681 Figure 9a we see a similar phase transition to that of Figure 8. The hybrid method is able to
 682 achieve perfect accuracy for ambient dimensions greater than 10, while [28] and the proposed
 683 method transition shortly thereafter. The SVD fails every time, but that is to be expected in this
 684 scenario. The elbow point method computes two lines that minimize the residual for the scree
 685 plot, and chooses dimension as the index of the singular value just larger than the intersection
 686 of those lines. A line cannot be fit to zero points, so the method will not select $k^* = 0$ or $k^* = n$

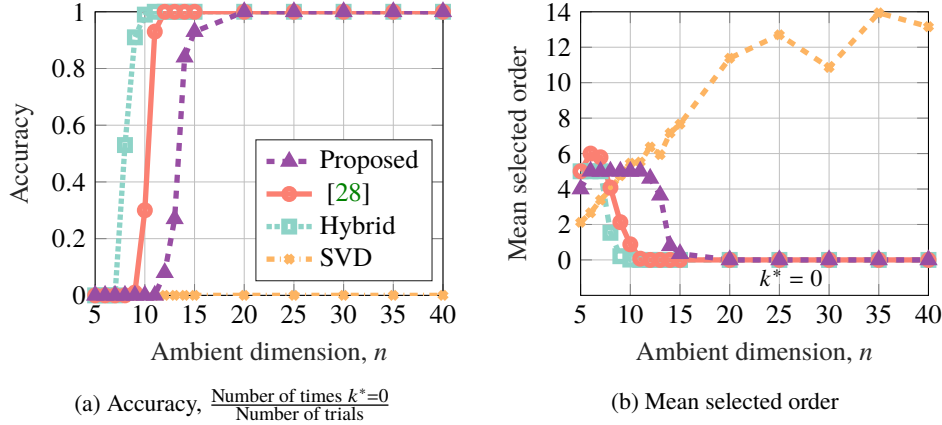


Fig. 9: Order-selection accuracy and mean selected order relative to the ambient dimension of the data when there is no common subspace. Results are from 100 Monte Carlo trials using the proposed order-selection rule (purple dashed line with triangle markers), the method of Santamaría *et al.* [28] (pink solid line with circle markers), the hybrid method (turquoise dotted line with square markers), and the elbow point of the SVD (orange dash-dotted line with circle markers). The data consists 50 points from $\coprod_{p \in \mathcal{P}} \text{Gr}(p, n)$ for $\mathcal{P} = \{3, 4, 5\}$ and $n = 5, 6, \dots, 15, 20, 25, \dots, 40$.

687 as a solution. However, in Figure 9b we see that the SVD is significantly overestimating the
688 dimension of the (non-existent) common subspace, so the poor performance is not an issue of
689 the method being unable to select 0 as the optimal dimension. When n is small the proposed
690 algorithm incorrectly identifies a relationship between the subspaces, but as the ambient
691 dimension grows the optimal order, $k^* = 0$, is selected with increasing accuracy. As noted
692 in discussion of Figure 7, the misidentifications in low dimensions are due to the minimum
693 similarity between the points and $\mathbf{U}^{*\perp}(k)$ being higher when $k \approx \max_i \{p_i\} \approx n$.

694 **8. Conclusions.** The recent trend of performing machine learning tasks on linear sub-
695 space data has created a need for flexible subspace averages, ones that can be computed
696 accurately and in a principled manner for subspaces of differing dimension. In response to
697 this need, we have proposed an algorithm to find the ℓ_∞ -center of mass using a subgradient
698 algorithm to solve the dual problem with respect to a point-to-set distance. We additionally
699 proposed a flexible data generation model to create subspaces of differing dimensions with
700 ground-truth for the GMEB that emulates realistic settings where an ℓ_∞ -average would be
701 appropriate. On this synthetic data, the proposed algorithm provides estimates of the GMEB
702 center with high accuracy. However, the high computational complexity means that an exist-
703 ing primal method can provide low-accuracy solutions more quickly for large data sets. One
704 direction for future expansion is to develop a core-set theory akin to that of [4] in order to es-
705 timate the GMEB on a subset of the data with theoretical accuracy guarantees. A related area
706 for further study is to develop an active-set approach for ℓ_∞ -averaging of mixed-dimensional
707 subspaces, à la John [14]. Active-set methods also attempt to minimize the cost function
708 over a subset of the data. However, the active-set approach looks for a subset of the data that
709 solves the original problem exactly, whereas the core-set technique computes error bounds
710 on the solution provided by *any* subset of a given size. One theoretical hurdle to achieving
711 an active-set method is a theorem on the minimum number of points required to define a

712 Grassmannian ball given a fixed Grassmann manifold and subspaces of differing dimensions.

713 Finally, we proposed a geometric order-fitting rule that estimates the best dimension for
 714 the common subspace. This rule fits the common dimensions of the subspaces in the support
 715 set of the minimum enclosing ball, which is appropriate for data where all subspace samples
 716 are assumed to be valid examples of the model of interest. We additionally implement a
 717 hybrid technique for estimating the dimension of the common subspace that modifies the
 718 order-selection rule of [28] for use with the ℓ_∞ -average. This hybrid method would not be
 719 possible for existing techniques that estimate the GMEB, because it uses the values of the
 720 dual variables as weights for an eigenvalue decomposition at each potential order. The hybrid
 721 approach outperforms the proposed technique and that of [28] when the ambient dimension
 722 is close to the subspace dimension of the data points.

723 A high-accuracy estimate of the GMEB center combined with an order-selection rule for
 724 the number of common dimensions results in a powerful technique for detecting and estimating
 725 similarity in a collection of subspaces. We anticipate that many practical applications will
 726 arise in the form of distributed large-scale problems, where the subspace averaging can be
 727 used for aggregation, for example the sparse subspace clustering of [1].

728 **Acknowledgments.** The authors would like to thank Emilie Renard for the stimulating
 729 discussions that improved the ideas presented here.

730

REFERENCES

- 731 [1] M. ABDOLALI, N. GILLIS, AND M. RAHMATI, *Scalable and robust sparse subspace clustering using randomized*
 732 *clustering and multilayer graphs*, Signal Processing, 163 (2019), pp. 166–180. 2, 23
- 733 [2] B. AFSARI, *Riemannian l^p center of mass: existence, uniqueness, and convexity*, Proceedings of the American
 734 Mathematical Society, 139 (2011), pp. 655–673. 1, 2
- 735 [3] M. ARNAUDON AND F. NIELSEN, *On approximating the Riemannian 1-center*, Computational Geometry, 46
 736 (2013), pp. 93–104. 1, 16
- 737 [4] M. BADOIU AND K. L. CLARKSON, *Smaller core-sets for balls*, in Proc. 14th ACM-SIAM Symposium on
 738 Discrete Algorithms, SIAM, 2003, pp. 801–802. 1, 23
- 739 [5] D. P. BERTSEKAS, *Nonlinear programming*, Journal of the Operational Research Society, 48 (1997), pp. 334–
 740 334. 7, 8
- 741 [6] A. BJÖRCK AND G. GOLUB, *Numerical methods for computing angles between linear subspaces*, Mathematics
 742 of Computation, 27 (1973), pp. 579–594. 3
- 743 [7] R. CHAKRABORTY AND B. C. VEMURI, *Recursive frechet mean computation on the grassmannian and its*
 744 *applications to computer vision*, in Proc. IEEE International Conference on Computer Vision, IEEE,
 745 2015, pp. 4229–4237. 2
- 746 [8] J.-M. CHANG, C. PETERSON, M. KIRBY, ET AL., *Feature patch illumination spaces and karcher compression for*
 747 *face recognition via grassmannians*, Advances in Pure Mathematics, 2 (2012), pp. 226–242. 2
- 748 [9] F. E. CURTIS, T. MITCHELL, AND M. L. OVERTON, *A BFGS-SQP method for nonsmooth, nonconvex, constrained*
 749 *optimization and its evaluation using relative minimization profiles*, Optimization Methods and Software,
 750 32 (2017), pp. 148–181. 8
- 751 [10] K. FISCHER AND B. GÄRTNER, *The smallest enclosing ball of balls: combinatorial structure and algorithms*,
 752 International Journal of Computational Geometry & Applications, 14 (2004), pp. 341–378. 1
- 753 [11] T. FRANZ, R. ZIMMERMANN, S. GÖRTZ, AND N. KARCHER, *Interpolation-based reduced-order modelling for*
 754 *steady transonic flows via manifold learning*, International Journal of Computational Fluid Dynamics, 28
 755 (2014), pp. 106–121. 2
- 756 [12] V. GARG, I. SANTAMARÍA, D. RAMÍREZ, AND L. L. SCHARF, *Subspace averaging and order determination for*
 757 *source enumeration*, IEEE Transactions on Signal Processing, 67 (2019), pp. 3028–3041. 13
- 758 [13] J.-B. HIRIART-URRUTY AND C. LEMARÉCHAL, *Convex Analysis and Minimization Algorithms I: Fundamentals*,
 759 vol. 305, Springer Science & Business Media, 2013. 8
- 760 [14] F. JOHN, *Extremum problems with inequalities as subsidiary conditions*, in Traces and Emergence of Nonlinear
 761 Programming, Springer, 2014, pp. 197–215. 23
- 762 [15] E. JURRUS, N. HODAS, N. BAKER, T. MARRINAN, AND M. D. HOOVER, *Adaptive visual sort and summary of*
 763 *micrographic images of nanoparticles for forensic analysis*, in Proc. IEEE Symposium on Technologies
 764 for Homeland Security, IEEE, 2016, pp. 1–6. 2
- 765 [16] A. V. KNYAZEV AND M. E. ARGENTATI, *Majorization for changes in angles between subspaces, Ritz values,*

- 766 *and graph Laplacian spectra*, SIAM Journal on Matrix Analysis and Applications, 29 (2006), pp. 15–32.
767 12
- 768 [17] P. KUMAR, J. S. MITCHELL, AND E. A. YILDIRIM, *Approximate minimum enclosing balls in high dimensions*
769 *using core-sets*, Journal of Experimental Algorithmics, 8 (2003), pp. 1–1. 1
- 770 [18] X. MA, M. KIRBY, C. PETERSON, AND L. SCHARF, *Self-organizing mappings on the grassmannian with*
771 *applications to data analysis in high dimensions*, Neural Computing and Applications, (2018), pp. 1–12.
772 2
- 773 [19] T. MARRINAN, J. BEVERIDGE, B. DRAPER, M. KIRBY, AND C. PETERSON, *Flag-based detection of weak gas*
774 *signatures in long-wave infrared hyperspectral image sequences.*, in Proc. SPIE Defense, Security, and
775 Sensing, International Society for Optics and Photonics, 2016. 1
- 776 [20] T. MARRINAN, B. DRAPER, J. R. BEVERIDGE, M. KIRBY, AND C. PETERSON, *Finding the subspace mean or*
777 *median to fit your need*, in Proc. IEEE Conference on Computer Vision and Pattern Recognition, IEEE,
778 2014, pp. 1082–1089. 16
- 779 [21] F. NIELSEN AND R. NOCK, *Approximating smallest enclosing balls with applications to machine learning*,
780 International Journal of Computational Geometry & Applications, 19 (2009), pp. 389–414. 1
- 781 [22] S. O'HARA AND B. A. DRAPER, *Scalable action recognition with a subspace forest*, in Proc. IEEE Conference
782 on Computer Vision and Pattern Recognition, IEEE, 2012, pp. 1210–1217. 2
- 783 [23] M. L. OVERTON AND R. S. WOMERSLEY, *Optimality conditions and duality theory for minimizing sums of the*
784 *largest eigenvalues of symmetric matrices*, Mathematical Programming, 62 (1993), pp. 321–357. 6, 7
- 785 [24] E. RENARD, P.-A. ABSIL, AND K. A. GALLIVAN, *Minimax center to extract a common subspace from multiple*
786 *datasets*, in Proc. 27th European Symposium on Artificial Neural Networks, Computational Intelligence
787 and Machine Learning, 2019. 1
- 788 [25] E. RENARD, K. A. GALLIVAN, AND P.-A. ABSIL, *A Grassmannian minimum enclosing ball approach for common*
789 *subspace extraction*, in Proc. International Conference on Latent Variable Analysis and Signal Separation,
790 Springer, 2018, pp. 69–78. 1, 5, 6, 8, 16, 17, 18
- 791 [26] Q. RENTMEESTERS, P. ABSIL, P. VAN DOOREN, K. GALLIVAN, AND A. SRIVASTAVA, *An efficient particle filtering*
792 *technique on the Grassmann manifold*, in Proc. IEEE International Conference on Acoustics, Speech and
793 Signal Processing, IEEE, 2010, pp. 3838–3841. 1, 16
- 794 [27] S. SALVADOR AND P. CHAN, *Determining the number of clusters/segments in hierarchical cluster-*
795 *ing/segmentation algorithms*, in Proc. 16th IEEE International Conference on Tools with Artificial
796 Intelligence, IEEE, 2004, pp. 576–584. 19
- 797 [28] I. SANTAMARÍA, L. L. SCHARF, C. PETERSON, M. KIRBY, AND J. FRANCO, *An order fitting rule for optimal*
798 *subspace averaging*, in Proc. IEEE Statistical Signal Processing Workshop, 2016, pp. 1–4. 2, 11, 13, 14,
799 15, 19, 20, 21, 22, 23
- 800 [29] A. N. SCHWICKERATH, *Linear models, signal detection, and the Grassmann manifold*, PhD thesis, Colorado
801 State University. Libraries, 2014. 4
- 802 [30] N. Z. SHOR, *Minimization Methods for Non-Differentiable Functions*, vol. 3, Springer Science & Business
803 Media, 2012. 6, 8
- 804 [31] K. SIM AND R. HARTLEY, *Removing outliers using the l_∞ norm*, in Proc. IEEE Conference on Computer Vision
805 and Pattern Recognition, IEEE, 2006. 2
- 806 [32] A. SRIVASTAVA AND E. KLASSEN, *Bayesian and geometric subspace tracking*, Advances in Applied Probability,
807 36 (2004), pp. 43–56. 2
- 808 [33] G. W. STEWART AND J.-G. SUN, *Matrix perturbation theory (computer science and scientific computing)*, (1990),
809 MR1061154, (1990). 2
- 810 [34] M. STEYVERS, *Multidimensional scaling*, Encyclopedia of Cognitive Science, (2006). 19
- 811 [35] P. TURAGA, A. VEERARAGHAVAN, A. SRIVASTAVA, AND R. CHELLAPPA, *Statistical computations on Grassmann*
812 *and Stiefel manifolds for image and video-based recognition*, IEEE Transactions on Pattern Analysis and
813 Machine Intelligence, 33 (2011), pp. 2273–2286. 1
- 814 [36] A. USCHMAJEV AND B. VANDEREYCKEN, *Greedy rank updates combined with riemannian descent methods*
815 *for low-rank optimization*, in 2015 International Conference on Sampling Theory and Applications
816 (SampTA), IEEE, 2015, pp. 420–424. 18
- 817 [37] Y.-C. WONG, *Differential geometry of Grassmann manifolds*, Proceedings of the National Academy of Sciences
818 of the United States of America, 57 (1967), p. 589. 3
- 819 [38] K. YE AND L.-H. LIM, *Schubert varieties and distances between subspaces of different dimensions*, SIAM
820 Journal on Matrix Analysis and Applications, 37 (2016), pp. 1176–1197. 2, 3, 4, 11
- 821 [39] E. A. YILDIRIM, *Two algorithms for the minimum enclosing ball problem*, SIAM Journal on Optimization, 19
822 (2008), pp. 1368–1391. 1
- 823 [40] M. ZHU AND A. GHODSI, *Automatic dimensionality selection from the scree plot via the use of profile likelihood*,
824 Computational Statistics & Data Analysis, 51 (2006), pp. 918–930. 19

Appendix A. GMEB dual subgradient algorithm.**Algorithm A.1** Algorithm to minimize Equation (3.8) with back-tracking line search

```

1: function GMEB( $\{\mathbf{X}_i\}_{i=1}^M, k, a, \eta, \zeta, \beta$ )
2:   input: Data:  $\{\mathbf{X}_i\}_{i=1}^M$ , Rank:  $k$ , Step size parameter:  $a$ , Stopping criteria:  $\eta$ , Step
   size threshold:  $\zeta$ , Growth parameter:  $\beta$ 
3:   output: Weights:  $\lambda^*$ , Minimax center:  $\mathbf{U}^*$ 
4:    $t \leftarrow 0$ 
5:    $\lambda^{(t)} \leftarrow [1/M, \dots, 1/M]^T \in \mathbb{R}^M$   $\triangleright \lambda^{(t)} \leftarrow \lambda^*(k-1)$  for warm-start
6:    $\mathbf{U}^{(t)} \leftarrow$  dominant  $k$  eigenvectors( $\sum_{i=1}^M \lambda_i^{(t)} \mathbf{X}_i \mathbf{X}_i^T$ )
7:    $\mathbf{g}^{(t)} \leftarrow -[d_{\text{Gr}(k,n)}(\mathbf{U}^{(t)}, \mathbf{X}_1), d_{\text{Gr}(k,n)}(\mathbf{U}^{(t)}, \mathbf{X}_2), \dots, d_{\text{Gr}(k,n)}(\mathbf{U}^{(t)}, \mathbf{X}_M)]^T$ 
8:    $f_{\text{primal}}(\mathbf{U}^{(t)}) \leftarrow \min_{i=1, \dots, M} \{-d_{\text{Gr}(k,n)}(\mathbf{U}^{(t)}, \mathbf{X}_i)\}$   $\triangleright$  Primal cost at iteration  $t$ 
9:    $f_{\text{dual}}(\lambda^{(t)}) \leftarrow \lambda^{(t)T} \mathbf{g}^{(t)}$   $\triangleright$  Dual cost at iteration  $t$ 
10:  while  $f_{\text{dual}}(\lambda^{(t)}) - f_{\text{primal}}(\mathbf{U}^{(t)}) > \eta$  and  $\max_{i=1, \dots, 10} \{f_{\text{dual}}(\lambda^{(t-i)}) - f_{\text{dual}}(\lambda^{(t)})\} > \eta$  do
11:     $t \leftarrow t + 1$ 
12:     $\alpha^{(t)} \leftarrow a/\sqrt{t}$ 
13:     $\lambda^{(t)} \leftarrow \lambda^{(t-1)} - \alpha^{(t)} \mathbf{g}^{(t-1)}$ ,  $\lambda^{(t)} \leftarrow \lambda^{(t)} / \|\lambda^{(t)}\|_1$ 
14:     $\mathbf{U}^{(t)} \leftarrow$  dominant  $k$  eigenvectors( $\sum_{i=1}^M \lambda_i^{(t)} \mathbf{X}_i \mathbf{X}_i^T$ )
15:     $\mathbf{g}^{(t)} \leftarrow -[d_{\text{Gr}(k,n)}(\mathbf{U}^{(t)}, \mathbf{X}_1), d_{\text{Gr}(k,n)}(\mathbf{U}^{(t)}, \mathbf{X}_2), \dots, d_{\text{Gr}(k,n)}(\mathbf{U}^{(t)}, \mathbf{X}_M)]^T$ 
16:     $\tilde{\alpha}^{(t)} \leftarrow \alpha^{(t)}$ 
17:     $\tilde{\lambda}^{(t)} \leftarrow \lambda^{(t)}$ 
18:     $f_{\text{dual}}(\tilde{\lambda}^{(t)}) \leftarrow \tilde{\lambda}^{(t)T} \mathbf{g}^{(t)}$ 
19:    while  $f_{\text{dual}}(\tilde{\lambda}^{(t)}) > f_{\text{dual}}(\lambda^{(t-1)})$  and  $\tilde{\alpha}^{(t)} > \zeta \alpha^{(t)}$  do  $\triangleright$  Back-tracking line search
20:       $a \leftarrow a/2$ 
21:       $\tilde{\alpha}^{(t)} \leftarrow a/\sqrt{t}$ 
22:       $\tilde{\lambda}^{(t)} \leftarrow \lambda^{(t-1)} - \tilde{\alpha}^{(t)} \mathbf{g}^{(t-1)}$ ,  $\tilde{\lambda}^{(t)} \leftarrow \tilde{\lambda}^{(t)} / \|\tilde{\lambda}^{(t)}\|_1$ 
23:       $\tilde{\mathbf{U}}^{(t)} \leftarrow$  dominant  $k$  eigenvectors( $\sum_{i=1}^M \tilde{\lambda}_i^{(t)} \mathbf{X}_i \mathbf{X}_i^T$ )
24:       $\tilde{\mathbf{g}}^{(t)} \leftarrow -[d_{\text{Gr}(k,n)}(\tilde{\mathbf{U}}^{(t)}, \mathbf{X}_1), d_{\text{Gr}(k,n)}(\tilde{\mathbf{U}}^{(t)}, \mathbf{X}_2), \dots, d_{\text{Gr}(k,n)}(\tilde{\mathbf{U}}^{(t)}, \mathbf{X}_M)]^T$ 
25:       $f_{\text{dual}}(\tilde{\lambda}^{(t)}) \leftarrow \tilde{\lambda}^{(t)T} \tilde{\mathbf{g}}^{(t)}$ 
26:      if  $f_{\text{dual}}(\tilde{\lambda}^{(t)}) \leq f_{\text{dual}}(\lambda^{(t-1)})$  then  $\triangleright$  Update variables if  $f_{\text{dual}}$  decreases
27:         $a \leftarrow \beta a$ 
28:         $\lambda^{(t)} \leftarrow \tilde{\lambda}^{(t)}$ 
29:         $\mathbf{U}^{(t)} \leftarrow \tilde{\mathbf{U}}^{(t)}$ 
30:         $\mathbf{g}^{(t)} \leftarrow \tilde{\mathbf{g}}^{(t)}$ 
31:       $f_{\text{primal}}(\mathbf{U}^{(t)}) \leftarrow \min_{i=1, \dots, M} \{-d_{\text{Gr}(k,n)}(\mathbf{U}^{(t)}, \mathbf{X}_i)\}$ 
32:       $f_{\text{dual}}(\lambda^{(t)}) \leftarrow \lambda^{(t)T} \mathbf{g}^{(t)}$ 
return  $\lambda^{(t)}, \mathbf{U}^{(t)}$ 

```
