

43 so we want to develop a reduced model that is also a function of the parameters.
 44 This is called parametric reduced order modeling (PROM), which is useful for design,
 45 control, optimization, uncertainty quantification, and inverse problems. Since most
 46 ROM methods are based on Petrov-Galerkin projection, which projects the model
 47 state space onto a low-dimensional subspace, one approach to PROM is to approximate
 48 the mapping from the parameters to such subspaces [8].

49 Subspace-valued mappings also arise in other areas of scientific computing. Active
 50 subspace methods [12] reduce the input space of a real-valued function to a low-
 51 dimensional subspace called the active subspace. For functional outputs, e.g. spatially
 52 varying fields or time series, such active subspaces become a function of space or time.
 53 Time-varying subspaces also arise in subspace tracking [13] and ROM [6].

54 **1.2. Previous methods.** A natural idea to solve this problem is to interpo-
 55 late subspaces as a function of the parameters. However, this is infeasible since the
 56 Grassmann manifold is not a vector space and linear combination is undefined. To
 57 circumvent this difficulty, [3] proposed a method that takes the interpolation to tangent
 58 spaces of the Grassmann manifold, which are vector spaces. It comes in three steps.
 59 Given a target parameter point, it chooses a few nearby parameter points and maps
 60 the associated subspaces to the tangent space of one of them via the Riemannian
 61 logarithm. Then the tangent vectors are interpolated as a function of the parameters,
 62 using any traditional interpolation method. Finally, the interpolated tangent vector
 63 is mapped back to the Grassmann manifold via the Riemannian exponential, which
 64 gives the predicted subspace. We will refer to this method as *subspace interpolation*.
 65 In fact, this three-step approach applies to any Riemannian manifold, as long as
 66 effective algorithms exist for the Riemannian exponential and logarithm maps [2, 4].
 67 This approach is extrinsic, i.e. referring to other sets and structures, which introduces
 68 distortions to the map.

69 Another type of method uses the Riemannian center of mass of weighted data
 70 points. The global or local Riemannian center of mass is the set of global or local
 71 minimizers of the sum of weighted squared Riemannian distances [1]. As before, the
 72 parameter-dependent weights can use any interpolation scheme such as splines [19] or
 73 Lagrange polynomials [40], both of which were introduced in the context of geodesic
 74 finite elements. Similarly, in the statistics literature, [37] proposed global and local
 75 regression models with predictors in a Euclidean space and random responses in a
 76 metric space. These methods are intrinsic, i.e. involving operations entirely on the
 77 manifold, so they avoid the limitations of mapping to a tangent space. However, their
 78 computation requires iterative algorithms for Riemannian optimization, and only local
 79 minimizers can be found. So far their uses are mostly for low-dimensional manifolds,
 80 with limited applications in PROM [35].

81 Zimmermann [49] reviewed interpolation methods on the Grassmann manifold and
 82 other matrix manifolds in the context of model reduction. More recently, he introduced
 83 Hermite interpolation of parameterized curves on Riemannian manifolds [50], which
 84 uses derivative data. All these methods are deterministic, while probabilistic methods
 85 for subspace approximation have not been explored in the literature.

86 **1.3. Contribution.** We propose a new Gaussian process (GP) model for the
 87 approximation of subspace-valued functions, which we call the Gaussian process sub-
 88 space (GPS) model. Instead of using differential geometric structures of the Grassmann
 89 manifold as in [3], the GPS uses matrix-variate Gaussian distributions on the Euclidean
 90 space to induce a probability model on the Grassmann manifold. Our method therefore
 91 yields a probabilistic prediction of the subspace response, with intrinsic characteri-

92 zation of its predictive mean and uncertainty. Specifically, the mean prediction is a
 93 k -subspace of the span of the observed subspaces, and the latter also covers most of
 94 the predictive uncertainty. This GP model is flexible and yet well-guided: it can be
 95 used with any correlation function on the parameter space, and the function form and
 96 hyperparameters can be optimized via specific model selection criteria.

97 The main advantages of our method are summarized as follows. (1) *Data efficient*:
 98 accurate prediction requires only a small sample size l , even when subspace dimension k
 99 and parameter dimension d are large. (2) *Computationally efficient*: its prediction cost
 100 does not depend on ambient dimension n , and thus it is suitable for large-scale problems
 101 and online computation. (3) *Flexible*: It is a flexible Bayesian nonparametric model
 102 that is robust against model misspecification. (4) It provides *uncertainty quantification*,
 103 which gives a probabilistic description of a predicted subspace.

104 In our observation, GPS is much more accurate than subspace interpolation [3],
 105 which is in turn much more accurate than other PROM methods [4, 36]. Such data
 106 efficiency can be attributed to two factors. First, our method is intrinsic, so it does not
 107 suffer from distortions due to pulling back the mapping to a tangent space. Second, it
 108 has clear rules for model selection, while the other methods are often subject to model
 109 misspecification, due to arbitrary choices of reference point, subsample points, and
 110 interpolation schemes.

111 **1.4. Related work.** The authors have worked on estimating functions whose
 112 domains or codomains are manifolds. For inputs on an unknown embedded submanifold,
 113 [44] proposed a GP model that attains the minimax-optimal convergence rate, without
 114 estimating the manifold. To allow for noisy inputs and better scalability, [21] first
 115 projects the input to random subspaces, and then applies a GP model. For inputs on
 116 a known embedded submanifold, [25] proposed an extrinsic GP, while [33] proposed
 117 an intrinsic GP, with heat kernel as the covariance function. For outputs on an
 118 embedded submanifold, [26] proposed a non-GP method, which applies an extrinsic
 119 local regression and then obtains manifold estimates via projection [46].

120 While our method extends GPs to mappings that take values in the Grassmann
 121 manifold, we are not the first to define GPs on Riemannian manifolds. Wrapped Gauss-
 122 ian process (WGP) regression [30] approximates mappings to a general Riemannian
 123 manifold, using distributions induced by Gaussian distributions on tangent spaces.
 124 However, this approach encounters problems when the manifold has a finite injectivity
 125 radius, as is the case for Grassmann manifolds. In particular, one cannot calculate the
 126 induced probability density function (PDF) on the manifold or the intrinsic mean. In
 127 contrast, our proposed approach produces analytic forms for predictive quantities that
 128 admit efficient computation, albeit restricted to Grassmann manifolds.

129 **1.5. Article structure and notations.** Section 2 provides basics of the algebra
 130 and statistics of some matrix manifolds. Section 3 presents the theoretical foundation
 131 of our GPS model, and Section 4 gives an algorithm for prediction. Section 5 discusses
 132 model selection for our model. Section 6 overviews ROM and discusses the use of
 133 GPS in PROM in the context of existing methods. Section 7 gives several numerical
 134 experiments: one to visualize the posterior process, and three to assess its accuracy in
 135 benchmark PROM problems. Section 8 concludes with a discussion on practical issues.
 136 Additional text is included in Supplementary Materials. An R package accompanying
 137 this paper is available at: <https://github.com/rudazhang/gpsr>.

138 *Notations.* Scalars are in lowercase, n, k, l, d ; vectors in boldface lowercase, $\mathbf{m}, \mathbf{x}_i, \boldsymbol{\theta}$;
 139 matrices in boldface uppercase, $\mathbf{M}, \mathbf{X}_i, \mathbf{K}_l$. Sets are in non-boldface uppercase, $\Theta, G_{k,n}$;
 140 subspaces in Fraktur script, $\mathfrak{X}, \mathfrak{M}$; equivalence classes in brackets, $[\mathbf{M}], [\mathbf{m}]$.

141 **2. Preliminaries.** Because we are building a probabilistic surrogate of subspace-
 142 valued mappings, it is helpful to review the algebra and statistics of the Grassmann
 143 manifold and some related matrix manifolds. For some basics of the algebra and
 144 differential geometry, see e.g. [7, 47]; for a reference on the statistics, see [11].

145 **2.1. Matrix manifolds.** Let $M_{n,k}$ be the set of all n -by- k real matrices, which
 146 can be identified as the Euclidean space $\mathbb{R}^{n \times k}$. The set of all full-rank n -by- k matrices
 147 is $M_{n,k}^* = \{\mathbf{M} \in M_{n,k} : \text{rank}(\mathbf{M}) = \min(n, k)\}$. When $k = n$, it coincides with the
 148 general linear group GL_n , which consists of full-rank order- n matrices.

149 The Stiefel manifold $V_{k,n}$ consists of all orthonormal k -frames in the Euclidean
 150 n -space: $V_{k,n} = \{\mathbf{X} \in M_{n,k}^* : \mathbf{X}^T \mathbf{X} = \mathbf{I}_k\}$, where $k \leq n$ and \mathbf{I}_k is the order- k
 151 identity matrix. The order of the subscripts is reversed by convention. When $k = n$,
 152 the Stiefel manifold coincides with the orthogonal group $O(n)$. Define projection
 153 $\pi : M_{n,k}^* \mapsto V_{k,n}$, such that for any $\mathbf{M} \in M_{n,k}^*$ with a thin singular value decomposition
 154 (SVD) $\mathbf{M} = \mathbf{V} \mathbf{\Sigma} \mathbf{U}^T$, $\mathbf{V} \in V_{k,n}$, $\mathbf{U} \in O(k)$, we have $\pi(\mathbf{M}) = \mathbf{V} \mathbf{U}^T$. Although the SVD
 155 is not unique, this mapping is uniquely defined.

156 The Grassmann manifold $G_{k,n}$ consists of all k -subspaces of the Euclidean n -space:
 157 $G_{k,n} = \{\text{span}(\mathbf{M}) : \mathbf{M} \in M_{n,k}^*\}$, where $\text{span}(\mathbf{M})$ denotes the subspace spanned by the
 158 columns of \mathbf{M} . Every element of $G_{k,n}$ is a subspace, which is often represented by a
 159 basis. For example, every $\mathbf{M} \in M_{n,k}^*$ represents $\mathfrak{M} = \text{span}(\mathbf{M})$, the column vectors of \mathbf{M}
 160 form a basis of \mathfrak{M} , and every element in its equivalence class $[\mathbf{M}] = \{\mathbf{M} \mathbf{A} : \mathbf{A} \in \text{GL}_k\}$
 161 represents \mathfrak{M} as well. We call \mathbf{M} a basis representation of \mathfrak{M} . In particular, every
 162 $\mathbf{X} \in V_{k,n}$ represents $\mathfrak{X} = \text{span}(\mathbf{X})$, and its column vectors form an orthonormal basis
 163 of \mathfrak{X} . We call \mathbf{X} a Stiefel representation of \mathfrak{X} .

164 The Grassmann manifold is often identified with the set of rank- k symmetric
 165 projection matrices $P_{k,n}$: let $\mathcal{S}(n)$ be the set of order- n symmetric matrices, define
 166 $P_{k,n} = \{\mathbf{P} \in \mathcal{S}(n) : \mathbf{P}^2 = \mathbf{P}, \text{rank}(\mathbf{P}) = k\}$. This identification is possible because
 167 $\text{span}(\cdot)$ is a bijection from $P_{k,n}$ to $G_{k,n}$. Given a Stiefel representation \mathbf{X} , a subspace \mathfrak{X}
 168 can thus be uniquely identified as $\mathbf{X} \mathbf{X}^T$. Due to this explicit identification, probability
 169 distributions on the Grassmann manifold can be induced through distributions on
 170 $P_{k,n}$, with the corresponding PDF being: $p : P_{k,n} \mapsto \mathbb{R}_{\geq 0}$, $\int_{P_{k,n}} p(\mathbf{P}) \mu(d\mathbf{P}) = 1$, where
 171 μ is the normalized invariant measure on $P_{k,n}$ under the group action of GL_n .

172 **2.2. Probability distributions.** Let $\mathcal{S}_+(n)$ be the set of order- n positive-definite
 173 matrices. Let $\mathbf{M} \in M_{n,k}$, $\mathbf{\Sigma}_1 \in \mathcal{S}_+(n)$, and $\mathbf{\Sigma}_2 \in \mathcal{S}_+(k)$. The n -by- k matrix-variate
 174 Gaussian distribution $N_{n,k}(\mathbf{M}; \mathbf{\Sigma}_1, \mathbf{\Sigma}_2)$ is the distribution of $\mathbf{Y} = \mathbf{\Sigma}_1^{1/2} \mathbf{Z} \mathbf{\Sigma}_2^{1/2} + \mathbf{M}$,
 175 where \mathbf{Z} is a random n -by- k matrix whose entries are independent standard Gaussian
 176 random variables. The vectorized matrix \mathbf{Y} is an (nk) -dimensional Gaussian random
 177 vector with a special form of covariance matrix: $\text{vec}(\mathbf{Y}) \sim N_{nk}(\text{vec}(\mathbf{M}), \mathbf{\Sigma}_2 \otimes \mathbf{\Sigma}_1)$,
 178 where $\text{vec}(\cdot)$ denotes vectorization of a matrix by stacking its columns, and \otimes is the
 179 Kronecker product.

180 The matrix angular central Gaussian distribution $\text{MACG}(\mathbf{\Sigma})$ is a probability
 181 distribution on $V_{k,n}$, with PDF $p(\mathbf{X}; \mathbf{\Sigma}) = z^{-1} |\mathbf{X}^T \mathbf{\Sigma}^{-1} \mathbf{X}|^{-n/2}$, where $|\cdot|$ denotes
 182 the determinant, normalizing constant $z = |\mathbf{\Sigma}|^{k/2}$, and parameter $\mathbf{\Sigma} \in \mathcal{S}_+(n)$. This
 183 parametric family contains the uniform distribution: since $p(\mathbf{X}; \mathbf{I}_n) = 1$, we have
 184 $\text{MACG}(\mathbf{I}_n) \sim \text{Uniform}$. The parameter of the MACG distribution is identified up to
 185 scaling: for all $\mathbf{\Sigma} \in \mathcal{S}_+(n)$ and $c \in \mathbb{R}_{>0}$, $\text{MACG}(\mathbf{\Sigma}) = \text{MACG}(c\mathbf{\Sigma})$.

186 Any probability distribution on $M_{n,k}$ or $V_{k,n}$ that is invariant under right-
 187 orthogonal transformation induces a probability distribution on $G_{k,n}$ [11, Thm 2.4.8]:
 188 let p be a PDF on $M_{n,k}$ such that $p(\mathbf{M}) = p(\mathbf{M} \mathbf{Q})$ for all $\mathbf{M} \in M_{n,k}$ and $\mathbf{Q} \in O(k)$, if

189 $\mathbf{M} \sim p$, let $\mathbf{X} = \pi(\mathbf{M}) \sim p_V$ and $\mathbf{X}\mathbf{X}^T \sim p_G$, then $p_V(\mathbf{X}) = p_V(\mathbf{X}\mathbf{Q})$ for all $\mathbf{Q} \in O(k)$,
 190 and $p_G(\mathbf{X}\mathbf{X}^T) = p_V(\mathbf{X})$. Because the MACG distribution on $V_{k,n}$ is invariant under
 191 right-orthogonal transformation, it defines a family of distributions on $G_{k,n}$ with the
 192 same PDF. We call it the MACG distribution on $G_{k,n}$.

193 These three distributions are related: let $\mathbf{M} \sim N_{n,k}(0; \boldsymbol{\Sigma}, \mathbf{I}_k)$ where $\boldsymbol{\Sigma} \in \mathcal{S}_+(n)$;
 194 let $\mathbf{X} = \pi(\mathbf{M})$, then $\mathbf{X} \sim \text{MACG}(\boldsymbol{\Sigma})$ and $\mathbf{X}\mathbf{X}^T \sim \text{MACG}(\boldsymbol{\Sigma})$. Due to this property,
 195 one can easily sample $\text{MACG}(\boldsymbol{\Sigma})$: generate $\mathbf{M} \sim N_{n,k}(0; \boldsymbol{\Sigma}, \mathbf{I}_k)$, and project it via π .

196 **3. Gaussian process subspace prediction.** We now present the proposed
 197 Gaussian Process Subspace (GPS) model. Because GP models take values in Euclidean
 198 spaces, they are not directly applicable to approximate subspace-valued mappings
 199 $f : \Theta \mapsto G_{k,n}$, where the codomain is the Grassmann manifold. Instead, we may
 200 find vector-valued mappings $\bar{f} : \Theta \mapsto \mathbb{R}^{nk}$ that are representations of f , in the sense
 201 that $f = \text{span} \circ \text{vec}^{-1} \circ \bar{f}$. Here, \circ denotes the composition of two mappings and
 202 $\text{vec}^{-1} : \mathbb{R}^{nk} \mapsto M_{n,k}$ denotes the ‘‘inverse’’ of $\text{vec}()$, that is, constructing a matrix
 203 columnwise from a vector. Such representations are not unique, and we denote the set
 204 of representations as $\bar{F} = \{\bar{f} : f = \text{span} \circ \text{vec}^{-1} \circ \bar{f}\}$. Now f can be identified with \bar{F} ,
 205 or equivalently, any distribution supported on \bar{F} .

206 GP models extend naturally to approximate distributions on a set of functions.
 207 Let $\mathfrak{X} = f(\boldsymbol{\theta})$ with a basis representation \mathbf{X} . Recall that \mathbf{X} has an equivalence class
 208 $[\mathbf{X}] = \{\mathbf{X}\mathbf{A} : \mathbf{A} \in \text{GL}_k\}$. Let $\mathbf{x} = \text{vec}(\mathbf{X})$, whose equivalence class can be written as
 209 $[\mathbf{x}] = \{\text{vec}(\mathbf{X}\mathbf{A}) : \mathbf{A} \in \text{GL}_k\}$. Assume that \bar{f} has a GP prior, we may assign equal
 210 likelihood to $[\mathbf{x}]$. We can then proceed to derive the posterior and the predictive
 211 distributions. In the following, we provide modeling details and analytical solutions
 212 for this approach.

213 **3.1. Model specification.** We start by specifying a prior for the representations.
 214 Without other information on f , an uninformative prior is for $f(\boldsymbol{\theta})$ to be uniformly
 215 distributed on $G_{k,n}$. We can achieve this by assigning $\bar{f}(\boldsymbol{\theta}) \sim N_{nk}(0, \mathbf{I}_{nk})$, the nk -
 216 dimensional standard Gaussian. To see this, let matrix $\mathbf{M} = \text{vec}^{-1}(\bar{f}(\boldsymbol{\theta}))$, then
 217 $\mathbf{M} \sim N_{n,k}(0; \mathbf{I}_n, \mathbf{I}_k)$ is a matrix-variate standard Gaussian; let subspace $\mathfrak{M} = \text{span}(\mathbf{M})$,
 218 then $\mathfrak{M} \sim \text{MACG}(\mathbf{I}_n) \sim \text{Uniform}$. We assign a correlation structure as follows. Let $k : \Theta \times \Theta \mapsto [-1, 1]$ be a correlation function, i.e. a positive definite kernel with $k(\boldsymbol{\theta}, \boldsymbol{\theta}) = 1$
 219 for all $\boldsymbol{\theta} \in \Theta$. For any finite collection of input points $\boldsymbol{\theta} = (\boldsymbol{\theta}_i)_{i=1}^l$, let $\mathbf{m}_i = \bar{f}(\boldsymbol{\theta}_i)$,
 220 and let \mathbf{K}_l be the order- l correlation matrix with entry $[\mathbf{K}_l]_{ij} = k(\boldsymbol{\theta}_i, \boldsymbol{\theta}_j)$. We assign
 221 the function values $\mathbf{m} = (\mathbf{m}_i)_{i=1}^l$ a prior joint distribution $\mathbf{m} \sim N_{nkl}(0, \mathbf{K}_l \otimes \mathbf{I}_{nk})$.
 222 Compactly, we can write this GP prior as $\bar{f} \sim \mathcal{GP}(0, k \otimes \mathbf{I}_{nk})$. This is the simplest
 223 covariance structure for \bar{f} .
 224

225 Without a likelihood function, this GP prior gives predictions as follows. Let $\boldsymbol{\theta}_*$
 226 be a target point and $\mathbf{m}_* = \bar{f}(\boldsymbol{\theta}_*)$. We have the prior joint distribution:

$$227 \quad (3.1) \quad (\mathbf{m}_*, \mathbf{m}) \sim N_{nk(l+1)}(0, \mathbf{K}_{l+1} \otimes \mathbf{I}_{nk})$$

228 where $\mathbf{K}_{l+1} = [1 \ \mathbf{k}_l^T; \mathbf{k}_l \ \mathbf{K}_l]$ and $\mathbf{k}_l = (k(\boldsymbol{\theta}_*, \boldsymbol{\theta}_i))_{i=1}^l$. If we write $\mathbf{K}_{22} = \mathbf{K}_l \otimes \mathbf{I}_{nk}$ and
 229 $\mathbf{K}_{12} = \mathbf{k}_l^T \otimes \mathbf{I}_{nk}$, by properties of multivariate Gaussian distributions, the conditional
 230 distribution of \mathbf{m}_* given \mathbf{m} can be written as:

$$231 \quad (3.2) \quad \begin{aligned} \mathbf{m}_* | \mathbf{m} &\sim N_{nk}(\mathbf{K}_{12}\mathbf{K}_{22}^{-1}\mathbf{m}, \mathbf{I}_{nk} - \mathbf{K}_{12}\mathbf{K}_{22}^{-1}\mathbf{K}_{12}^T) \\ &= N_{nk}\left(\sum_{i=1}^l [\mathbf{K}_l^{-1}\mathbf{k}_l]_i \mathbf{m}_i, (1 - \mathbf{k}_l^T \mathbf{K}_l^{-1} \mathbf{k}_l) \mathbf{I}_{nk}\right) \end{aligned}$$

232 We assign equal likelihood to the equivalence class of representations. Assume
 233 that we have function evaluations $\mathfrak{X}_i = f(\boldsymbol{\theta}_i)$ with Stiefel representations $\mathbf{X}_i \in V_{k,n}$.
 234 Let $\mathbf{x}_i = \text{vec}(\mathbf{X}_i)$ and $[\mathbf{x}_i] = \{\text{vec}(\mathbf{X}_i \mathbf{A}) : \mathbf{A} \in \text{GL}_k\}$. For $\mathbf{m}_i = \bar{f}(\boldsymbol{\theta}_i)$, the likelihood
 235 function gives:

$$236 \quad (3.3) \quad L(\mathbf{m}_i | \mathfrak{X}_i) = 1(\mathbf{m}_i \in [\mathbf{x}_i])$$

237 The posterior distribution of \mathbf{m} given observations $\mathfrak{X} = (\mathfrak{X}_i)_{i=1}^l$ is derived from
 238 the prior and the likelihood (3.3) via Bayes' rule:

$$239 \quad (3.4) \quad p(\mathbf{m} | \mathfrak{X}) \propto \exp \left\{ -\frac{1}{2} \mathbf{m}^T (\mathbf{K}_l \otimes \mathbf{I}_{nk})^{-1} \mathbf{m} \right\} \prod_{i=1}^l 1(\mathbf{m}_i \in [\mathbf{x}_i])$$

240 **3.2. Predictive distributions.** The predictive distribution of \mathbf{m}_* given observa-
 241 tions \mathfrak{X} is obtained by integrating the conditional distribution (3.2) over the posterior
 242 distribution (3.4). We summarize the result as follows:

243 **THEOREM 3.1.** *Let $\mathbf{X} = [\mathbf{X}_1 \cdots \mathbf{X}_l]$ be the matrix that combines \mathbf{X}_i by columns,
 244 and $\mathbb{X} = \text{diag}(\mathbf{X}_i)_{i=1}^l$ be the matrix with \mathbf{X}_i as diagonal blocks. Let $\varepsilon^2 = 1 - \mathbf{k}_l^T \mathbf{K}_l^{-1} \mathbf{k}_l$
 245 and $\mathbf{v} = \mathbf{K}_l^{-1} \mathbf{k}_l$. If $\mathbf{v} \in \mathbb{R}_{\neq 0}^l$,¹ let $\mathbf{D}_v = \text{diag}(\mathbf{v})$ and $\tilde{\mathbf{K}}_l = (\mathbf{D}_v \mathbf{K}_l \mathbf{D}_v)^{-1}$. The predictive
 246 distribution of \mathbf{m}_* is:*

$$247 \quad \mathbf{m}_* | \mathfrak{X} \sim N_{nk}(0, \mathbf{I}_k \otimes \boldsymbol{\Sigma})$$

$$248 \quad (3.5) \quad \boldsymbol{\Sigma} = \varepsilon^2 \mathbf{I}_n + \mathbf{X} [\mathbb{X}^T (\tilde{\mathbf{K}}_l \otimes \mathbf{I}_n) \mathbb{X}]^{-1} \mathbf{X}^T$$

250 The proof is quite lengthy and thus deferred to [section SM1](#). This theorem shows that,
 251 given observations: (1) the matrix $\mathbf{M}_* = \text{vec}^{-1}(\mathbf{m}_*)$ has a matrix-variate Gaussian
 252 distribution, $\mathbf{M}_* | \mathfrak{X} \sim N_{n,k}(0; \boldsymbol{\Sigma}, \mathbf{I}_k)$; and (2) the subspace $\mathfrak{M}_* = \text{span}(\mathbf{M}_*)$ has an
 253 MACG distribution, $\mathfrak{M}_* | \mathfrak{X} \sim \text{MACG}(\boldsymbol{\Sigma})$ (see [subsection 2.2](#)).

254 The predictive distributions admit an intuitive interpretation. Since $\boldsymbol{\Sigma}$ is positive
 255 semi-definite, there is an eigenvalue decomposition (EVD) $\boldsymbol{\Sigma} = \mathbf{Q} \text{diag}(\boldsymbol{\lambda}) \mathbf{Q}^T$, where
 256 $\boldsymbol{\lambda} \in \mathbb{R}_{\geq 0}^n$ are in decreasing order and $\mathbf{Q} \in O(n)$. Therefore we can simulate $\mathbf{M}_* | \mathfrak{X}$ as
 257 $\mathbf{M}_* = \boldsymbol{\Sigma}^{1/2} \mathbf{Z} = \mathbf{Q} \text{diag}(\boldsymbol{\lambda})^{1/2} \mathbf{Q}^T \mathbf{Z}$, where $\mathbf{Z} \in M_{n,k}$ is a random matrix of standard
 258 Gaussians. The columns of \mathbf{Z} are scaled by the square root of the eigenvalue in each
 259 eigenspace; therefore the range (i.e. column space) of \mathbf{M}_* is more likely to align with
 260 the top eigenspaces of $\boldsymbol{\Sigma}$. Recall that $\mathfrak{M}_* = \text{span}(\mathbf{M}_*)$. We have the following results.
 261 (1) The global Riemannian center of mass of $\mathfrak{M}_* | \mathfrak{X}$ is $\text{span}(\mathbf{V})$, where \mathbf{V} is the first k
 262 columns of \mathbf{Q} . (2) The uncertainty of $\mathfrak{M}_* | \mathfrak{X}$ is compactly described by the eigenvalues
 263 $\boldsymbol{\lambda}$: the larger an eigenvalue is, the more important is the associated eigenspace; and
 264 the mean prediction is more useful if $(\lambda_i)_{i=k+1}^n$ are small relative to $(\lambda_i)_{i=1}^k$.

265 A main feature of our GP model is that, while its construction involves the
 266 extrinsic Euclidean space \mathbb{R}^{nk} of basis representations of subspaces, its predictive
 267 distribution is intrinsic to the Grassmann manifold $G_{k,n}$. In particular, our model does
 268 not involve tangent spaces or the Riemannian exponential, and thus it is not subject to
 269 the distortions associated with applying local tangent approximations. Moreover, the
 270 function space explored by the GPS is much broader than the existing interpolation

¹The condition of no zero entry in $\mathbf{v}(\boldsymbol{\theta})$ holds almost everywhere in Θ , but it breaks most notably
 when predicting at sample points: $\mathbf{v}(\boldsymbol{\theta}_i) = \mathbf{e}_i$, which means $\boldsymbol{\Sigma}$ is singular at sample points and close
 to singular nearby. In these cases, one needs to be careful with matrix inversion in implementing the
 prediction algorithm in [section 4](#).

271 methods, so our model is more flexible and robust to model misspecification. Perhaps
 272 surprisingly, the GPS has closed-form expressions for its predictive distributions, which
 273 enables efficient computation for subspace prediction and uncertainty quantification.

274 While [Theorem 3.1](#) is concerned with point predictions on the Grassmann manifold,
 275 our GPS model also induces joint distributions on $G_{k,n}$ and can be used to generate
 276 random subspace-valued functions (see [section SM2](#)).

277 **4. Prediction algorithm.** From [Theorem 3.1](#) and the discussion thereafter we
 278 see that, to compute the predictive distribution, one needs the EVD of Σ . Even with
 279 Σ available, the EVD would cost $\mathcal{O}(n^3)$, which is intractable for large n . Here we give
 280 an efficient method to compute this.

281 **4.1. Efficient EVD of Σ .** Denote $\Pi = \mathbb{X}^T(\tilde{\mathbf{K}}_l \otimes \mathbf{I}_n)\mathbb{X}$ and $\tilde{\Sigma} = \mathbf{X}\Pi^{-1}\mathbf{X}^T$.
 282 We note that $\tilde{\mathbf{K}}_l, \Pi > 0$ and $\tilde{\Sigma} \geq 0$. Let $r = \text{rank}(\mathbf{X}) \leq \min(n, kl)$, then $\tilde{\Sigma}$ also
 283 has rank r and therefore r positive eigenvalues. From the form of $\tilde{\Sigma}$, we see that its
 284 top- r eigenvectors span the range of \mathbf{X} . Let $\mathbf{X} = \tilde{\mathbf{V}}\tilde{\mathbf{R}}\tilde{\mathbf{P}}^T$ be a rank-revealing QR
 285 decomposition, such that $\tilde{\mathbf{V}} \in V_{r,n}$ has r orthonormal columns, $\tilde{\mathbf{R}} \in M_{r,kl}$ is upper
 286 triangular, and $\tilde{\mathbf{P}}$ is a permutation matrix. Denote order- r matrix $\mathbf{S} = \tilde{\mathbf{V}}^T\tilde{\Sigma}\tilde{\mathbf{V}}$ and let
 287 $\mathbf{S} = \mathbf{Q}\text{diag}(\tilde{\lambda})\mathbf{Q}^T$ be an EVD where $\tilde{\lambda}$ is descending and $\mathbf{Q} \in O(r)$. Let $\mathbf{V} = \tilde{\mathbf{V}}\mathbf{Q}$ and
 288 let $\mathbf{Q} = (\mathbf{V}, \mathbf{V}_\perp) \in O(n)$ be an orthogonal completion. Let $\check{\lambda} = (\tilde{\lambda}, \mathbf{0}_{n-r})$ where $\mathbf{0}_{n-r}$
 289 is the vector of zeros with length $n - r$. Then we have an EVD: $\tilde{\Sigma} = \mathbf{Q}\text{diag}(\check{\lambda})\mathbf{Q}^T$.
 290 Because $\Sigma = \tilde{\Sigma} + \varepsilon^2\mathbf{I}_n$, we have an EVD of Σ :

$$291 \quad (4.1) \quad \Sigma = \mathbf{Q}\text{diag}(\check{\lambda} + \varepsilon^2\mathbf{1}_n)\mathbf{Q}^T$$

292 Here $\mathbf{1}_n$ is the vector of ones with length n . We see that, for a complete probabilistic
 293 prediction, we only need a rank-revealing QR of \mathbf{X} , an EVD of \mathbf{S} , and ε^2 . For the
 294 mean prediction, we only need the top- k eigenvectors of \mathbf{S} .

295 We can simplify the computation of \mathbf{S} as follows. Note that $\tilde{\mathbf{V}}^T\mathbf{X} = \tilde{\mathbf{R}}\tilde{\mathbf{P}}^T$ and
 296 $\tilde{\mathbf{P}}^{-1} = \tilde{\mathbf{P}}^T$. Because $\mathbf{S} = \tilde{\mathbf{V}}^T\tilde{\Sigma}\tilde{\mathbf{V}}$ and $\tilde{\Sigma} = \mathbf{X}\Pi^{-1}\mathbf{X}^T$, we have $\mathbf{S} = \tilde{\mathbf{R}}(\tilde{\mathbf{P}}\Pi\tilde{\mathbf{P}}^T)^{-1}\tilde{\mathbf{R}}^T$.
 297 Let order- (kl) Gram matrix $\square = \mathbf{X}^T\mathbf{X}$, which has a block matrix structure $\square =$
 298 $[\square_{ij}]_{i,j=1}^l$ with $\square_{ij} = \mathbf{X}_i^T\mathbf{X}_j$. Note that Π similarly has a block matrix structure
 299 $\Pi = [\Pi_{ij}]_{i,j=1}^l$ with $\Pi_{ij} = \tilde{k}_{ij}\square_{ij}$, where $\tilde{k}_{ij} = [\tilde{\mathbf{K}}_l]_{i,j}$. The construction of Π can be
 300 written in a compact form: $\Pi = \square \circ (\tilde{\mathbf{K}}_l \otimes \mathbf{J}_k)$, where \circ denotes the Hadamard product
 301 and $\mathbf{J}_k = \mathbf{1}_k\mathbf{1}_k^T$ is the order- k matrix of ones. Let $\tilde{\Pi} = \tilde{\mathbf{P}}\Pi\tilde{\mathbf{P}}^T$ and let $\tilde{\Pi} = \mathbf{L}\mathbf{L}^T$ be a
 302 Cholesky decomposition, where $\mathbf{L} \in M_{kl,kl}$ is lower triangular. Let $\tilde{\mathbf{L}} = \mathbf{L}^{-1}\tilde{\mathbf{R}}^T \in M_{kl,r}$
 303 by solving linear equations, which is also lower triangular, then we have $\mathbf{S} = \tilde{\mathbf{L}}^T\tilde{\mathbf{L}}$.

304 We formally describe the prediction procedure in two parts: [Algorithm 4.1](#) only
 305 needs to be done once, and [Algorithm 4.2](#) is needed for each prediction.

Algorithm 4.1 GPS: Preprocessing

Input: observation $\mathbf{X} = [\mathbf{X}_1 \ \cdots \ \mathbf{X}_l]$.

- 1: Compute Gram matrix: $\square \leftarrow \mathbf{X}^T\mathbf{X}$.
- 2: Rank-revealing QR: $\mathbf{X} = \tilde{\mathbf{V}}\tilde{\mathbf{R}}\tilde{\mathbf{P}}^T$.

Output: Gram matrix \square ; global basis $\tilde{\mathbf{V}}$; upper triangular $\tilde{\mathbf{R}}$; pivoting $\tilde{\mathbf{P}}$.

306 **4.2. Computational cost.** Here we analyze the computational cost of each step
 307 in floating point operations (flops), accurate up to the dominant term. In [Algorithm 4.1](#),
 308 line 1 takes nk^2l^2 flops; line 2 takes $\mathcal{O}(nklr)$ flops, and if $r \approx kl$, this requires about

Algorithm 4.2 GPS: Prediction

Require: correlation function $k(\cdot, \cdot)$; preprocessing output $(\square, \tilde{\mathbf{V}}, \tilde{\mathbf{R}}, \tilde{\mathbf{P}})$.

Input: sample $(\theta_i)_{i=1}^l$; target θ_* ; truncation size $t \in \{k, k+1, \dots, r\}$.

1: Construct correlation matrix and vector: $k_{ij} \leftarrow k(\theta_i, \theta_j)$, $k_i \leftarrow k(\theta_*, \theta_i)$.

2: Solve linear equations: $\mathbf{v} \leftarrow \text{solve}(\mathbf{K}, \mathbf{k})$, $\hat{\mathbf{K}} \leftarrow \text{solve}(\mathbf{K}, \text{diag}(\mathbf{v})^{-1})$.

3: Construct matrix: $\mathbf{\Pi} \leftarrow [\mathbf{\Pi}_{ij}]_{i,j=1}^l$, where $\mathbf{\Pi}_{ij} \leftarrow v_i^{-1} \hat{k}_{ij} \square_{ij}$.

4: Cholesky decomposition: $\tilde{\mathbf{P}}\mathbf{\Pi}\tilde{\mathbf{P}}^T = \mathbf{L}\mathbf{L}^T$.

5: Solve linear equations: $\tilde{\mathbf{L}} \leftarrow \text{solve}(\mathbf{L}, \tilde{\mathbf{R}}^T)$

6: Cross product: $\mathbf{S} \leftarrow \tilde{\mathbf{L}}^T \tilde{\mathbf{L}}$.

7: Truncated EVD: $\mathbf{S} = \hat{\mathbf{V}} \text{diag}(\hat{\boldsymbol{\lambda}}) \hat{\mathbf{V}}^T$, where $\hat{\boldsymbol{\lambda}}$ has length t .

8: Compute noise variance: $\varepsilon^2 \leftarrow 1 - \mathbf{k}^T \mathbf{v}$.

Output: principal directions $\mathbf{V} = \tilde{\mathbf{V}}\hat{\mathbf{V}}$; principal variances $\hat{\boldsymbol{\lambda}}$; noise variance ε^2 .

Note: May return $\tilde{\mathbf{V}}$ and $\hat{\mathbf{V}}$ instead of \mathbf{V} to avoid matrix multiplication.

309 $4nk^2l^2$ flops using the Householder QR with column pivoting [16]. In Algorithm 4.2,
 310 line 1 evaluates the correlation function $l^2/2$ times; line 2 takes $l^3/3$ flops for Cholesky
 311 decomposition, and $2l^3$ for forward and back substitution; line 3 takes $k^2l^2/2$ flops;
 312 line 4 takes $k^3l^3/3$ flops; line 5 takes $k^3l^3/3 - (kl-r)^3/3$ flops, due to the upper
 313 triangular structure in $\tilde{\mathbf{R}}$; line 6 takes $r^3/3 + (kl-r)r^2$ flops, due to the lower triangular
 314 structure in $\tilde{\mathbf{L}}$; line 7 takes $\mathcal{O}(r^2t)$ with classical or randomized algorithms [22]; and
 315 line 8 takes $2l$ flops. Note that \mathbf{K} and its Cholesky decomposition can be reused for
 316 future predictions. Overall, with $n > kl$ and assuming $r \approx kl$ and $t = k$, Algorithm 4.1
 317 gives an overhead cost of about $5nk^2l^2$ flops if we use the Householder QR with column
 318 pivoting, and Algorithm 4.2 gives a cost of about k^3l^3 flops per prediction.

319 An alternative version of Algorithm 4.2 is to conduct a truncated singular value
 320 decomposition: $\tilde{\mathbf{L}} = \hat{\mathbf{V}} \text{diag}(\hat{\boldsymbol{\sigma}}) \mathbf{W}^T$, and then return $\hat{\mathbf{V}}$ and $\hat{\boldsymbol{\lambda}} = \hat{\boldsymbol{\sigma}}^2$. Although this
 321 avoids the cross product in line 6 and thus saves about $k^3l^3/3$ flops, truncated
 322 SVD can take a significant amount of time and eliminate the saving. Theoretically,
 323 the truncated SVD takes $\mathcal{O}(rkl t)$ with classical algorithms, and $\mathcal{O}(rkl \log t)$ with
 324 randomized algorithms [22]. But in practice, the truncated SVD appears to be more
 325 costly than the truncated EVD. Since truncated SVD gives a less accurate result than
 326 truncated EVD, we consider Algorithm 4.2 as the reference version.

327 Note that the matrix multiplication $\mathbf{V} = \tilde{\mathbf{V}}\hat{\mathbf{V}}$ takes $2nrk$ flops for $t = k$, which
 328 would dominate the prediction cost if $n > kl^2/2$. However, this cost can be avoided
 329 if \mathbf{V} is not explicitly needed. In PROM problems, to compute an order- k reduced
 330 matrix $\mathbf{A}_k = \mathbf{V}^T \mathbf{A} \mathbf{V}$, one may precompute an order- r matrix $\mathbf{A}_r = \tilde{\mathbf{V}}^T \mathbf{A} \tilde{\mathbf{V}}$, and
 331 then compute $\mathbf{A}_k = \tilde{\mathbf{V}}^T \mathbf{A}_r \tilde{\mathbf{V}}$. Since \mathbf{A} is usually sparse, the cost of a matrix-vector
 332 multiplication $\mathbf{A} \mathbf{x}$ is usually $T_{\text{mult}} = \mathcal{O}(n)$. Then this approach has an overhead cost
 333 of $2nk^2l^2 + klT_{\text{mult}}$ flops, and only takes about $2k^3l^2$ flops per prediction.

334 **5. Model selection.** To make predictions with a GP model, we need to specify
 335 a covariance function; this is called model selection. Although the kernel $k(\cdot, \cdot)$ can be
 336 arbitrary, it is often specified in a form that depends on some hyperparameters [38,
 337 Ch. 4]. For example, the squared exponential (SE) kernel is:

$$338 \quad (5.1) \quad k(\boldsymbol{\theta}, \boldsymbol{\theta}'; \boldsymbol{\beta}) = \prod_{i=1}^d \exp \left[-\frac{(\theta_i - \theta'_i)^2}{2\beta_i^2} \right]$$

339 where length-scales $\boldsymbol{\beta} = (\beta_i)_{i=1}^d$ are the hyperparameters. GP models with the SE
 340 kernel are smooth, and the length-scales can be understood as characteristic distances
 341 along each parameter dimension before the function values become uncorrelated.

342 One can set the hyperparameters to optimize a certain criterion, see e.g. [38,
 343 Sec 5.4] and [41, Sec 3.3]. For GPS, we recommend minimizing the leave-one-out
 344 cross validation (LOOCV) predictive error, measured in Riemannian distances. (Other
 345 distances between subspaces may be used as well, but we choose Riemannian distance
 346 for concreteness.) In this section we analyze and give an algorithm to compute this
 347 criterion. Section SM3 provides a procedure to compute its gradient, and section SM4
 348 discusses some alternative criteria.

349 *A rule-of-thumb length-scale.* In our experience, the predictive performance of
 350 GPS is not very sensitive to hyperparameters, so one may use certain default values
 351 to trade accuracy for reduced computational cost. For the SE kernel, one may set the
 352 length-scales to $3d^{3/2}/l$ relative to the parameter ranges, and expect good predictions.

353 **5.1. LOOCV predictive error.** To measure predictive error, we need a score
 354 of dissimilarity for pairs of subspaces. There are many metrics defined on the Grass-
 355 mann manifold, see e.g. [45] for a list. Among them, the most commonly used is the
 356 Riemannian distance, which is the length of the shortest curves connecting two points
 357 in a Riemannian manifold. The Riemannian distance between subspaces $\mathfrak{X}, \mathfrak{Y} \in G_{k,n}$
 358 is the 2-norm of their principal angles, which can be computed as: [7]

$$359 \quad (5.2) \quad d_g(\mathfrak{X}, \mathfrak{Y}) = \|\arccos \boldsymbol{\sigma}(\mathbf{X}^T \mathbf{Y})\|$$

360 Here, $\mathbf{X}, \mathbf{Y} \in V_{k,n}$ are representations of the subspaces, and $\boldsymbol{\sigma}(\cdot)$ denotes the singular
 361 values of a matrix. Let \mathbf{V}_{-i} represent the mean prediction for target $\boldsymbol{\theta}_i$, using the
 362 remaining data points $(\boldsymbol{\theta}_j, \mathbf{X}_j)_{j \neq i}$. The LOOCV predictive error can be defined as:

$$363 \quad (5.3) \quad L_{\text{LOO}} = \sum_{i=1}^l d_g^2(\mathbf{X}_i, \mathbf{V}_{-i}) = \sum_{i=1}^l \sum_{j=1}^k (\arccos \sigma_j(\mathbf{X}_i^T \mathbf{V}_{-i}))^2$$

364 Here we use the sum of squared errors for its smoothness and, with a slight abuse of
 365 notation, we replace the subspaces with their Stiefel representations.

366 **5.2. Efficient computation of L_{LOO} .** To compute the LOOCV predictive error
 367 in (5.3), we need $\mathbf{X}_i^T \mathbf{V}_{-i}$. First we derive a form of \mathbf{V}_{-i} . Analogous to (3.5), for the
 368 leave-one-out prediction we have:

$$369 \quad (5.4) \quad \boldsymbol{\Sigma}_{-i} = \varepsilon_{-i}^2 \mathbf{I}_n + \mathbf{X}_{-i} [\mathbb{X}_{-i}^T (\tilde{\mathbf{K}}_{-i} \otimes \mathbf{I}_n) \mathbb{X}_{-i}]^{-1} \mathbf{X}_{-i}^T$$

370 Here, all the quantities are defined without the i -th observation. Similar to the analysis
 371 in subsection 4.1, denote $\boldsymbol{\Pi}_{-i} = \mathbb{X}_{-i}^T (\tilde{\mathbf{K}}_{-i} \otimes \mathbf{I}_n) \mathbb{X}_{-i}$ and $\tilde{\boldsymbol{\Sigma}}_{-i} = \mathbf{X}_{-i} (\boldsymbol{\Pi}_{-i})^{-1} \mathbf{X}_{-i}^T$. Let
 372 $r_{-i} = \text{rank}(\mathbf{X}_{-i})$, then the top- r_{-i} eigenvectors of $\tilde{\boldsymbol{\Sigma}}_{-i}$ span the range of \mathbf{X}_{-i} , which
 373 is a subset of the range of \mathbf{X} . Recall that $\mathbf{X} = \tilde{\mathbf{V}} \tilde{\mathbf{R}} \tilde{\mathbf{P}}^T$ is a rank-revealing QR. Let
 374 $\mathbf{S}_{-i} = \tilde{\mathbf{V}}^T \tilde{\boldsymbol{\Sigma}}_{-i} \tilde{\mathbf{V}}$ and let $\mathbf{S}_{-i} \approx \mathring{\mathbf{V}}_{-i} \text{diag}(\mathring{\boldsymbol{\lambda}}_{-i}) \mathring{\mathbf{V}}_{-i}^T$ be a rank- k truncated EVD, then
 375 $\mathring{\mathbf{V}} \mathring{\mathbf{V}}_{-i}$ are the top- k eigenvectors of $\tilde{\boldsymbol{\Sigma}}_{-i}$. Since \mathbf{V}_{-i} consists of the top- k eigenvectors
 376 of $\boldsymbol{\Sigma}_{-i}$ and $\boldsymbol{\Sigma}_{-i} = \varepsilon_{-i}^2 \mathbf{I}_n + \tilde{\boldsymbol{\Sigma}}_{-i}$, we have $\mathbf{V}_{-i} = \mathring{\mathbf{V}} \mathring{\mathbf{V}}_{-i}$.

377 To avoid big matrix multiplication, let $\tilde{\mathbf{C}} = \tilde{\mathbf{V}}^T \mathbf{X} = \tilde{\mathbf{R}} \tilde{\mathbf{P}}^T$, which has the form
 378 $\tilde{\mathbf{C}} = [\tilde{\mathbf{C}}_1 \ \cdots \ \tilde{\mathbf{C}}_l]$ where $\tilde{\mathbf{C}}_i = \tilde{\mathbf{V}}^T \mathbf{X}_i \in M_{r,k}$. We have $\mathbf{X}_i^T \mathbf{V}_{-i} = \mathbf{X}_i^T \mathring{\mathbf{V}} \mathring{\mathbf{V}}_{-i} = \tilde{\mathbf{C}}_i^T \mathring{\mathbf{V}}_{-i}$.
 379 Similarly, let $\tilde{\mathbf{C}}_{-i} = \tilde{\mathbf{V}}^T \mathbf{X}_{-i} = [\cdots \ \tilde{\mathbf{C}}_j \ \cdots]_{j \neq i}$, and we have $\mathbf{S}_{-i} = \tilde{\mathbf{C}}_{-i} (\boldsymbol{\Pi}_{-i})^{-1} \tilde{\mathbf{C}}_{-i}^T$.

380 We can express $\mathbf{\Pi}_{-i}$ using entries of \mathbf{K}^{-1} . Let $\mathbf{K}_{-i} = [k_{pq}]_{p,q \neq i}$ and $\mathbf{k}_{-i} = (k_{pi})_{p \neq i}$.
 381 Let $\bar{\mathbf{K}} = \mathbf{K}^{-1}$, $\bar{\mathbf{K}}_{-i} = [\bar{k}_{pq}]_{p,q \neq i}$, and $\bar{\mathbf{k}}_{-i} = (\bar{k}_{pi})_{p \neq i}$. We can write $\mathbf{v}_{-i} = (\mathbf{K}_{-i})^{-1} \mathbf{k}_{-i}$
 382 as $\mathbf{v}_{-i} = -\bar{\mathbf{k}}_{-i} / \bar{k}_{ii}$ and $(\mathbf{K}_{-i})^{-1} = \bar{\mathbf{K}}_{-i} - \bar{k}_{ii} \mathbf{v}_{-i} \mathbf{v}_{-i}^T$ (see for example [38, Sec. 5.4.2]).
 383 With $\tilde{\mathbf{K}}_{-i} = (\mathbf{D}_{\mathbf{v}_{-i}} \mathbf{K}_{-i} \mathbf{D}_{\mathbf{v}_{-i}})^{-1}$ and $\mathbf{D}_{\mathbf{v}_{-i}} = \text{diag}(\mathbf{v}_{-i})$, we have $\tilde{\mathbf{K}}_{-i} = \bar{k}_{ii}^{-1} \Delta_{-i}$
 384 where $\Delta_{-i} = [\bar{k}_{pq} \bar{k}_{ii} / (\bar{k}_{ip} \bar{k}_{iq}) - 1]_{p,q \neq i}$. Now we have $\mathbf{\Pi}_{-i} = \bar{k}_{ii}^{-1} \square_{-i} \circ (\Delta_{-i} \otimes \mathbf{J}_k)$.
 385 The computation of \mathbf{S}_{-i} follows subsection 4.1. Since we are only concerned with
 386 the eigenvectors of \mathbf{S}_{-i} , with a little abuse of notation, we redefine $\mathbf{\Pi}_{-i}$ without the
 387 term \bar{k}_{ii}^{-1} . We describe the overall procedure in Algorithm 5.1.

Algorithm 5.1 LOOCV Predictive Error

Require: correlation function k ; sample $(\boldsymbol{\theta}_i)_{i=1}^l$; preprocessing output $(\square, \tilde{\mathbf{C}} = \tilde{\mathbf{R}}\tilde{\mathbf{P}}^T)$.

Input: hyperparameters $\boldsymbol{\beta}$.

- 1: Construct inverse correlation matrix: $\bar{\mathbf{K}} \leftarrow \text{solve}(\mathbf{K})$, where $k_{ij} \leftarrow k(\boldsymbol{\theta}_i, \boldsymbol{\theta}_j; \boldsymbol{\beta})$.
- 2: **for** i in $1, \dots, l$ **do**
- 3: Construct: $\mathbf{\Pi} \leftarrow [\mathbf{\Pi}_{pq}]_{p,q \neq i}$, where $\mathbf{\Pi}_{pq} \leftarrow \delta_{pq} \square_{pq}$, $\delta_{pq} \leftarrow \bar{k}_{pq} \bar{k}_{ii} / (\bar{k}_{ip} \bar{k}_{iq}) - 1$.
- 4: Construct: $\mathbf{S} \leftarrow \tilde{\mathbf{L}}^T \tilde{\mathbf{L}}$, where $\mathbf{\Pi} = \mathbf{L}\mathbf{L}^T$, $\tilde{\mathbf{L}} \leftarrow \text{solve}(\mathbf{L}, \tilde{\mathbf{C}}_{-i}^T)$.
- 5: Truncated EVD: $\mathbf{S} = \mathring{\mathbf{V}} \text{diag}(\mathring{\boldsymbol{\lambda}}) \mathring{\mathbf{V}}^T$, where $\mathring{\boldsymbol{\lambda}}$ has length k .
- 6: Compute singular values: $\boldsymbol{\sigma} \leftarrow \boldsymbol{\sigma}(\tilde{\mathbf{C}}_{-i}^T \mathring{\mathbf{V}})$.
- 7: Compute squared error: $\epsilon_i \leftarrow \sum_{j=1}^k \arccos(\sigma_j)^2$
- 8: **end for**

Output: LOOCV predictive error $L_{\text{LOO}} = \sum_{i=1}^l \epsilon_i$.

388 **5.3. Computational cost.** In terms of computation, Algorithm 5.1 is approx-
 389 imately l repetitions of Algorithm 4.2, so it costs about $k^3 l^4$ flops per evaluation.
 390 This means that evaluating the LOOCV error takes about the same time as making l
 391 predictions. Because such evaluation needs to be repeated until numerical optimization
 392 converges, hyperparameter training may be a significant part of the overall cost. In
 393 practice, we recommend setting a very rough convergence threshold: for parameters
 394 with a range of one, a threshold of 0.01 is sufficient for the length-scale. If the problem
 395 has multiple parameters, they may be scaled into comparable ranges and share the
 396 same length-scale. If multiple hyperparameters are to be trained, gradient-based opti-
 397 mization methods (see section SM3) can be more efficient than just using the LOOCV
 398 error. To minimize the number of iterations, one may also set a restrictive range and,
 399 if applicable, a good initial value for the hyperparameters; for example, $\pm 30\%$ of the
 400 aforementioned rule-of-thumb length-scale, with initial value at the midpoint.

401 **6. Application in model reduction.** In this section, we review the general
 402 setup of model reduction, and compare the GPS with other methods for PROM.

403 **6.1. Reduced order modeling.** To simplify the narrative, consider a system of
 404 ordinary differential equations (ODEs) that is first-order, linear and time-invariant,
 405 with multiple input and output:

$$406 \quad (6.1) \quad \Sigma : \begin{cases} \mathbf{E}\dot{\mathbf{x}} = \mathbf{A}\mathbf{x} + \mathbf{B}\mathbf{u} \\ \mathbf{y} = \mathbf{C}\mathbf{x} \end{cases}$$

407 With system dimension n , input dimension p , and output dimension q , this system
 408 is defined by constant matrices $\mathbf{E}, \mathbf{A} \in M_{n,n}$, $\mathbf{B} \in M_{n,p}$, and $\mathbf{C} \in M_{q,n}$. The state \mathbf{x} ,
 409 input \mathbf{u} , and output \mathbf{y} are all functions of time, with dimension n , p , and q respectively.

410 We assume $\mathbf{x}(0) = \mathbf{0}$; any fixed initial condition \mathbf{x}_0 can be included in the input as an
 411 impulse $\mathbf{x}_0\delta(t)$. In general, the ODE system Σ may represent a physical or artificial
 412 system modeled by a PDE system, which is discretized in space, and linearized around
 413 a stationary trajectory. The system dimension n typically scales with the size of a
 414 spatial grid, and for a large-scale problem, usually we have $n > 10^5$.

415 Projection-based model reduction constructs a reduced-order model (ROM) as:

$$416 \quad (6.2) \quad \Sigma_r : \begin{cases} \mathbf{E}_r \dot{\mathbf{x}}_r = \mathbf{A}_r \mathbf{x}_r + \mathbf{B}_r \mathbf{u} \\ \mathbf{y}_r = \mathbf{C}_r \mathbf{x}_r \end{cases}$$

417 Let $\mathbf{V}, \mathbf{W} \in V_{k,n}$ be orthonormal bases of k -dimensional subspaces, the reduced system
 418 matrices are defined as $\mathbf{E}_r = \mathbf{W}^T \mathbf{E} \mathbf{V}$, $\mathbf{A}_r = \mathbf{W}^T \mathbf{A} \mathbf{V}$, $\mathbf{B}_r = \mathbf{W}^T \mathbf{B}$, and $\mathbf{C}_r = \mathbf{C} \mathbf{V}$.
 419 Therefore we have $\mathbf{E}_r, \mathbf{A}_r \in M_{k,k}$, $\mathbf{B}_r \in M_{k,p}$, and $\mathbf{C}_r \in M_{q,k}$. If the reduced bases
 420 \mathbf{V} and \mathbf{W} are the same, this framework is called the Galerkin projection; otherwise,
 421 it is called the Petrov-Galerkin projection. Usually we would want a reduced system
 422 dimension $k \leq 50$. Because simulation time and model storage scale at least linearly
 423 with system dimension, they are reduced by several orders of magnitude via ROM.

424 **6.2. Error measures.** To measure the error introduced by a ROM, one choice
 425 is the \mathcal{L}_2 state error for a given input. The \mathcal{L}_2 metric of square-integrable functions
 426 on the interval $[0, T]$, discretized into J parts of length δt , can be approximated as:

$$427 \quad (6.3) \quad \|\mathbf{x} - \hat{\mathbf{x}}\|_{\mathcal{L}_2}^2 = \int_0^T \|\mathbf{x}(t) - \hat{\mathbf{x}}(t)\|_2^2 dt \approx \sum_{i=1}^J \|\mathbf{x}(t_i) - \hat{\mathbf{x}}(t_i)\|_2^2 \delta t$$

428 Relative \mathcal{L}_2 state error is the \mathcal{L}_2 error of the state of a ROM, divided by the \mathcal{L}_2 norm
 429 of the state of the original system. With approximated state $\hat{\mathbf{x}} = \mathbf{V} \mathbf{x}_r$, we have:

$$430 \quad (6.4) \quad e(\mathbf{x}, \mathbf{x}_r)_{\mathcal{L}_2} = \frac{\|\mathbf{x} - \mathbf{V} \mathbf{x}_r\|_{\mathcal{L}_2}}{\|\mathbf{x}\|_{\mathcal{L}_2}}$$

431 Another error measure is the \mathcal{H}_2 metric, defined as the largest possible amplitude
 432 of the output error given any unit-energy input: with $\|\mathbf{y}\|_{\mathcal{L}_\infty} = \sup_{t \geq 0} \|\mathbf{y}(t)\|_\infty$,

$$433 \quad (6.5) \quad \|\Sigma - \Sigma_r\|_{\mathcal{H}_2} = \sup_{\mathbf{u} \in \mathcal{L}_2} \frac{\|\mathbf{y} - \mathbf{y}_r\|_{\mathcal{L}_\infty}}{\|\mathbf{u}\|_{\mathcal{L}_2}}$$

434 The \mathcal{H}_2 error of a ROM is, in a sense, more comprehensive than the \mathcal{L}_2 state error.
 435 Relative \mathcal{H}_2 error is the \mathcal{H}_2 error divided by the \mathcal{H}_2 norm of the original system:

$$436 \quad (6.6) \quad e(\Sigma, \Sigma_r)_{\mathcal{H}_2} = \frac{\|\Sigma - \Sigma_r\|_{\mathcal{H}_2}}{\|\Sigma\|_{\mathcal{H}_2}}$$

437 The \mathcal{H}_2 norms can be obtained analytically via the controllability Gramian, which
 438 can be computed by solving the Lyapunov equations [39].

439 **6.3. Methods for ROM.** To compute a reduced basis for the Galerkin projection,
 440 a widely-used classic method is called the proper orthogonal decomposition (POD),
 441 originally proposed for turbulent flow analysis by [28]. This method takes a collection of
 442 system states $\mathbf{x}(t_i)$ at discrete times $\{t_i\}_{i=1}^m$, called snapshots, which may be obtained
 443 via simulation or experimental measurements. Let \mathbf{X} be the matrix that stacks the
 444 snapshots as column vectors, then the POD basis \mathbf{V} corresponds to the left singular

445 vectors of \mathbf{X} associated with the largest k singular values. This means that the POD
 446 basis minimizes the \mathcal{L}_2 error of snapshot reconstruction, which is an appealing property
 447 of POD. Besides providing a reduced basis, POD also associates each basis vector with
 448 the corresponding singular value, which can be used to determine basis dimension k .
 449 For large-scale systems, the number of snapshots required is far less than the system
 450 dimension, and usually $m = \mathcal{O}(10^3)$.

451 Another class of ROM methods are interpolatory [5], which approximate the
 452 transfer function of the original system using rational interpolation. The transfer
 453 function of the system Σ is defined as $\mathbf{H}(s) = \mathbf{C}(s\mathbf{E} - \mathbf{A})^{-1}\mathbf{B}$. Here, $\mathbf{H} : \mathbb{C} \mapsto M_{q,p}(\mathbb{C})$
 454 is a complex matrix-valued function of a complex frequency variable. These methods
 455 interpolate the transfer function at an arbitrary number of points and up to an arbitrary
 456 number of derivatives along certain tangent directions. Among such methods, the
 457 iterative rational Krylov algorithm (IRKA) introduced by [20] has seen great success.
 458 It iteratively searches for an order- k rational function that approximates the transfer
 459 function, until it satisfies the tangential interpolation conditions. If IRKA converges,
 460 the converged point locally minimizes the \mathcal{H}_2 error in the space of order- k rational
 461 functions. IRKA constructs a ROM in state space via the two-sided Petrov-Galerkin
 462 projection, that is, the reduced bases \mathbf{V} and \mathbf{W} are different.

463 Besides POD and interpolatory methods, there are other ROM methods such
 464 as balanced truncation [31], most common in systems and control theory. There are
 465 effective ROM methods for systems more general than (6.1) as well, such as DEIM [9]
 466 for nonlinear systems and DMD [42] for black-box systems.

467 **6.4. Methods for PROM.** Our discussion so far assumes that the full model Σ
 468 in (6.1) is constant. In a more general class of problems, Σ is parametric, such that
 469 the system matrices \mathbf{E} , \mathbf{A} , \mathbf{B} , and \mathbf{C} depend on a set of parameters $\boldsymbol{\theta} \in \Theta \subset \mathbb{R}^d$. This
 470 dependency can be nonlinear in general, and the dimension d of the parameter space
 471 varies greatly with the problem. There are many methods for PROM, and we refer
 472 the readers to [8] for a comprehensive review.

473 One approach is to construct a single basis that works well for the entire parametric
 474 set of systems. For example, given local reduced bases $(\mathbf{V}_i)_{i=1}^l$ obtained for a sample
 475 of the parameter space, one can concatenate them into a global basis $\mathbf{V} = [\mathbf{V}_1 \cdots \mathbf{V}_l]$.
 476 However, this increases the dimension of the reduced subspace, and therefore the size
 477 and simulation time of the ROM.

478 Another approach is to consider a projection-based ROM method as a mapping
 479 that associates each parameter point with a reduced subspace. Given local reduced
 480 subspaces at a parameter sample, one may approximate this subspace-valued mapping
 481 and predict reduced subspaces at other parameter points. This fits the problem in
 482 Section 1, and includes subspace interpolation and our GPS method. Compared with
 483 using a global basis, this keeps the ROM small and often more reliable [3].

484 Instead of interpolating subspaces, [36] proposed a method that directly inter-
 485 polates the reduced models: it first applies a congruence transformation to the local
 486 reduced models, and then interpolates the model matrices element-wise. We will refer
 487 to this approach as *matrix interpolation*. Influenced by this work, [4] proposed a
 488 method that interpolates the transformed matrices on a relevant matrix manifold, e.g.
 489 the general linear group, in a procedure analogous to subspace interpolation. We will
 490 refer to this approach as *manifold interpolation*.

491 An idea bridging the global and local approaches is parameter domain partitioning:
 492 one can partition the parameter space into small regions and apply a PROM method
 493 within each. This idea has been adopted in many papers, see e.g. [2, 14].

494 **6.5. Comparing PROM methods.** Here we compare the GPS with other
 495 methods in model reduction, in terms of speed, accuracy, and property preservation.

496 **6.5.1. Speed vs. local bases.** Our method is typically much faster than methods
 497 for computing local reduced bases. Consider the computation of a local POD basis
 498 given m snapshots at one parameter point. The cost is dominated by a truncated SVD
 499 of the n -by- m snapshot matrix, which takes $\mathcal{O}(nmk)$ time. To compare the costs, take
 500 the rocket injector example in [29], where $n \approx 10^5$, $m = 10^3$, $k = 45$, $l = 30$. We have
 501 $(nmk)/(k^3l^3) \approx 1.83$. Considering the constant factor in truncated SVD, in this case
 502 our method is about an order of magnitude faster than computing a local POD basis.
 503 Because the cost of computing snapshots dominates the overall POD procedure, this
 504 implies a clear advantage in using our method to approximate local POD bases.

505 The cost of computing a pair of local IRKA bases is less straightforward to analyze
 506 [5]. Every iteration needs to solve $2k$ systems of linear equations, each with a different
 507 coefficient matrix of order n that cannot be reused across iterations. The number of
 508 iterations depends on the initial values provided to the algorithm, and the algorithm
 509 needs to be restarted if it does not converge after a predefined maximum number of
 510 iterations. Depending on the problem, IRKA can take longer than the POD procedure.

511 **6.5.2. Speed vs. interpolatory methods.** Subspace interpolation [3] uses the
 512 Riemannian exponential and logarithm of the Grassmann manifold, both involving a
 513 thin SVD of an n -by- k matrix, which scales with $\mathcal{O}(nk^2)$. Since its prediction does not
 514 have a special factorization structure (as the GPS does), it takes another $2nk^2 + kT_{\text{mult}}$
 515 flops to compute a reduced matrix, where T_{mult} denotes the cost of a matrix-vector
 516 multiplication. The prediction cost can be greatly reduced if the problem has only one
 517 parameter and one uses linear interpolation [43]. In general, the prediction scales with
 518 n and is slow for large-scale problems.

519 Matrix interpolation [36] and manifold interpolation [4] directly interpolate local
 520 ROMs so their prediction costs do not depend on n , and therefore they are considered
 521 as suitable for online computation.

522 In comparison, our algorithm turns the truncated EVD of the order- n matrix Σ
 523 into one of the order- kl matrix \mathbf{S} , and the prediction cost is instead dominated by the
 524 construction of \mathbf{S} , which is carried out efficiently via matrix decomposition and linear
 525 solvers. Thus, the prediction cost also does not depend on n .

526 **Table 1** compares the computational costs of these methods in detail. This table
 527 does not include the generation of reduced bases at a sample of the parameter space, a
 528 step required by all these methods. Generating a reduced basis can be computationally
 529 expensive depending on the ROM method in use, which limits the sample size l .

530 **6.5.3. Accuracy.** All three interpolation methods lack a clear rule for model
 531 selection, i.e. selecting the reference point, other interpolation points, and the interpo-
 532 lation scheme. This often leads to model misspecification which undermines accuracy.
 533 Moreover, interpolation on tangent spaces of Riemannian manifolds, such as subspace
 534 and manifold interpolation, are extrinsic to the underlying manifolds. As explained in
 535 [section SM6](#), when points further away from the reference point are used, the true
 536 mapping becomes more distorted on the tangent space and thus harder to approximate.
 537 A similar concern is addressed in [50] Sec. 3. Therefore, these methods cannot use
 538 more than a handful of points at a time, and have limited potential to extend to
 539 higher-dimensional parameter spaces.

540 Our method has specific model selection criteria which make it data efficient, so
 541 a small sample size is enough to give accurate results. Besides, the GPS is intrinsic

TABLE 1
Interpolatory methods for PROM: flop counts of the dominant terms.

	Preprocess	Subspace	ROM	Training	Reference
GPS	$5nk^2l^2$	k^3l^3	$2k^3l^2$	k^3l^4	this paper
Subspace-Int	$10nk^2l^2$	$8nk^2$	$2nk^2$	†	[3]
Matrix-Int	$6nk^2l^2$	-	$2k^2l$	†	[36]
Manifold-Int	nk^2l^2	-	$\mathcal{O}(k^3l)^*$	†	[4]

* Coefficient usually on the scale of 50 due to matrix logarithm / exponential, which can be numerically unstable [23].

† Optimal choice of reference ROM and interpolation scheme is an open problem.

542 to the Grassmann manifold, so it does not incur extra approximation error and its
 543 accuracy improves with sample size.

544 **6.5.4. Preservation of properties.** Another important issue in ROM is the
 545 preservation of system properties, such as stability, passivity, and contractivity. Al-
 546 though stability is not guaranteed for the reduced models generated by our method,
 547 from Section 7 we will see that, it is still observed in most cases, simply because our
 548 method can accurately approximate the subspace map of local ROMs.

549 7. Numerical experiments.

550 **7.1. Visualization of GP subspace prediction.** The simplest type of subspace-
 551 valued functions have the form $f : \mathbb{R} \mapsto G_{1,2}$, which maps a real number to a
 552 one-dimensional linear subspace in the plane. The Grassmann manifold $G_{1,2}$ can
 553 be identified as the unit circle, treating antipodal points as equivalent (Figure 1a).
 554 Therefore, such a function f can be plotted on the surface of a cylinder (Figure 1b),
 555 which helps us visualize the posterior process of the GPS model.

556 Specifically, let f be a covering map such that $f(\theta)$ is the subspace with angle
 557 $\alpha = \theta \bmod \pi$. This can be plotted as a double helix on the cylinder. To approximate
 558 this function with the proposed GPS model, suppose we observe sample points $\theta_i = c_i\pi$,
 559 where c_i are seven equal-distanced points between 0.2 and 1.8. For the correlation
 560 function k , we use the SE kernel, and set the length-scale β by minimizing the LOOCV
 561 predictive error. In this example, $\beta = 2.8 \approx 0.9\pi$. To visualize predictive uncertainty,
 562 we plot the 95% posterior predictive intervals (PI) from Theorem 3.1. We also include
 563 results from subspace interpolation for comparison. As suggested by the authors of [3],
 564 for every target parameter we use the nearest n_r sampled points for the interpolation
 565 (where $n_r = 3$ and 4 in Figure 1), among which the nearest sampled point is used as
 566 the reference point. We use Lagrange interpolation for the tangent vectors.

567 We see that, with only seven data points, the predictive mean function of GPS
 568 closely tracks the true function within the range of sampled parameter points. Fur-
 569 thermore, the uncertainties from our model also well-cover the truth: the posterior
 570 predictive intervals contain the true subspace values for all $\theta \in [0, 2\pi]$. Note that
 571 as the target point moves away from the sample points, the predictive distribution
 572 degenerates to the prior, the uniform distribution on $G_{1,2}$. Subspace interpolation,
 573 on the other hand, yields noticeably poorer predictions compared to GPS for both
 574 $n_r = 3$ and $n_r = 4$. As a deterministic interpolation approach, it also does not provide
 575 a quantification of interpolation uncertainty. This shows that, for this example, the
 576 proposed GPS model uses sample data more effectively to yield better predictions

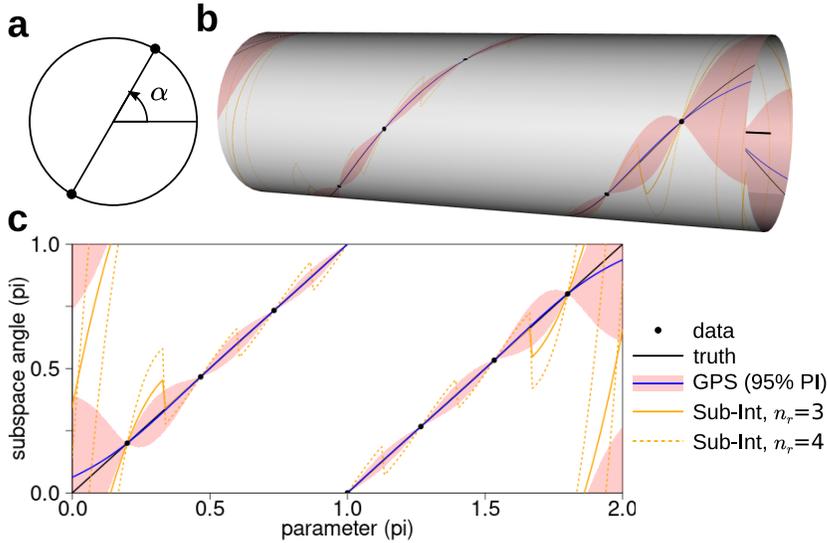


FIG. 1. Visualization of the GPS model. (a) Every 1d subspace in the plane can be uniquely identified by either a pair of antipodal points on a circle, or an angle $\alpha \in [0, \pi)$. (b) Posterior process of the GPS model on the surface of a cylinder. (c) Same as (b) but as a 2d plot. True function (black line), data (black points), GPS predictive mean (blue curve), 95% predictive interval (red shade). Orange curves are predictions from subspace interpolation: $n_r = 3$ (solid), $n_r = 4$ (dotted).

577 with uncertainty quantification.

578 **7.2. Anemometer: approximating local POD bases.** Here we consider
 579 a benchmark problem for PROM known as the anemometer [32], a type of micro-
 580 electromechanical system (MEMS) device that measures the flow speed of its sur-
 581 roundings. Such a device needs to be calibrated under different flow conditions for
 582 its temperature response. However, an accurate representation of the device needs to
 583 resolve the coupled fluid and thermodynamics, and can be very time-consuming to
 584 compute. It is therefore useful to apply PROM methods.

585 Specifically, a convection-diffusion equation is discretized into a linear ODE system
 586 as (6.1), with system dimension $n = 29,008$ and input and output dimensions $p = q = 1$.
 587 The matrix \mathbf{A} depends on one parameter $\theta \in [0, 1]$ representing fluid velocity and is
 588 not symmetric in general, while $\mathbf{E}, \mathbf{B}, \mathbf{C}$ are constants. The input map \mathbf{B} represents a
 589 heat source, and the output map \mathbf{C} gives the temperature difference of two nodes.

590 To build a parametric reduced-order model (PROM), we first construct local
 591 POD bases at a sample of the parameter space, and then use the mean prediction
 592 of GPS to estimate the reduced subspaces at other parameter points. As before, we
 593 use the SE kernel, with a length-scale that minimizes the LOOCV predictive error.
 594 The subspace-valued mappings being approximated in this problem have very high
 595 dimensional codomains: because the dimension of $G_{k,n}$ is $k(n - k)$, with $k = 20$ and
 596 $k = 40$, the manifold dimensions here are 579,760 and 1,158,720 respectively.

597 For comparison, we also estimate the reduced subspaces using subspace interpola-
 598 tion, with the same setup as in the visualization example. For manifold interpolation
 599 [4], we use the same setup for subspace interpolation. For matrix interpolation [36], we
 600 use the nearest sampled point as the reference point and, as suggested by the authors,
 601 we use linear interpolation for the reduced system matrices. We include results for

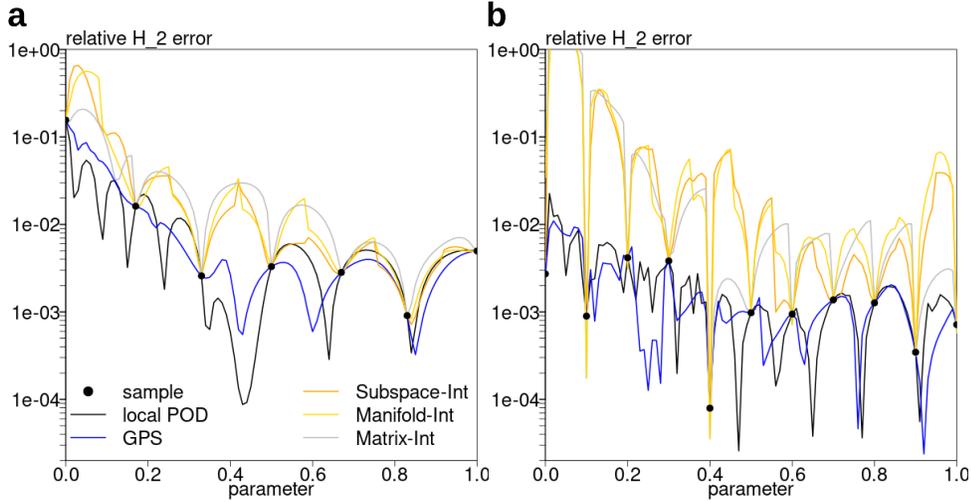


FIG. 2. Anemometer, relative \mathcal{H}_2 error: (a) $k = 20$; (b) $k = 40$. Training data shown as points. The \mathcal{H}_2 error curve of local POD is wiggly because it minimizes the \mathcal{L}_2 state error.

602 local POD bases as a reference level we would like to match.

603 **Figure 2a** shows the relative \mathcal{H}_2 errors using these methods, with subspace di-
 604 mension $k = 20$. Here we use a sample of seven equal-distanced points from 0 to 1.
 605 GPS uses a length-scale $\beta = 0.36$, selected via LOOCV. The results for subspace and
 606 manifold interpolation use $n_r = 3$; the results are similar for $n_r = 4$ or 5. We see that
 607 the three existing interpolation methods perform similarly, and the errors tend to
 608 blow up in between sample points. In comparison, the proposed GPS model yields
 609 much lower errors: the relative \mathcal{H}_2 error is comparable to that for the local POD (the
 610 reference level). Note that the goal here is not to perfectly match the error curve of
 611 local POD, but to keep the error as low as possible; in this sense, the GPS model
 612 appears to provide noticeable improvements over existing methods.

613 **Figure 2b** shows the results for $k = 40$. Here we use a sample of 11 equal-distanced
 614 points from 0 to 1. GPS uses a length-scale $\beta = 0.25$. Setup for the interpolation
 615 methods are unchanged. We see that, even with the increased sample size, all three
 616 interpolation methods fail to keep a low error level. While matrix interpolation
 617 occasionally does better than the other two, this is probably not generalizable due to
 618 the linear interpolation scheme. In comparison, our method again yields much lower
 619 errors, and maintains a similar level of accuracy as the local POD.

620 **Figure 3** shows the relative \mathcal{L}_2 state errors using these methods. Local POD is
 621 omitted from these plots since its relative \mathcal{L}_2 state error is practically zero. The error
 622 curves of the three interpolation methods are qualitatively similar, with subspace
 623 interpolation better than manifold interpolation, which is in turn better than matrix
 624 interpolation. In comparison, the GPS again yields much lower errors: for $k = 20$, the
 625 average error is about two orders of magnitude lower than that of subspace interpolation;
 626 for $k = 40$, it is about three orders of magnitude lower. This improvement can be
 627 attributed to the more flexible and intrinsic nature of the GPS model, which allows
 628 for more effective use of sample data.

629 Measured computation time for this problem is provided in [section SM5](#).

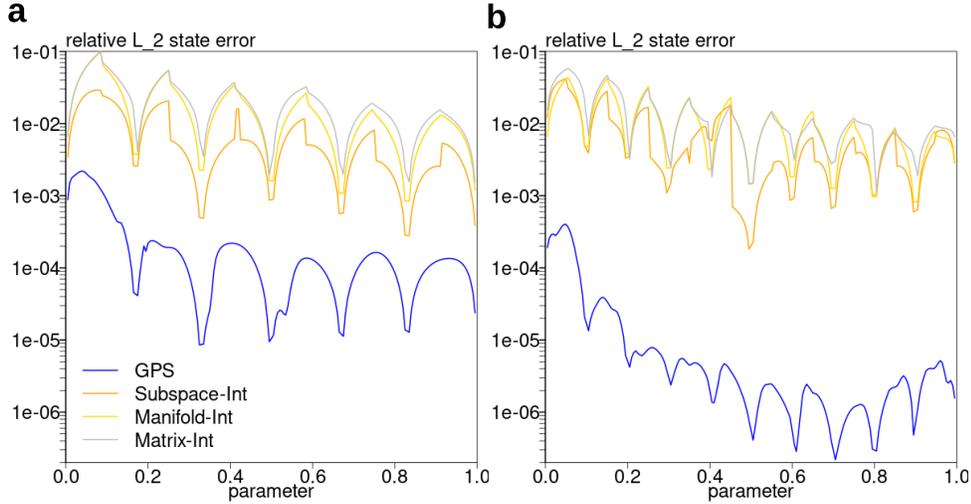


FIG. 3. Anemometer, relative \mathcal{L}_2 state error: (a) $k = 20$; (b) $k = 40$.

630 **7.3. Microthruster: approximating local IRKA bases.** Here we consider
 631 another benchmark problem for PROM known as the microthruster [34], an array of
 632 solid propellant microthrusters on a chip. To find an optimal design of array geometry
 633 and driving circuit, many simulations need to be carried out, which can be prohibitive
 634 with large-scale models. The use of PROM is therefore justified.

635 Specifically, the numerical model discretizes a heat transfer equation into a linear
 636 ODE system as (6.1), with system dimension $n = 4,257$, input dimension $p = 1$, and
 637 output dimension $q = 7$. The input \mathbf{B} represents the electrical circuit, and the output
 638 \mathbf{C} gives the temperature at seven nodes. The convection boundary conditions are
 639 parameterized into three parameters, each within the range $[1, 10^4]$, and affect the
 640 symmetric system matrix \mathbf{A} on the diagonal. To simplify comparison, we fix the three
 641 parameters to always be the same, and take the base-10 logarithm of their original
 642 values, so we have one parameter $\theta \in [0, 4]$.

643 For this problem, we use IRKA to construct reduced bases at the sample points.
 644 Because IRKA uses two different bases \mathbf{V} and \mathbf{W} , for a parametric system this means
 645 that each parameter is associated with a pair of subspaces, and we may construct
 646 a PROM by approximating a mapping for the form $(\mathfrak{A}, \mathfrak{B})(\theta)$. Since our proposed
 647 method only handles mappings that output one subspace, we proceed by modeling
 648 the pair of subspaces separately. This inevitably leaves some information in the data
 649 unused, and there may be methods that can improve upon this work-around. Setup
 650 for the interpolation methods are the same as in the anemometer example.

651 Figure 4 shows the relative \mathcal{H}_2 errors using these methods, with subspace dimension
 652 $k = 10$. Here we use a sample of six points: $\theta = 0.17, 0.94, 1.7, 2.47, 3.23, 4$. GPS uses a
 653 length-scale $\beta = 1.4$ for basis \mathbf{V} , and $\beta = 2.56$ for basis \mathbf{W} . The result for subspace
 654 interpolation uses $n_r = 3$; the other values of n_r give results with larger errors. We
 655 see that, while subspace interpolation matches the error curve of local IRKA (the
 656 reference level) quite well in some parts of the parameter space, its error blows up in
 657 an unsmooth region in between. These errors are noticeably larger for manifold and
 658 matrix interpolation, so we cropped them out of the plot. To contrast, the proposed
 659 GPS method instead tracks the local IRKA error curve smoothly across the parameter

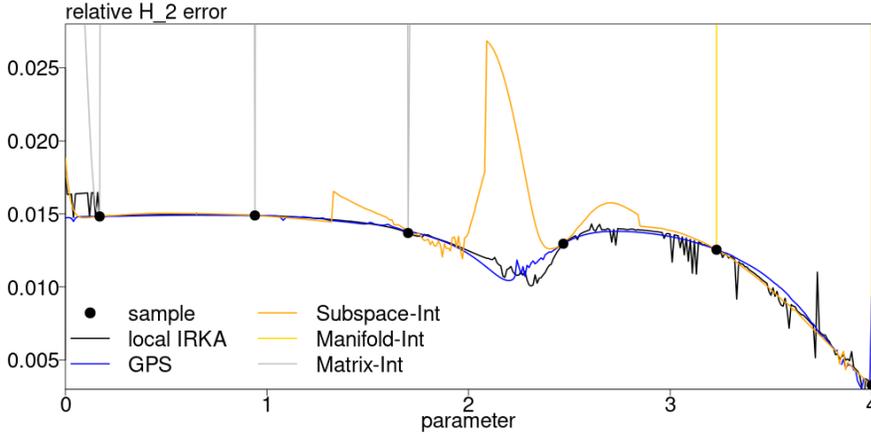


FIG. 4. *Microthruster*, relative \mathcal{H}_2 error. $k = 10$. Training data are shown as points. The error curve of local IRKA is more level than local POD in Figure 2 because it minimizes the \mathcal{H}_2 error.

660 space, yielding much lower errors than existing interpolation methods.

661 For this problem, many of the ROMs generated by manifold interpolation are
 662 complex-valued, due to the matrix logarithm that computes the tangent vectors.
 663 Moreover, many ROMs generated by manifold and matrix interpolation are unstable,
 664 which means that the \mathcal{H}_2 errors are infinite. Although our method and subspace
 665 interpolation do not guarantee the stability of reduced models, because they seem to
 666 accurately approximate the reduced subspaces, unstable ROMs appear less often. We
 667 discuss issues specific to approximating IRKA bases in section SM7.

668 **7.4. Anemometer: 3-parameter case.** To compare the methods in a PROM
 669 problem with multiple parameters, here we consider the three-parameter version
 670 of the anemometer [32]. The parameters include specific heat $c \in [0, 1]$, thermal
 671 conductivity $\kappa \in [1, 2]$, and fluid velocity $v \in [0.1, 2]$. The system matrices have the
 672 form $\mathbf{E} = \mathbf{E}_s + c\mathbf{E}_f$ and $\mathbf{A} = \mathbf{A}_{d,s} + \kappa\mathbf{A}_{d,f} + cv\mathbf{A}_c$, while \mathbf{B} and \mathbf{C} are constant. Other
 673 aspects of the problem are unchanged.

674 To sample the parameter space, we first use the maximin Latin hypercube sampling
 675 (LHS) to obtain a training set, and then use the sequential maximin design to obtain
 676 a testing set, see e.g. [17, Ch. 4]. Maximin LHS generates a random set of points that
 677 are spread out in the parameter space and well-distanced from each other. Sequential
 678 maximin design generates another with similar properties, but also well-distanced from
 679 the given training set.

680 The setup for the PROM methods remain unchanged from the 1-parameter case,
 681 except the interpolation scheme for the three interpolation methods. Since Lagrange
 682 and linear interpolations do not apply to multiple parameters, we use the radial
 683 basis function (RBF) method described in [2, p. 278]. Specifically, a multiquadric
 684 RBF is applied entrywise to interpolate the tangent vectors in subspace and manifold
 685 interpolation as well as the matrices in matrix interpolation. For subspace interpolation,
 686 horizontal projection is applied to maintain validity of the interpolated tangent vector.

687 Table 2 compares the mean relative \mathcal{H}_2 -errors, with training sample sizes $l = 14, 18,$
 688 or 21, and testing sample size 100. For each training sample, GPS uses a length-scale
 689 $\beta = 1.05, 0.85,$ or 0.7, respectively. Notice that, with subspace dimension $k = 20,$
 690 the mean relative \mathcal{H}_2 -error of local POD is about 5.5%, which is not particularly low. It

TABLE 2

Mean relative \mathcal{H}_2 -error for 3-parameter anemometer, $k = 20$, varying sample size.

	$l = 14$		$l = 18$		$l = 21$	
local POD	5.55%	(1)*	5.46%	(1)	5.69%	(1)
GPS	6.49%	(1.169)	5.80%	(1.062)	5.14%	(0.903)
Subspace-Int	8.34%	(1.503)	7.38%	(1.352)	6.19%	(1.148)
Manifold-Int	16.6%	(2.986)	13.8%	(2.524)	12.7%	(2.232)
Matrix-Int	49.7%	(8.962)	44.2%	(8.104)	45.5%	(8.003)

* Relative errors to local POD are shown in parentheses.

TABLE 3

Mean relative \mathcal{L}_2 state error for 3-parameter anemometer, $k = 20$, varying sample size.

	$l = 14$		$l = 18$		$l = 21$	
local POD	7.98e-13	(0)*	8.36e-13	(0)	8.77e-13	(0)
GPS	1.24e-2	(0.437)	6.42e-3	(0.273)	5.55e-3	(0.250)
Subspace-Int	2.85e-2	(1)	2.35e-2	(1)	2.22e-2	(1)

* Relative errors to subspace interpolation are shown in parentheses.

691 is clear that our method is able to maintain the error level of local POD with as few
 692 as 18 training points. In comparison, the error increase in subspace interpolation is
 693 several times higher in all cases. Manifold interpolation is much less accurate than the
 694 previous two methods, while matrix interpolation is the least accurate.

695 Similarly, Table 3 compares the mean relative \mathcal{L}_2 state errors. Manifold and matrix
 696 interpolation are excluded because they cannot reconstruct the state vector. Since
 697 local POD minimizes the \mathcal{L}_2 state error by construction, its error level is practically
 698 zero. With $l = 14$, our method has a relative error of about 1%, less than half that of
 699 subspace interpolation. This ratio drops as sample size gets larger. Overall, the GPS
 700 method is much more data efficient than subspace interpolation in this multi-parameter
 701 setting, again owing to its flexibility and intrinsic nature.

702 Measured computation time for this problem is provided in section SM5.

703 **8. Concluding remarks.** In this paper we propose a new GP model for proba-
 704 bilistic approximation of subspace-valued functions. A key application of this model
 705 is parametric reduced order modeling. We show that the GPS model gives accurate
 706 predictions even with small sample sizes, and because its prediction cost does not
 707 depend on system dimension n , it is typically faster than subspace interpolation in
 708 PROM problems. In the following, we discuss several topics on the use of the GPS.

709 *Prediction speed.* Since the prediction cost of our method is cubic in subspace
 710 dimension k and sample size l , it is best to keep them small for fast computation. To
 711 keep k small, one needs to choose a ROM method that is best suited for the relevant
 712 error measure. For example, POD is optimal in \mathcal{L}_2 error of snapshot reconstruction,
 713 while IRKA is locally optimal in \mathcal{H}_2 error. To keep l small, one needs to choose an
 714 efficient method for parameter sampling. One may consider adaptive sampling and
 715 sparse grids [8], or experimental design methods in statistics [41, 17].

716 *Handling higher-dimensional parameter spaces.* When parameter dimension d is
 717 large, even with the $l = 10d$ rule of thumb for GP models [27], l can quickly become

718 very large. Fortunately, there are some methods to cap the l^3 scaling. One approach
 719 is to use local approximate GP [18], where for each target point only a subsample
 720 of mostly nearby points are used in the prediction. Another approach is covariance
 721 tapering [15] or compactly supported kernels [24], where the kernel becomes zero
 722 beyond a certain distance, so that the covariance matrix is sparse and sparse matrix
 723 algorithms can be used to speed up computation. Both are in a similar spirit to
 724 parameter domain partitioning.

725 *Prediction uncertainty.* The uncertainty in subspace predictions, quantified by the
 726 eigenvalues of Σ , serves as a diagnostic tool for prediction confidence. It can also guide
 727 parameter sampling: one can put extra sample points in regions with high prediction
 728 uncertainty. This could lead to efficient adaptive sampling methods [17].

729 *Variation of subspace dimension.* In some cases it can be desirable to let k vary
 730 with the parameters, e.g. to attain a fixed ROM accuracy. Since the Grassmann
 731 manifold requires a fixed k , it acts as a Procrustean bed and limits all methods based
 732 on it, including the GPS. We recommend setting k to the highest value in the sample.

733 **Acknowledgments.** The authors thank Johan Guilleminot, Akil Narayan, Daniel
 734 Tartakovsky, and the anonymous reviewers for helpful feedback during this work. We
 735 also thank the INFORMS QSR community for their recognition.

736

REFERENCES

- 737 [1] B. AFSARI, *Riemannian L^p center of mass: Existence, uniqueness, and convexity*, Proc. Amer.
 738 Math. Soc., 139 (2011), pp. 655–673.
- 739 [2] D. AMSALLEM, *Interpolation on manifolds of CFD-based fluid and finite element-based structural*
 740 *reduced-order models for on-line aeroelastic predictions*, phdthesis, Stanford Univ., 2010.
- 741 [3] D. AMSALLEM AND C. FARHAT, *Interpolation method for adapting reduced-order models and*
 742 *application to aeroelasticity*, AIAA J., 46 (2008), pp. 1803–1813.
- 743 [4] D. AMSALLEM AND C. FARHAT, *An online method for interpolating linear parametric reduced-*
 744 *order models*, SIAM J. Sci. Comput., 33 (2011), pp. 2169–2198.
- 745 [5] A. C. ANTOULAS, C. A. BEATTIE, AND S. GUGERCIN, *Interpolatory Methods for Model Reduction*,
 746 SIAM, 2020.
- 747 [6] H. BABAEI, *An observation-driven time-dependent basis for a reduced description of transient*
 748 *stochastic systems*, Proc. R. Soc. A, 475 (2019).
- 749 [7] T. BENDOKAT, R. ZIMMERMANN, AND P. A. ABSIL, *A Grassmann manifold handbook: Basic*
 750 *geometry and computational aspects*. arXiv, 2020.
- 751 [8] P. BENNER, S. GUGERCIN, AND K. WILLCOX, *A survey of projection-based model reduction*
 752 *methods for parametric dynamical systems*, SIAM Rev., 57 (2015), pp. 483–531.
- 753 [9] S. CHATURANTABUT AND D. C. SORENSEN, *Nonlinear model reduction via discrete empirical*
 754 *interpolation*, SIAM J. Sci. Comput., 32 (2010), pp. 2737–2764.
- 755 [10] J. CHEN, S. MAK, V. R. JOSEPH, AND C. ZHANG, *Function-on-function kriging, with applications*
 756 *to three-dimensional printing of aortic tissues*, Technometrics, 63 (2021), pp. 384–395.
- 757 [11] Y. CHIKUSE, *Statistics on Special Manifolds*, Springer-Verlag, New York, 2003.
- 758 [12] P. G. CONSTANTINE, *Active Subspaces: Emerging Ideas for Dimension Reduction in Parameter*
 759 *Studies*, SIAM, Philadelphia, PA, 2015.
- 760 [13] A. EDELMAN, T. A. ARIAS, AND S. T. SMITH, *The geometry of algorithms with orthogonality*
 761 *constraints*, SIAM J. Matrix Anal. Appl., 20 (1998), pp. 303–353.
- 762 [14] J. L. EFTANG, A. T. PATERA, AND E. M. RÖNQUIST, *An "hp" certified reduced basis method for*
 763 *parametrized elliptic partial differential equations*, SIAM J. Sci. Comput., 32 (2010).
- 764 [15] R. FURRER, M. G. GENTON, AND D. NYCHKA, *Covariance tapering for interpolation of large*
 765 *spatial datasets*, J. Comput. Graph. Statist., 15 (2006), pp. 502–523.
- 766 [16] G. H. GOLUB AND C. F. VAN LOAN, *Matrix Computations*, Johns Hopkins University Press,
 767 Baltimore, MD, 2013.
- 768 [17] R. B. GRAMACY, *Surrogates: Gaussian Process Modeling, Design, and Optimization for the*
 769 *Applied Sciences*, Chapman and Hall/CRC, 2020.
- 770 [18] R. B. GRAMACY AND D. W. APLEY, *Local Gaussian process approximation for large computer*
 771 *experiments*, J. Comput. Graph. Statist., 24 (2015), pp. 561–578.

- 772 [19] P. GROHS, *Quasi-interpolation in Riemannian manifolds*, IMA J. Numer. Anal., 33 (2013),
773 pp. 849–874.
- 774 [20] S. GUGERCIN, A. C. ANTOULAS, AND C. BEATTIE, *\mathcal{H}_2 model reduction for large-scale linear*
775 *dynamical systems*, SIAM J. Matrix Anal. Appl., 30 (2008), pp. 609–638.
- 776 [21] R. GUHANIYOGI AND D. B. DUNSON, *Compressed Gaussian process for manifold regression*, J.
777 Mach. Learn. Res., 17 (2016), pp. 1–26.
- 778 [22] N. HALKO, P. G. MARTINSSON, AND J. A. TROPP, *Finding structure with randomness: Prob-*
779 *abilistic algorithms for constructing approximate matrix decompositions*, SIAM Rev., 53
780 (2011), pp. 217–288.
- 781 [23] N. J. HIGHAM, *The scaling and squaring method for the matrix exponential revisited*, SIAM J.
782 Matrix Anal. Appl., 26 (2005), pp. 1179–1193.
- 783 [24] C. G. KAUFMAN, D. BINGHAM, S. HABIB, K. HEITMANN, AND J. A. FRIEMAN, *Efficient emu-*
784 *lators of computer experiments using compactly supported correlation functions, with an*
785 *application to cosmology*, Ann. Appl. Stat., 5 (2011), pp. 2470–2492.
- 786 [25] L. LIN, N. MU, P. CHEUNG, AND D. DUNSON, *Extrinsic Gaussian processes for regression and*
787 *classification on manifolds*, Bayesian Anal., 14 (2019), pp. 887–906.
- 788 [26] L. LIN, B. S. THOMAS, H. ZHU, AND D. B. DUNSON, *Extrinsic local regression on manifold-valued*
789 *data*, J. Am. Stat. Assoc., 112 (2017), pp. 1261–1273.
- 790 [27] J. L. LOEPPKY, J. SACKS, AND W. J. WELCH, *Choosing the sample size of a computer experiment:*
791 *A practical guide*, Technometrics, 51 (2009), pp. 366–376.
- 792 [28] J. L. LUMLEY, *The structure of inhomogeneous turbulent flows*, Atmospheric Turbulence and
793 Radio Wave Propagation, (1967).
- 794 [29] S. MAK, C.-L. SUNG, X. WANG, S.-T. YEH, Y.-H. CHANG, V. R. JOSEPH, V. YANG, AND
795 C. F. J. WU, *An efficient surrogate model for emulation and physics extraction of large*
796 *eddy simulations*, J. Am. Stat. Assoc., 113 (2018), pp. 1443–1456.
- 797 [30] A. MALLASTO AND A. FERAGEN, *Wrapped Gaussian process regression on Riemannian manifolds*,
798 in Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, 2018.
- 799 [31] B. C. MOORE, *Principal component analysis in linear systems: Controllability, observability,*
800 *and model reduction*, IEEE Trans. Autom. Control, 26 (1981), pp. 17–32.
- 801 [32] MORWIKI COMMUNITY, *Anemometer*. Model Order Reduction Wiki, 2018.
- 802 [33] M. NIU, P. CHEUNG, L. LIN, Z. DAI, N. LAWRENCE, AND D. DUNSON, *Intrinsic Gaussian*
803 *processes on complex constrained domains*, J. R. Stat. Soc. Ser. B. Stat. Methodol., 81
804 (2019), pp. 603–627.
- 805 [34] OBERWOLFACH BENCHMARK COLLECTION, *Thermal model*. Model Order Reduction Wiki, 2018.
- 806 [35] M. OULGHELOU, C. ALLERY, AND R. MOSQUERA, *Parametric reduced order models based on a*
807 *riemannian barycentric interpolation*, Int. J. Numer. Methods Eng., (2021).
- 808 [36] H. PANZER, J. MOHRING, R. EID, AND B. LOHMANN, *Parametric model order reduction by*
809 *matrix interpolation*, at - Automatisierungstechnik, 58 (2010), pp. 475–484.
- 810 [37] A. PETERSEN AND H.-G. MÜLLER, *Fréchet regression for random objects with Euclidean predic-*
811 *tors*, Ann. Statist., 47 (2019), pp. 691–719.
- 812 [38] C. E. RASMUSSEN AND C. K. I. WILLIAMS, *Gaussian Processes for Machine Learning*, MIT
813 Press, Cambridge, MA, 2006.
- 814 [39] J. SAAK, M. KÖHLER, AND P. BENNER, *M-M.E.S.S. - the matrix equation sparse solver library*.
815 Zenodo, Apr. 2021.
- 816 [40] O. SANDER, *Geodesic finite elements of higher order*, IMA J. Numer. Anal., 36 (2016), pp. 238–
817 266.
- 818 [41] T. J. SANTNER, B. J. WILLIAMS, AND W. I. NOTZ, *The Design and Analysis of Computer*
819 *Experiments*, Springer, New York, NY, 2018.
- 820 [42] P. J. SCHMID, *Dynamic mode decomposition of numerical and experimental data*, J. Fluid
821 Mech., 656 (2010), pp. 5–28.
- 822 [43] N. T. SON, *A real time procedure for affinely dependent parametric model order reduction*
823 *using interpolation on Grassmann manifolds*, Int. J. Numer. Methods Eng., 93 (2013),
824 pp. 818–833.
- 825 [44] Y. YANG AND D. B. DUNSON, *Bayesian manifold regression*, Ann. Statist., 44 (2016), pp. 876–905.
- 826 [45] K. YE AND L.-H. LIM, *Schubert varieties and distances between subspaces of different dimensions*,
827 SIAM J. Matrix Anal. Appl., 37 (2016), pp. 1176–1197.
- 828 [46] R. ZHANG, *Newton retraction as approximate geodesics on submanifolds*. arXiv, June 2020.
- 829 [47] R. ZHANG AND R. GHANEM, *Normal-bundle bootstrap*, SIAM J. Math. Data Sci., 3 (2021),
830 pp. 573–592.
- 831 [48] R. ZHANG, P. WINGO, R. DURAN, K. ROSE, J. BAUER, AND R. GHANEM, *Environmental*
832 *economics and uncertainty: Review and a machine learning outlook*, in Oxford Research
833 Encyclopedia of Environmental Science, Oxford University Press, Aug. 2020.

- 834 [49] R. ZIMMERMANN, *Manifold interpolation and model reduction*. arXiv, 2019.
835 [50] R. ZIMMERMANN, *Hermite interpolation and data processing errors on Riemannian matrix*
836 *manifolds*, SIAM J. Sci. Comput., 42 (2020), pp. A2593–A2619.

1 **SUPPLEMENTARY MATERIALS: GAUSSIAN PROCESS SUBSPACE**
2 **PREDICTION**
3 **FOR MODEL REDUCTION***

4 RUDA ZHANG[†], SIMON MAK[‡], AND DAVID DUNSON[§]

5 **SM1. Proof of Theorem 3.1.** We see that the posterior $p(\mathbf{m}|\mathfrak{X})$ in (3.4) takes
6 positive values in $\prod_{i=1}^l [\mathbf{x}_i]$, where $[\mathbf{x}_i] = \{\text{vec}(\mathbf{X}_i \mathbf{A}) : \mathbf{A} \in \text{GL}_k\}$. Because GL_k is a
7 full-measure subset of $M_{k,k}$, we can replace $[\mathbf{x}_i]$ with $\{\text{vec}(\mathbf{X}_i \mathbf{A}) : \mathbf{A} \in M_{k,k}\}$ without
8 changing the posterior. Note that the latter equals $\mathfrak{X}_i^k = \prod_{j=1}^k \{\mathbf{X}_i \mathbf{c} : \mathbf{c} \in \mathbb{R}^k\}$, so the
9 support of the posterior can be written as: $S = \prod_{i=1}^l \mathfrak{X}_i^k$.

10 The predictive distribution of \mathbf{m}_* given observations \mathfrak{X} is obtained by integrating
11 the conditional distribution (3.2) over the posterior distribution (3.4), that is:

$$\circledast := p(\mathbf{m}_*|\mathfrak{X}) = \int_S p(\mathbf{m}_*|\mathbf{m}) p(\mathbf{m}|\mathfrak{X}) d\mathbf{m}$$

13 Every $\mathbf{m} \in S$ can be written as $\mathbf{m} = (\mathbf{m}_i)_{i=1}^l$, where $\mathbf{m}_i = \text{vec}(\mathbf{X}_i \mathbf{A}_i)$, $\mathbf{A}_i \in M_{k,k}$.
14 Let $\mathbf{m}_{:ji}$ and $\mathbf{a}_{:ji}$ be the j -th column of \mathbf{M}_i and \mathbf{A}_i respectively, then $\mathbf{m}_{:ji} = \mathbf{X}_i \mathbf{a}_{:ji}$.
15 Because \mathbf{X}_i has orthonormal columns, we have:

$$(SM1.1) \quad d\mathbf{m} = \prod_{i=1}^l d\mathbf{m}_i = \prod_{i=1}^l \prod_{j=1}^k d\mathbf{m}_{:ji} = \prod_{i=1}^l \prod_{j=1}^k d(\mathbf{X}_i \mathbf{a}_{:ji}) = \prod_{i=1}^l \prod_{j=1}^k d\mathbf{a}_{:ji} = d\mathbf{a}$$

17 Here, $\mathbf{a} = \text{vec}(\mathcal{A}) \in \mathbb{R}^{kkl}$ and \mathcal{A} is the $k \times k \times l$ array with frontal slices \mathbf{A}_i . Replacing
18 the integration domain S with \mathbb{R}^{kkl} , we have:

$$\circledast \propto \int_{\mathbb{R}^{kkl}} p(\mathbf{m}_*|\mathbf{m}) p(\mathbf{m}|\mathfrak{X}) d\mathbf{a}$$

20 Let $N(\mathbf{x}; \boldsymbol{\mu}, \boldsymbol{\Sigma})$ denote the value at \mathbf{x} of the Gaussian PDF with mean $\boldsymbol{\mu}$ and
21 covariance matrix $\boldsymbol{\Sigma}$. From (3.2), we have:

$$(SM1.2) \quad p(\mathbf{m}_*|\mathbf{m}) = N_{nk}(\mathbf{m}_*; \mathbf{K}_{12} \mathbf{K}_{22}^{-1} \mathbf{m}, \mathbf{I}_{nk} - \mathbf{K}_{12} \mathbf{K}_{22}^{-1} \mathbf{K}_{12}^T) \\ \propto \exp\left(-\frac{1}{2}(\mathbf{m}_* - \mathbf{K}_{12} \mathbf{K}_{22}^{-1} \mathbf{m})^T \mathbf{S}^\dagger (\mathbf{m}_* - \mathbf{K}_{12} \mathbf{K}_{22}^{-1} \mathbf{m})\right)$$

25 Here, $\mathbf{S} = \mathbf{I}_{nk} - \mathbf{K}_{12} \mathbf{K}_{22}^{-1} \mathbf{K}_{12}^T$ and \dagger denotes the Moore–Penrose inverse. In particular,
26 this allows \mathbf{S} to be singular. By computation rules of the Kronecker product:

$$(SM1.3) \quad \mathbf{S} = \mathbf{I}_{nk} - (\mathbf{k}_l^T \otimes \mathbf{I}_{nk})(\mathbf{K}_l \otimes \mathbf{I}_{nk})^{-1}(\mathbf{k}_l^T \otimes \mathbf{I}_{nk})^T \\ = \mathbf{I}_{nk} - (\mathbf{k}_l^T \mathbf{K}_l^{-1} \mathbf{k}_l) \otimes \mathbf{I}_{nk} = (1 - \mathbf{k}_l^T \mathbf{K}_l^{-1} \mathbf{k}_l) \mathbf{I}_{nk} = \varepsilon^2 \mathbf{I}_{nk}$$

*Submitted to the editors.

Funding: This work was supported by the National Science Foundation grants DMS-1638521 and CSSI-2004571, and the United States Office of Naval Research grant N00014-21-1-2510-01.

[†]Department of Mathematics, Duke University, Durham, NC 27710 USA (ruda.zhang@duke.edu)

[‡]Department of Statistical Science, Duke University, Durham, NC 27710 USA (sm769@duke.edu)

[§]Department of Mathematics and Department of Statistical Science, Duke University, Durham, NC 27710 USA (dunson@duke.edu)

30 We denote noise variance $\varepsilon^2 = 1 - \mathbf{k}_l^T \mathbf{K}_l^{-1} \mathbf{k}_l$. Note that $\varepsilon^2 \in [0, 1]$, and $\varepsilon^2 = 0$ if and
 31 only if $\boldsymbol{\theta}_* \in (\boldsymbol{\theta}_i)_{i=1}^l$. In the following, we assume $\varepsilon^2 \neq 0$. Since $\mathbf{m} = (\text{vec}(\mathbf{X}_i \mathbf{A}_i))_{i=1}^l$,
 32 by computation rules of the Kronecker product, we can write $\mathbf{K}_{12} \mathbf{K}_{22}^{-1} \mathbf{m}$ as:

$$33 \quad (\mathbf{k}_l^T \otimes \mathbf{I}_{nk})(\mathbf{K}_l \otimes \mathbf{I}_{nk})^{-1} \mathbf{m} = ((\mathbf{k}_l^T \mathbf{K}_l^{-1}) \otimes \mathbf{I}_{nk}) \mathbf{m} = \sum_{i=1}^l (\mathbf{k}_l^T \mathbf{K}_l^{-1} \mathbf{e}_i) \text{vec}(\mathbf{X}_i \mathbf{A}_i)$$

34 Since $\text{vec}()$ is a linear operator, we have

$$35 \quad \mathbf{K}_{12} \mathbf{K}_{22}^{-1} \mathbf{m} = \text{vec} \left(\sum_{i=1}^l (\mathbf{k}_l^T \mathbf{K}_l^{-1} \mathbf{e}_i) \mathbf{X}_i \mathbf{A}_i \right) = \text{vec} \left(\sum_{i=1}^l \mathbf{X}_i (\mathbf{k}_l^T \mathbf{K}_l^{-1} \mathbf{e}_i) \mathbf{I}_k \mathbf{A}_i \right)$$

36 Let $\mathbf{A}_{(13 \times 2)}$ be the matricization of \mathcal{A} by combining the matrices \mathbf{A}_i by rows. Recall
 37 that \mathbf{X} combines \mathbf{X}_i by columns, we have

$$38 \quad \mathbf{K}_{12} \mathbf{K}_{22}^{-1} \mathbf{m} = \text{vec}(\mathbf{X}(\text{diag}(\mathbf{K}_l^{-1} \mathbf{k}_l) \otimes \mathbf{I}_k) \mathbf{A}_{(13 \times 2)}) = \text{vec}(\tilde{\mathbf{X}} \mathbf{A}_{(13 \times 2)})$$

39 Here, $\tilde{\mathbf{X}} = \mathbf{X}(\text{diag}(\mathbf{K}_l^{-1} \mathbf{k}_l) \otimes \mathbf{I}_k)$. By the ‘‘vec trick’’ of the Kronecker product, we have

$$40 \quad (\text{SM1.4}) \quad \mathbf{K}_{12} \mathbf{K}_{22}^{-1} \mathbf{m} = (\mathbf{I}_k \otimes \tilde{\mathbf{X}}) \text{vec}(\mathbf{A}_{(13 \times 2)}) = (\mathbf{I}_k \otimes \tilde{\mathbf{X}}) \mathbf{a}_{(13 \times 2)}$$

41 Here $\mathbf{a}_{(13 \times 2)} = \text{vec}(\mathbf{A}_{(13 \times 2)})$. Substituting (SM1.3) and (SM1.4) into (SM1.2), we have

$$42 \quad (\text{SM1.5}) \quad p(\mathbf{m}_* | \mathbf{m}) \propto \exp \left(-\frac{1}{2} \varepsilon^{-2} \|\mathbf{m}_* - (\mathbf{I}_k \otimes \tilde{\mathbf{X}}) \mathbf{a}_{(13 \times 2)}\|^2 \right)$$

43 From (3.4), $p(\mathbf{m} | \mathfrak{X}) \propto \exp\{-\frac{1}{2} \mathbf{m}^T (\mathbf{K}_l^{-1} \otimes \mathbf{I}_{nk}) \mathbf{m}\}$, where $\mathbf{m} = (\text{vec}(\mathbf{X}_i \mathbf{A}_i))_{i=1}^l$.
 44 Note that matrix inverse and the Kronecker product commute. Let $\bar{k}_{ij} = [\mathbf{K}_l^{-1}]_{ij}$.
 45 Expand the Kronecker product and use properties of the trace, we have

$$46 \quad \mathbf{m}^T (\mathbf{K}_l^{-1} \otimes \mathbf{I}_{nk}) \mathbf{m}$$

$$47 \quad = \sum_{i=1}^l \sum_{j=1}^l \bar{k}_{ij} \text{vec}(\mathbf{X}_i \mathbf{A}_i)^T \text{vec}(\mathbf{X}_j \mathbf{A}_j) = \sum_{i=1}^l \sum_{j=1}^l \bar{k}_{ij} \text{tr}((\mathbf{X}_i \mathbf{A}_i)^T (\mathbf{X}_j \mathbf{A}_j))$$

$$48 \quad = \text{tr} \left(\sum_{i=1}^l \sum_{j=1}^l \bar{k}_{ij} (\mathbf{X}_i \mathbf{A}_i)^T (\mathbf{X}_j \mathbf{A}_j) \right) = \text{tr} \left(\sum_{i=1}^l \sum_{j=1}^l (\mathbf{X}_i \mathbf{A}_i)^T (\bar{k}_{ij} \mathbf{I}_n) (\mathbf{X}_j \mathbf{A}_j) \right)$$

$$49 \quad$$

50 Let $(\mathbf{X}_i \mathbf{A}_i)_{i=1}^l$ be the matrix combining $\mathbf{X}_i \mathbf{A}_i$ by rows. Let $\mathbb{X} = \text{diag}(\mathbf{X}_i)_{i=1}^l$, then
 51 $\mathbb{X} \mathbf{A}_{(13 \times 2)} = (\mathbf{X}_i \mathbf{A}_i)_{i=1}^l$. Reconstruct a Kronecker product, we have

$$52 \quad \mathbf{m}^T (\mathbf{K}_l^{-1} \otimes \mathbf{I}_{nk}) \mathbf{m} = \text{tr}([\mathbf{X}_i \mathbf{A}_i]_{i=1}^l)^T (\mathbf{K}_l^{-1} \otimes \mathbf{I}_n) [\mathbf{X}_j \mathbf{A}_j]_{j=1}^l$$

$$53 \quad = \text{tr}(\mathbf{A}_{(13 \times 2)}^T \mathbb{X}^T (\mathbf{K}_l^{-1} \otimes \mathbf{I}_n) \mathbb{X} \mathbf{A}_{(13 \times 2)})$$

$$54 \quad$$

55 Let $\check{\mathbf{Q}} = \mathbb{X}^T (\mathbf{K}_l^{-1} \otimes \mathbf{I}_n) \mathbb{X}$. With the ‘‘vec trick’’, we have

$$56 \quad \mathbf{m}^T (\mathbf{K}_l^{-1} \otimes \mathbf{I}_{nk}) \mathbf{m} = \text{tr}(\mathbf{A}_{(13 \times 2)}^T \check{\mathbf{Q}} \mathbf{A}_{(13 \times 2)}) = \text{vec}(\mathbf{A}_{(13 \times 2)})^T \text{vec}(\check{\mathbf{Q}} \mathbf{A}_{(13 \times 2)})$$

$$57 \quad (\text{SM1.6}) \quad = \text{vec}(\mathbf{A}_{(13 \times 2)})^T (\mathbf{I}_k \otimes \check{\mathbf{Q}}) \text{vec}(\mathbf{A}_{(13 \times 2)}) = \mathbf{a}_{(13 \times 2)}^T (\mathbf{I}_k \otimes \check{\mathbf{Q}}) \mathbf{a}_{(13 \times 2)}$$

59 So the posterior distribution has the form:

$$60 \quad (\text{SM1.7}) \quad p(\mathbf{m}|\mathfrak{X}) \propto \exp\left\{-\frac{1}{2}\mathbf{a}_{(13 \times 2)}^T(\mathbf{I}_k \otimes \check{\square})\mathbf{a}_{(13 \times 2)}\right\}$$

61 Substitute (SM1.5) and (SM1.7) into \circledast , we have:

$$62 \quad \circledast \propto \int_{\mathbb{R}^{kkl}} \exp\left(-\frac{1}{2}\left[\varepsilon^{-2}\|\mathbf{m}_* - (\mathbf{I}_k \otimes \tilde{\mathbf{X}})\mathbf{a}_{(13 \times 2)}\|^2 + \mathbf{a}_{(13 \times 2)}^T(\mathbf{I}_k \otimes \check{\square})\mathbf{a}_{(13 \times 2)}\right]\right) d\mathbf{a}$$

63 Note that we can expand the inner product to have:

$$64 \quad \|\mathbf{m}_* - (\mathbf{I}_k \otimes \tilde{\mathbf{X}})\mathbf{a}_{(13 \times 2)}\|^2 = \|\mathbf{m}_*\|^2 - 2\mathbf{m}_*^T(\mathbf{I}_k \otimes \tilde{\mathbf{X}})\mathbf{a}_{(13 \times 2)} + \mathbf{a}_{(13 \times 2)}^T(\mathbf{I}_k \otimes (\tilde{\mathbf{X}}^T \tilde{\mathbf{X}}))\mathbf{a}_{(13 \times 2)}$$

65 Denote $\Sigma_c^{-1} = \mathbf{I}_k \otimes (\varepsilon^{-2}\tilde{\mathbf{X}}^T \tilde{\mathbf{X}} + \check{\square})$ and $\mathbf{m}_c^T \Sigma_c^{-1} = \varepsilon^{-2}\mathbf{m}_*^T(\mathbf{I}_k \otimes \tilde{\mathbf{X}})$. Because $d\mathbf{a} =$
66 $d\mathbf{a}_{(13 \times 2)}$, we have

$$67 \quad \circledast \propto \int_{\mathbb{R}^{kkl}} \exp\left(-\frac{1}{2}\varepsilon^{-2}\|\mathbf{m}_*\|^2 + \mathbf{m}_c^T \Sigma_c^{-1} \mathbf{a}_{(13 \times 2)} - \frac{1}{2}\mathbf{a}_{(13 \times 2)}^T \Sigma_c^{-1} \mathbf{a}_{(13 \times 2)}\right) d\mathbf{a}_{(13 \times 2)}$$

$$68 \quad = \det(2\pi \Sigma_c)^{1/2} \exp\left(-\frac{1}{2}\varepsilon^{-2}\|\mathbf{m}_*\|^2 + \frac{1}{2}\mathbf{m}_c^T \Sigma_c^{-1} \mathbf{m}_c\right)$$

69

70 With the definitions of Σ_c^{-1} and $\mathbf{m}_c^T \Sigma_c^{-1}$, we have

$$71 \quad \varepsilon^{-2}\|\mathbf{m}_*\|^2 - \mathbf{m}_c^T \Sigma_c^{-1} \mathbf{m}_c = \varepsilon^{-2}\|\mathbf{m}_*\|^2 - (\mathbf{m}_c^T \Sigma_c^{-1})(\Sigma_c^{-1})^{-1}(\mathbf{m}_c^T \Sigma_c^{-1})^T$$

$$72 \quad = \varepsilon^{-2}\|\mathbf{m}_*\|^2 - \varepsilon^{-4}\mathbf{m}_*^T(\mathbf{I}_k \otimes \tilde{\mathbf{X}})(\mathbf{I}_k \otimes (\varepsilon^{-2}\tilde{\mathbf{X}}^T \tilde{\mathbf{X}} + \check{\square}))^{-1}(\mathbf{I}_k \otimes \tilde{\mathbf{X}})^T \mathbf{m}_*$$

$$73 \quad = \mathbf{m}_*^T \left(\varepsilon^{-2}\mathbf{I}_{nk} - \varepsilon^{-4}\mathbf{I}_k \otimes (\tilde{\mathbf{X}}(\varepsilon^{-2}\tilde{\mathbf{X}}^T \tilde{\mathbf{X}} + \check{\square})^{-1}\tilde{\mathbf{X}}^T)\right) \mathbf{m}_* = \mathbf{m}_*^T (\mathbf{I}_k \otimes \Sigma^\dagger) \mathbf{m}_*$$

74

75 Here we define

$$76 \quad (\text{SM1.8}) \quad \Sigma^\dagger = \varepsilon^{-2}\mathbf{I}_n - \varepsilon^{-4}\tilde{\mathbf{X}}(\varepsilon^{-2}\tilde{\mathbf{X}}^T \tilde{\mathbf{X}} + \check{\square})^{-1}\tilde{\mathbf{X}}^T$$

77 Because Σ_c does not depend on \mathbf{m}_* but $\circledast = p(\mathbf{m}_*|\mathfrak{X})$, we have

$$78 \quad (\text{SM1.9}) \quad \circledast \propto \exp\left(-\frac{1}{2}\mathbf{m}_*^T (\mathbf{I}_k \otimes \Sigma^\dagger) \mathbf{m}_*\right)$$

79 This means that the predictive distribution is

$$80 \quad (\text{SM1.10}) \quad \circledast := p(\mathbf{m}_*|\mathfrak{X}) = N_{nk}(\mathbf{m}_*; \mathbf{0}, \mathbf{I}_k \otimes \Sigma)$$

81 Now we simplify Σ . Recall that $\tilde{\mathbf{X}} = \mathbf{X}(\text{diag}(\mathbf{K}_l^{-1}\mathbf{k}_l) \otimes \mathbf{I}_k)$. Let $\mathbf{v} = \mathbf{K}_l^{-1}\mathbf{k}_l$. Using
82 the definition and properties of the Kronecker product, we have the following:

$$83 \quad \tilde{\mathbf{X}} = \mathbf{X}(\text{diag}(\mathbf{v}) \otimes \mathbf{I}_k) = (\mathbf{v}^T \otimes \mathbf{I}_n)\mathbb{X}$$

$$84 \quad \tilde{\mathbf{X}}^T \tilde{\mathbf{X}} = \mathbb{X}^T(\mathbf{v}^T \otimes \mathbf{I}_n)^T(\mathbf{v}^T \otimes \mathbf{I}_n)\mathbb{X} = \mathbb{X}^T[(\mathbf{v}\mathbf{v}^T) \otimes \mathbf{I}_n]\mathbb{X}$$

86 Recall that $\check{\square} = \mathbb{X}^T(\mathbf{K}_l^{-1} \otimes \mathbf{I}_n)\mathbb{X}$, from (SM1.8) and the above, we have

$$87 \quad \Sigma^\dagger = \varepsilon^{-2}\mathbf{I}_n - \varepsilon^{-4}\tilde{\mathbf{X}}\{\varepsilon^{-2}\mathbb{X}^T[(\mathbf{v}\mathbf{v}^T) \otimes \mathbf{I}_n]\mathbb{X} + \mathbb{X}^T(\mathbf{K}_l^{-1} \otimes \mathbf{I}_n)\mathbb{X}\}^{-1}\tilde{\mathbf{X}}^T$$

$$88 \quad = \varepsilon^{-2}\mathbf{I}_n - \varepsilon^{-4}\tilde{\mathbf{X}}\{\mathbb{X}^T[(\varepsilon^{-2}\mathbf{v}\mathbf{v}^T + \mathbf{K}_l^{-1}) \otimes \mathbf{I}_n]\mathbb{X}\}^{-1}\tilde{\mathbf{X}}^T$$

90 Let $\mathbf{D}_\mathbf{v} = \text{diag}(\mathbf{v})$, then $\tilde{\mathbf{X}} = \mathbf{X}(\mathbf{D}_\mathbf{v} \otimes \mathbf{I}_k)$. For simplicity we assume \mathbf{v} has no zero entries,
 91 which is almost always true, so that $\mathbf{D}_\mathbf{v}$ is invertible. Since $\mathbb{X}(\mathbf{D}_\mathbf{v}^{-1} \otimes \mathbf{I}_k) = (\mathbf{D}_\mathbf{v}^{-1} \otimes \mathbf{I}_n)\mathbb{X}$,

$$\begin{aligned} 92 \quad \Sigma^\dagger &= \varepsilon^{-2}\mathbf{I}_n - \varepsilon^{-4}\mathbf{X}(\mathbf{D}_\mathbf{v} \otimes \mathbf{I}_k)\{\mathbb{X}^T[(\varepsilon^{-2}\mathbf{v}\mathbf{v}^T + \mathbf{K}_l^{-1}) \otimes \mathbf{I}_n]\mathbb{X}\}^{-1}(\mathbf{D}_\mathbf{v} \otimes \mathbf{I}_k)\mathbf{X}^T \\ 93 \quad &= \varepsilon^{-2}\mathbf{I}_n - \varepsilon^{-4}\mathbf{X}\{(\mathbf{D}_\mathbf{v} \otimes \mathbf{I}_k)^{-1}\mathbb{X}^T[(\varepsilon^{-2}\mathbf{v}\mathbf{v}^T + \mathbf{K}_l^{-1}) \otimes \mathbf{I}_n]\mathbb{X}(\mathbf{D}_\mathbf{v} \otimes \mathbf{I}_k)^{-1}\}^{-1}\mathbf{X}^T \\ 94 \quad &= \varepsilon^{-2}\mathbf{I}_n - \varepsilon^{-4}\mathbf{X}\{\mathbb{X}^T(\mathbf{D}_\mathbf{v}^{-1} \otimes \mathbf{I}_n)[(\varepsilon^{-2}\mathbf{v}\mathbf{v}^T + \mathbf{K}_l^{-1}) \otimes \mathbf{I}_n](\mathbf{D}_\mathbf{v}^{-1} \otimes \mathbf{I}_n)\mathbb{X}\}^{-1}\mathbf{X}^T \\ 95 \quad &= \varepsilon^{-2}\mathbf{I}_n - \varepsilon^{-4}\mathbf{X}\{\mathbb{X}^T[(\varepsilon^{-2}\mathbf{D}_\mathbf{v}^{-1}\mathbf{v}\mathbf{v}^T\mathbf{D}_\mathbf{v}^{-1} + \mathbf{D}_\mathbf{v}^{-1}\mathbf{K}_l^{-1}\mathbf{D}_\mathbf{v}^{-1}) \otimes \mathbf{I}_n]\mathbb{X}\}^{-1}\mathbf{X}^T \\ 96 \quad &= \varepsilon^{-2}\mathbf{I}_n - \varepsilon^{-4}\mathbf{X}\{\mathbb{X}^T[(\varepsilon^{-2}\mathbf{1}_l\mathbf{1}_l^T + \mathbf{D}_\mathbf{v}^{-1}\mathbf{K}_l^{-1}\mathbf{D}_\mathbf{v}^{-1}) \otimes \mathbf{I}_n]\mathbb{X}\}^{-1}\mathbf{X}^T \end{aligned}$$

98 Define $\Omega = \varepsilon^{-2}\mathbf{1}_l\mathbf{1}_l^T + \mathbf{D}_\mathbf{v}^{-1}\mathbf{K}_l^{-1}\mathbf{D}_\mathbf{v}^{-1}$, then we have

$$99 \quad (\text{SM1.11}) \quad \Sigma^\dagger = \varepsilon^{-2}\mathbf{I}_n - \varepsilon^{-4}\mathbf{X}[\mathbb{X}^T(\Omega \otimes \mathbf{I}_n)\mathbb{X}]^{-1}\mathbf{X}^T$$

100 For now, let us assume Σ is invertible, then we can apply the Woodbury identity:

$$101 \quad (\mathbf{A} + \mathbf{C}\mathbf{B}\mathbf{C}^T)^{-1} = \mathbf{A}^{-1} - \mathbf{A}^{-1}\mathbf{C}(\mathbf{B}^{-1} + \mathbf{C}^T\mathbf{A}^{-1}\mathbf{C})^{-1}\mathbf{C}^T\mathbf{A}^{-1}$$

102 where we substitute $\mathbf{A} = \varepsilon^{-2}\mathbf{I}_n$, $\mathbf{B} = -[\mathbb{X}^T(\Omega \otimes \mathbf{I}_n)\mathbb{X}]^{-1}$, and $\mathbf{C} = \varepsilon^{-2}\mathbf{X}$. This gives:

$$103 \quad \Sigma = \varepsilon^2\mathbf{I}_n - \mathbf{X}[-\mathbb{X}^T(\Omega \otimes \mathbf{I}_n)\mathbb{X} + \varepsilon^{-2}\mathbf{X}^T\mathbf{X}]^{-1}\mathbf{X}^T$$

104 We note that generalizations of the Woodbury identity to the Moore–Penrose inverse
 105 usually do not have a simple formula. Since $\mathbf{X} = (\mathbf{1}_l^T \otimes \mathbf{I}_n)\mathbb{X}$, we have:

$$106 \quad \mathbf{X}^T\mathbf{X} = \mathbb{X}^T(\mathbf{1}_l^T \otimes \mathbf{I}_n)^T(\mathbf{1}_l^T \otimes \mathbf{I}_n)\mathbb{X} = \mathbb{X}^T[(\mathbf{1}_l\mathbf{1}_l^T) \otimes \mathbf{I}_n]\mathbb{X}$$

107 Let $\tilde{\mathbf{K}}_l = (\mathbf{D}_\mathbf{v}\mathbf{K}_l\mathbf{D}_\mathbf{v})^{-1}$, then $\Omega = \varepsilon^{-2}\mathbf{1}_l\mathbf{1}_l^T + \tilde{\mathbf{K}}_l$. We have:

$$\begin{aligned} 108 \quad \Sigma &= \varepsilon^2\mathbf{I}_n - \mathbf{X}[-\mathbb{X}^T(\Omega \otimes \mathbf{I}_n)\mathbb{X} + \varepsilon^{-2}\mathbb{X}^T[(\mathbf{1}_l\mathbf{1}_l^T) \otimes \mathbf{I}_n]\mathbb{X}]^{-1}\mathbf{X}^T \\ 109 \quad &= \varepsilon^2\mathbf{I}_n + \mathbf{X}\{\mathbb{X}^T[(\Omega - \varepsilon^{-2}\mathbf{1}_l\mathbf{1}_l^T) \otimes \mathbf{I}_n]\mathbb{X}\}^{-1}\mathbf{X}^T \\ 110 \quad (\text{SM1.12}) \quad &= \varepsilon^2\mathbf{I}_n + \mathbf{X}[\mathbb{X}^T(\tilde{\mathbf{K}}_l \otimes \mathbf{I}_n)\mathbb{X}]^{-1}\mathbf{X}^T \end{aligned}$$

112 Here the second term is positive semi-definite, so the overall matrix is nonsingular.
 113 Applying the Woodbury identity again we can verify that the inverse of (SM1.12)
 114 matches (SM1.11), therefore Σ is indeed invertible.

115 With (SM1.10) and (SM1.12), we complete the proof.

116 SM2. Joint distributions and random functions on Grassmann manifold.

117 In the main text we focus on point predictions on the Grassmann manifold, which is
 118 enough for PROM purposes. But more generally, our GP model induces a family of joint
 119 distributions on Grassmann manifolds, and can be used to generate random subspace-
 120 valued functions. Neither of these problems have been explored in the literature.

121 From section 3, we see that for any finite collection of parameter points $\boldsymbol{\theta} = (\boldsymbol{\theta}_i)_{i=1}^l$,
 122 our GP model gives a collection of random points on the Grassmann manifold $\mathfrak{M}_i =$
 123 $\text{span}(\text{vec}^{-1}(\bar{f}(\boldsymbol{\theta}_i)))$, whose marginal distributions are uniform: $\mathfrak{M}_i \sim \text{Uniform}(G_{k,n})$.
 124 For each $i \in \{2, \dots, l\}$, let $\Sigma_{\leq i}$ be defined by $\boldsymbol{\theta}_{\leq i} = (\boldsymbol{\theta}_j)_{j=1}^i$ and $\mathfrak{M}_{< i} = (\mathfrak{M}_j)_{j=1}^{i-1}$ as in
 125 (3.5). Then we have conditional distributions $\mathfrak{M}_i | \mathfrak{M}_{< i} \sim \text{MACG}(\Sigma_{\leq i})$. Combining the
 126 marginal and conditional distributions, we have a joint distribution on the Grassmann
 127 manifold, parameterized by $\boldsymbol{\theta}$:

$$128 \quad (\text{SM2.1}) \quad (\mathfrak{M}_i)_{i=1}^l \sim \text{Uniform}(G_{k,n}) \prod_{i=2}^l \text{MACG}(\Sigma_{\leq i})$$

129 GPS can be used to generate random subspace-valued functions. Suppose that $\boldsymbol{\theta}$
 130 is a sample grid to evaluate the random function, then we can use (SM2.1) to generate
 131 a sample path sequentially. The method to sample $\text{MACG}(\boldsymbol{\Sigma})$, including the uniform
 132 distribution, is implied in subsection 2.2, which requires $\boldsymbol{\Sigma}^{1/2}$. If we compute the
 133 EVD of $\boldsymbol{\Sigma}$ as in section 4, then we have $\boldsymbol{\Sigma}^{1/2} = \mathbf{V} \text{diag}(\sqrt{\sigma_i^2 + \varepsilon^2})_{i=1}^r \mathbf{V}^T + \varepsilon \mathbf{I}_n$. We
 134 summarize the overall sampling procedure in Algorithm SM2.1.

Algorithm SM2.1 GPS: Sampling a Random Subspace-valued Function

Require: correlation function $k(\cdot, \cdot)$.

Input: sample grid $(\boldsymbol{\theta}_i)_{i=1}^l$.

- 1: Generate random matrix: $\mathbf{Z} \in M_{n,k}$, $z_{ij} \sim N(0, 1)$.
- 2: Orthonormalization: $\mathbf{X}_1 \leftarrow \pi(\mathbf{Z})$.
- 3: **for** i in $2, \dots, l$ **do**
- 4: Generate random matrix: $\mathbf{Z} \in M_{n,k}$, $z_{ij} \sim N(0, 1)$.
- 5: Run Algorithms 4.1 and 4.2 with arguments $\mathbf{X}_{<i}$ and $(\boldsymbol{\theta}_{<i}, \boldsymbol{\theta}_i, r)$.
- 6: Matrix multiplication: $\mathbf{M} \leftarrow \mathbf{V} \text{diag}(\sqrt{\hat{\lambda} + \varepsilon^2} - \varepsilon) \mathbf{V}^T \mathbf{Z} + \varepsilon \mathbf{Z}$
- 7: Orthonormalization: $\mathbf{X}_i \leftarrow \pi(\mathbf{M})$.
- 8: **end for**

Output: Stiefel representations of subspaces $(\mathbf{X}_i)_{i=1}^l$.

Note: Projection $\pi(\mathbf{M}) = \mathbf{U} \mathbf{W}^T$, where $\mathbf{M} = \mathbf{U} \text{diag}(\sigma) \mathbf{W}^T$ is a thin SVD.

135 **SM3. Gradient of LOOCV predictive error.** The gradient of the LOOCV
 136 predictive error can also be computed. Denote $d_i = d_g(\mathbf{X}_i, \mathbf{V}_{-i})$ and let ∂ denote the
 137 partial derivative with respect to a scalar hyperparameter. With (5.3) and chain rule:

$$138 \quad (\text{SM3.1}) \quad \partial L_{\text{LOO}} = \sum_{i=1}^l \partial d_i^2 = -2 \sum_{i=1}^l \sum_{j=1}^k (\arccos \sigma_j) (1 - \sigma_j^2)^{-1/2} \partial \sigma_j$$

139 Here, $\sigma_j = \sigma_j(\mathbf{X}_i^T \mathbf{V}_{-i}) = \sigma_j(\tilde{\mathbf{C}}_i^T \mathring{\mathbf{V}}_{-i})$. Let $\tilde{\mathbf{C}}_i^T \mathring{\mathbf{V}}_{-i} = \hat{\mathbf{V}} \text{diag}(\boldsymbol{\sigma}) \hat{\mathbf{W}}^T$ be a thin SVD.
 140 Using the derivative of a singular value, see for example [SM3, p. 170], we have:

$$141 \quad (\text{SM3.2}) \quad \partial \sigma_j = \hat{\mathbf{v}}_j^T (\tilde{\mathbf{C}}_i^T \partial \mathring{\mathbf{V}}) \hat{\mathbf{w}}_j = \hat{\mathbf{v}}_j^T \tilde{\mathbf{C}}_i^T (\partial \mathring{\mathbf{V}}) \hat{\mathbf{w}}_j$$

142 Recall that $\mathring{\mathbf{V}}$ consists of the top- k eigenvectors of \mathbf{S}_{-i} . Let $(\mathring{\lambda}_p, \mathring{\mathbf{v}}_p)$ be the p -th
 143 eigenpair of \mathbf{S}_{-i} , $p = 1, \dots, k$. Using the derivative of an eigenvector of a symmetric
 144 matrix, see for example [SM2, Thm 8.9], we have:

$$145 \quad (\text{SM3.3}) \quad \partial \mathring{\mathbf{v}}_p = (\mathring{\lambda}_p \mathbf{I} - \mathbf{S}_{-i})^\dagger (\partial \mathbf{S}_{-i}) \mathring{\mathbf{v}}_p$$

146 Let $\mathbf{S}_{-i} = \mathring{\mathbf{Q}} \text{diag}(\mathring{\boldsymbol{\lambda}}) \mathring{\mathbf{Q}}^T$ be an EVD, then we have $(\mathring{\lambda}_p \mathbf{I} - \mathbf{S}_{-i})^\dagger = \mathring{\mathbf{V}} \text{diag}\{(\mathring{\lambda}_p -$
 147 $\mathring{\lambda}_q)^{-1}\}_{q=1}^r \mathring{\mathbf{V}}^T$. Recall that $\mathbf{S}_{-i} = \tilde{\mathbf{C}}_{-i} (\boldsymbol{\Pi}_{-i})^{-1} \tilde{\mathbf{C}}_{-i}^T$, we have:

$$148 \quad (\text{SM3.4}) \quad \partial \mathbf{S}_{-i} = -\tilde{\mathbf{C}}_{-i} (\boldsymbol{\Pi}_{-i})^{-1} (\partial \boldsymbol{\Pi}_{-i}) (\boldsymbol{\Pi}_{-i})^{-1} \tilde{\mathbf{C}}_{-i}^T$$

149 Recall that $\boldsymbol{\Pi}_{-i} = \square_{-i} \circ (\Delta_{-i} \otimes \mathbf{J}_k)$, $\Delta_{-i} = [\bar{k}_{pq} \bar{k}_{ii} / (\bar{k}_{ip} \bar{k}_{iq}) - 1]_{p,q \neq i}$, and $\bar{\mathbf{K}} = \mathbf{K}^{-1}$,
 150 we have:

$$151 \quad \partial \boldsymbol{\Pi}_{-i} = \square_{-i} \circ [(\partial \Delta_{-i}) \otimes \mathbf{J}_k]$$

$$152 \quad (\text{SM3.5}) \quad \partial [\Delta_{-i}]_{pq} = ([\Delta_{-i}]_{pq} + 1) (\bar{k}_{pq}^{-1} \partial \bar{k}_{pq} + \bar{k}_{ii}^{-1} \partial \bar{k}_{ii} - \bar{k}_{ip}^{-1} \partial \bar{k}_{ip} - \bar{k}_{iq}^{-1} \partial \bar{k}_{iq})$$

$$153 \quad \partial \bar{k}_{pq} = [\partial \mathbf{K}^{-1}]_{pq} = [-\mathbf{K}^{-1} (\partial \mathbf{K}) \mathbf{K}^{-1}]_{pq}$$

155 Combining (SM3.1)–(SM3.5), we can compute the partial derivative ∂L_{LOO} of the
 156 LOOCV predictive error with respect to any hyperparameter, as long as we can
 157 compute the partial derivative ∂k of the correlation function. For the SE kernel in (5.1)
 158 for example, $\partial k / \partial \beta_i = (\theta_i - \theta'_i)^2 \beta_i^{-3} k$. We omit a formal algorithm for the gradient
 159 computation, since it is straightforward given these equations.

160 We point out one way to speed up the evaluation of (SM3.3) and (SM3.4). Because
 161 the eigenvalues of \mathbf{S}_{-i} decline rapidly, we have

$$162 \text{ (SM3.6)} \quad (\mathring{\lambda}_p \mathbf{I} - \mathbf{S}_{-i})^\dagger \approx \mathring{\mathbf{V}} \text{diag}\{(\mathring{\lambda}_p - \mathring{\lambda}_q)^{-1}\}_{q=1}^\tau \mathring{\mathbf{V}}^T + \mathring{\lambda}_p^{-1}(\mathbf{I} - \mathring{\mathbf{V}}\mathring{\mathbf{V}}^T)$$

163 This approximation is accurate for any $p \in \{1, \dots, k\}$, as long as $\tau - k$ is reasonably
 164 large; for example, we can set $\tau = 2k$. To compute the approximation we only need the
 165 top τ eigenpairs of \mathbf{S}_{-i} . Since $r \approx kl > 2k$, the truncated EVD can be substantially
 166 faster than a full EVD. Algorithm SM3.1 gives an efficient procedure to compute $\partial \mathring{\mathbf{v}}_p$
 167 approximately given $\mathring{\mathbf{v}}_p$ and $\partial \mathbf{\Pi}_{-i}$.

Algorithm SM3.1 Approximate Computation of Derivative of an Eigenvector

Note: This procedure evaluates $\partial \mathring{\mathbf{v}}_p$ via (SM3.3) and (SM3.4) given $(\mathring{\mathbf{v}}_p, \partial \mathbf{\Pi}_{-i})$.

Require: $(\mathbf{L}, \tilde{\mathbf{L}}, \tilde{\mathbf{V}}, \tilde{\boldsymbol{\lambda}})$ from Algorithm 5.1.

- 1: $\mathbf{v} \leftarrow \text{solve}(\mathbf{L}^T, \tilde{\mathbf{L}}\tilde{\mathbf{v}}_p)$
 - 2: $\mathbf{v} \leftarrow \text{solve}(\mathbf{L}, (\partial \mathbf{\Pi}_{-i})\mathbf{v})$
 - 3: $\mathbf{u} \leftarrow -\mathring{\mathbf{V}}^T(\mathbf{L}\mathbf{v})$
 - 4: $\mathbf{w} \leftarrow \text{diag}\left\{(\mathring{\lambda}_p - \mathring{\lambda}_q)^{-1} - \mathring{\lambda}_p^{-1}\right\}_{q=1}^\tau$
 - 5: $\partial \mathring{\mathbf{v}}_p \leftarrow \mathring{\mathbf{V}}\mathbf{w} + \mathring{\lambda}_p^{-1}\mathbf{v}$
-

168 If the gradient is computed along with the LOOCV error, the additional cost is
 169 dominated by (1) the extended truncated EVD of \mathbf{S}_{-i} for l times and (2) the evaluation
 170 of Algorithm SM3.1 for kl times. Since the additional cost of truncated EVD takes
 171 about $\mathcal{O}(k^2 l^2 (\tau - k))$ flops, with $\tau = 2k$, part (1) takes about $\mathcal{O}(k^3 l^3)$ flops. Since
 172 Algorithm SM3.1 takes about $12k^2 l^2$ flops, part (2) takes about $12k^3 l^3$ flops. The
 173 overall additional cost is about $12k^3 l^3 + \mathcal{O}(k^3 l^3)$ flops per gradient evaluation, where
 174 the coefficient of the second term is determined by the truncated EVD algorithm.
 175 Compared with the $k^3 l^4$ flops for LOOCV error evaluation, the additional cost is at a
 176 similar level, depending on l .

177 **SM4. Other model selection criteria.** There are other model selection criteria
 178 for GP models in general. One popular possibility is to choose the hyperparameters to
 179 maximize the marginal likelihood with the GP integrated out. However, this approach
 180 is less robust to model and prior misspecification than CV. Another useful criteria is
 181 the LOOCV predictive probability density. We derived the analytical forms of both
 182 criteria for our model, and tried them for the numerical examples in this paper. In
 183 all cases, the marginal likelihood prefers infinite length-scales, inducing a singular
 184 covariance matrix. While the LOOCV predictive probability density can select a good
 185 length-scale for the visualization problem in subsection 7.1, it also prefers infinite
 186 length-scales in other problems, probably because $n \gg k$. We explain such behavior in
 187 this section.

188 The marginal likelihood of data is defined as the likelihood of data integrated
 189 over the prior. Recall that $\mathbf{x} = (\mathbf{x}_i)_{i=1}^l$, $\mathbf{x}_i = \text{vec}(\mathbf{X}_i)$, $\mathbf{X}_i \in V_{k,n}$, $\mathbf{m} = (\mathbf{m}_i)_{i=1}^l$, and
 190 $\mathbf{m}_i \in \mathbb{R}^{nk}$. Let $\mathfrak{M} = (\mathfrak{M}_i)_{i=1}^l$ and $\mathfrak{M}_i = \text{span}(\mathbf{m}_i)$, we can write the marginal likelihood

191 as:

$$192 \quad (\text{SM4.1}) \quad p(\mathbf{x}) = \int_{\mathbb{R}^{nkl}} p(\mathbf{m})L(\mathbf{x}|\mathfrak{M}) \, d\mathbf{m}$$

193 But from (3.3) we have likelihood $L(\mathbf{x}_i|\mathfrak{M}_i) = 1(\mathbf{x}_i \in [\mathbf{m}_i]) = 1(\mathbf{m}_i \in [\mathbf{x}_i])$, so the
 194 integrant in (SM4.1) only takes positive values for $\mathbf{m} \in \prod_{i=1}^l [\mathbf{x}_i]$, which is a measure-
 195 zero subset of the integration domain \mathbb{R}^{nkl} . This means that the marginal likelihood is
 196 identically zero.

197 Alternatively, we may modify the definition of marginal likelihood to only integrate
 198 over the support S of a singular likelihood, and define a modified marginal likelihood
 199 as:

$$200 \quad (\text{SM4.2}) \quad \tilde{p}(\mathbf{x}) = \int_S p(\mathbf{m})L(\mathbf{x}|\mathfrak{M}) \, d\mathbf{m}$$

201 PROPOSITION SM4.1. Let $\check{\square} = \mathbb{X}^T(\mathbf{K}_l^{-1} \otimes \mathbf{I}_n)\mathbb{X}$. The log modified marginal likeli-
 202 hood of data is:

$$203 \quad (\text{SM4.3}) \quad \log \tilde{p}(\mathbf{x}) = -\frac{1}{2}(n-k)kl \log(2\pi) - \frac{k}{2}(n \log |\mathbf{K}_l| + \log |\check{\square}|)$$

204 *Proof of Proposition SM4.1.* As in the proof of Theorem 3.1, the support of the
 205 likelihood can be written as $S = \prod_{i=1}^l \mathfrak{X}_i^k$, a linear subspace of \mathbb{R}^{nkl} where $\prod_{i=1}^l [\mathbf{x}_i]$ is
 206 a full-measure subset. Substituting prior joint distribution $\mathbf{m} \sim N_{nkl}(0, \mathbf{K}_l \otimes \mathbf{I}_{nk})$ into
 207 (SM4.2), we have:

$$208 \quad \tilde{p}(\mathbf{x}) = \int_S N_{nkl}(\mathbf{m}; 0, \mathbf{K}_l \otimes \mathbf{I}_{nk}) \prod_{i=1}^l 1(\mathbf{m}_i \in [\mathbf{x}_i]) \, d\mathbf{m}$$

210 With the same reasoning that leads to (SM1.1), let $\mathbf{m}_i = \text{vec}(\mathbf{X}_i \mathbf{A}_i)$, then we can
 211 change the integration domain to \mathbb{R}^{kkl} and replace $d\mathbf{m}$ with $d\mathbf{a}$, which gives:

$$212 \quad \tilde{p}(\mathbf{x}) = \int_{\mathbb{R}^{kkl}} N_{nkl}(\mathbf{m}; 0, \mathbf{K}_l \otimes \mathbf{I}_{nk}) \, d\mathbf{a}$$

$$213 \quad = \int_{\mathbb{R}^{kkl}} \det(2\pi \mathbf{K}_l \otimes \mathbf{I}_{kn})^{-1/2} \exp\left(-\frac{1}{2} \mathbf{m}^T (\mathbf{K}_l^{-1} \otimes \mathbf{I}_{kn}) \mathbf{m}\right) \, d\mathbf{a}$$

215 With (SM1.6), let $\check{\square} = \mathbb{X}^T(\mathbf{K}_l^{-1} \otimes \mathbf{I}_n)\mathbb{X}$ and because $d\mathbf{a} = d\mathbf{a}_{(13 \times 2)}$, we have:

$$216 \quad \tilde{p}(\mathbf{x}) = \int_{\mathbb{R}^{kkl}} \det(2\pi \mathbf{K}_l \otimes \mathbf{I}_{kn})^{-1/2} \exp\left(-\frac{1}{2} \mathbf{a}_{(13 \times 2)}^T (\mathbf{I}_k \otimes \check{\square}) \mathbf{a}_{(13 \times 2)}\right) \, d\mathbf{a}_{(13 \times 2)}$$

218 With Gaussian integral $\int_{\mathbb{R}^n} \exp\left(-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu})\right) \, d\mathbf{x} = \det(2\pi \boldsymbol{\Sigma})^{1/2}$, we have:

$$219 \quad \tilde{p}(\mathbf{x}) = \det(2\pi \mathbf{K}_l \otimes \mathbf{I}_{kn})^{-1/2} \det(2\pi(\mathbf{I}_k \otimes \check{\square}))^{-1/2}$$

$$220 \quad = (2\pi)^{-nkl/2} \det(\mathbf{K}_l)^{-nk/2} (2\pi)^{kkl/2} \det(\check{\square})^{-k/2}$$

$$221 \quad = (2\pi)^{-(n-k)kl/2} \det(\mathbf{K}_l)^{-nk/2} \det(\check{\square})^{-k/2}$$

223 Taking a logarithm gives the result in (SM4.3). \square

224 PROPOSITION SM4.2. Maximizing the modified marginal likelihood $\tilde{p}(\mathbf{x})$ leads to
 225 a singular covariance matrix \mathbf{K}_l .

226 *Proof of Proposition SM4.2.* With Proposition SM4.1, we have

$$227 \quad -\log \tilde{p}(\mathbf{x}) \propto h(\boldsymbol{\beta}) := n \log |\mathbf{K}_l| + \log |\mathbb{X}^T (\mathbf{K}_l^{-1} \otimes \mathbf{I}_n) \mathbb{X}|$$

229 Maximizing $\tilde{p}(\mathbf{x})$ is equivalent to minimizing the objective function $h(\boldsymbol{\beta})$. Let $\mathbf{Q} =$
 230 $(\mathbb{X}, \mathbb{X}_\perp)$ be an orthogonal completion of \mathbb{X} , then $|\mathbf{K}_l|^n = |\mathbf{K}_l \otimes \mathbf{I}_n| = |\mathbf{K}_l^{-1} \otimes$
 231 $\mathbf{I}_n|^{-1} = |\mathbf{Q}(\mathbf{K}_l^{-1} \otimes \mathbf{I}_n)\mathbf{Q}|^{-1}$. Let $\mathbf{B} = \mathbf{Q}^T(\mathbf{K}_l^{-1} \otimes \mathbf{I}_n)\mathbf{Q}$, with block structure $\mathbf{B} =$
 232 $[\mathbf{B}_{11} \ \mathbf{B}_{12}; \mathbf{B}_{12}^T \ \mathbf{B}_{22}]$ where \mathbf{B}_{11} is order- kl , then we have:

$$233 \quad h(\boldsymbol{\beta}) = \log \frac{|\mathbb{X}^T (\mathbf{K}_l^{-1} \otimes \mathbf{I}_n) \mathbb{X}|}{|\mathbf{Q}^T (\mathbf{K}_l^{-1} \otimes \mathbf{I}_n) \mathbf{Q}|} = \log \frac{|\mathbf{B}_{11}|}{|\mathbf{B}|}$$

235 Note that \mathbf{B} is positive semi-definite and so is \mathbf{B}_{11} . By the determinant properties
 236 of a block matrix, we have $|\mathbf{B}| = |\mathbf{B}_{11}| |\mathbf{C}_2|$, where $\mathbf{C}_2 = \mathbf{B}_{22} - \mathbf{B}_{12}^T \mathbf{B}_{11}^{-1} \mathbf{B}_{12}$. By
 237 the inverse properties of a block matrix, \mathbf{C}_2^{-1} is the trailing principal submatrix of
 238 $\mathbf{B}^{-1} = \mathbf{Q}^T (\mathbf{K}_l \otimes \mathbf{I}_n) \mathbf{Q}$. Therefore,

$$239 \quad h(\boldsymbol{\beta}) = \log(|\mathbf{C}_2|^{-1}) = \log |\mathbf{C}_2^{-1}| = \log |\mathbb{X}_\perp^T (\mathbf{K}_l \otimes \mathbf{I}_n) \mathbb{X}_\perp|$$

241 As \mathbf{K}_l tends to singularity, so does $|\mathbb{X}_\perp^T (\mathbf{K}_l \otimes \mathbf{I}_n) \mathbb{X}_\perp|$, which means the objective
 242 function $h(\boldsymbol{\beta})$ drops to negative infinity. Therefore, minimizing $h(\boldsymbol{\beta})$ selects a singular
 243 \mathbf{K}_l . \square

244 With an SE kernel, increasing length-scales drives \mathbf{K}_l to singularity. By Proposi-
 245 tion SM4.2, maximizing the modified marginal likelihood gives infinite length-scales.

246 Another model selection criteria is the log LOOCV predictive probability density.
 247 Because the predictive distribution of our GPS model is MACG($\boldsymbol{\Sigma}$), we have:

$$248 \quad \log p_{\text{LOO}} = \sum_{i=1}^l \log p_{\text{MACG}}(\mathbf{X}_i; \boldsymbol{\Sigma}_{-i}) = -\frac{1}{2} \sum_{i=1}^l (k \log |\boldsymbol{\Sigma}_{-i}| + n \log |\mathbf{X}_i^T (\boldsymbol{\Sigma}_{-i})^{-1} \mathbf{X}_i|)$$

250 Here, $\boldsymbol{\Sigma}_{-i}$ is defined in (5.4), predicting the i -th sample point using the other points.
 251 Similar to the proof of Proposition SM4.2, let $\mathbf{Q}_i = (\mathbf{X}_i, \mathbf{X}_{i\perp})$ be an orthogonal
 252 completion of \mathbf{X}_i , let $\mathbf{B} = \mathbf{Q}_i^T (\boldsymbol{\Sigma}_{-i})^{-1} \mathbf{Q}_i$, and let \mathbf{B}_{11} be its leading principal submatrix
 253 of order k , then

$$254 \quad \log p_{\text{LOO}} = -\frac{1}{2} \sum_{i=1}^l \log \frac{|\mathbf{B}_{11}|^n}{|\mathbf{B}|^k}$$

256 Note that both \mathbf{B} and \mathbf{B}_{11} are positive semi-definite, and of orders n and k respectively.
 257 As length-scale increases, both determinants increase. When n is not way larger than k ,
 258 as in our visualization example on $G_{1,2}$, the LOOCV predictive probability density can
 259 select a good length-scale. But when n is much larger than k , as in our example PROM
 260 problems, the numerator is less influential than the denominator, and maximizing
 261 p_{LOO} gives infinite length-scales.

262 **SM5. Computation time for the anemometer examples.** In Table 1, we
 263 compared the computational costs of the GPS and three other methods for PROM.
 264 The time complexities are broken down into various stages, measured in floating point
 265 operations, and are accurate up to the dominant terms. Coefficients are provided for
 266 all items except one. Therefore, this is the most general result for cost comparison.

TABLE SM1
Computation time for 1-parameter anemometer, $k = 20$.

	Preprocess	Subspace	ROM	Training
local POD	59s	7.5s	1.7s	-
GPS	0.6s	1.1s	1.7s	0.74s
Subspace-Int	2.0s	4.6s	1.0s	-
Manifold-Int	0.1s	-	1.1s	-
Matrix-Int	1.3s	-	0.1s	-

TABLE SM2
Computation time for 3-parameter anemometer, $l = 18$.

	Preprocess	Subspace	ROM	Training
local POD	71s	8.9s	1.9s	-
GPS	2.9s	3.2s	1.7s	4.2s
Subspace-Int	15s	10s	1.0s	-
Manifold-Int	0.54s	-	1.4s	-
Matrix-Int	2.2s	-	0.16s	-

267 Measured computation time depends on many factors besides the algorithm such
 268 as computer hardware, programming platform, algorithm implementation, and other
 269 processing commands apart from the main algorithm. It also depends on system
 270 dimension n , subspace dimension k , and sample size l . Such measurements can thus
 271 be misleading, and we do not provide them in the main text.

272 **Tables SM1** and **SM2** are typical computation times for the anemometer examples
 273 in **section 7**. In both cases, $n = 29,008$, $k = 20$, and we use $m = 50$ snapshots to
 274 generate the POD bases. Simulation time is included in the tables as the preprocessing
 275 step for local POD. For **Table SM1**, $l = 7$, parameter dimension $d = 1$, and the number
 276 of predictions is 101. For **Table SM2**, $l = 18$, parameter dimension $d = 3$ but an
 277 isotropic lengthscale is used, and the number of predictions is 118.

278 **SM6. A limitation of interpolation on tangent space.** In general, subspace
 279 interpolation is more accurate than the other two interpolation methods. But when: (1)
 280 sample size l is small; (2) subspace dimension k is large; or (3) parameter dimension d
 281 is large, the accuracy of all these methods can be unsatisfactory. [SM1] Sec. 9.6 also
 282 noted that the accuracy of matrix interpolation deteriorates between sample points
 283 when k increases, and gave a tentative explanation. Here we give an explanation of why
 284 interpolation on tangent spaces of a manifold, which includes subspace and manifold
 285 interpolation, fails in these situations.

286 When a point p' on a complete Riemannian manifold \mathcal{M} is pulled back to the
 287 tangent space $T_p\mathcal{M}$ of a reference point p via the exponential map, the preimage
 288 $\exp_p^{-1}(p')$ contains an infinite number of tangent vectors. The Riemannian logarithm
 289 $\log_p(p')$ is defined as the smallest tangent vector within this preimage, which lies
 290 in a star-shaped neighborhood of zero called the injectivity domain $ID(p)$. When a
 291 continuous map $f : \Theta \mapsto \mathcal{M}$ is pulled back to $T_p\mathcal{M}$, the preimage $(\exp_p^{-1} \circ f)(\Theta)$
 292 may have a connected component in $ID(p)$, which can be approximated given enough
 293 sample points. But this component will be increasingly distorted as it approaches the

294 boundary of $ID(p)$, called the tangent cut locus $TCL(p)$. This phenomenon can be
 295 observed, for example, in an azimuthal equidistant projection of the Earth. If the
 296 preimage only has connected components that intersects $TCL(p)$ or beyond, then the
 297 map cannot be approximated on $T_p\mathcal{M}$ by continuous maps interpolating points in
 298 $ID(p)$. As l decreases, d increases, or k increases, all sample points become further
 299 away from each other, and their Riemannian logarithms move closer to the tangent
 300 cut locus for any reference point. And as d or k increases, the map is more likely to
 301 cross the cut locus of any reference point. Therefore, the map becomes more difficult
 302 to approximate on the tangent space in these situations.

303 **SM7. On approximating local IRKA bases.** The microthruster example is
 304 just to showcase the accuracy of our proposed method when combined with a ROM
 305 method based on two-sided projection. The specific combination with IRKA may
 306 have several potential issues. First, IRKA only provides a local optimal ROM, and
 307 there may be an abundance of them depending on the dimensions of the full and
 308 the reduced model. Therefore, different runs of IRKA may give very different pairs
 309 of reduced subspaces, This is reflected in Figure 4, as the error curve of local IRKA
 310 is occasionally unsmooth. But for a method that approximates a subspace-valued
 311 mapping to work well, the true mapping needs to be well-defined and smooth in general.
 312 Second, a continuous trajectory of local \mathcal{H}_2 -optimal ROMs may not be all stable, which
 313 is possible because IRKA may converge to unstable ROMs. In fact, stability may break
 314 multiple times as parameter varies. Finally, there may not be a continuous trajectory
 315 of local \mathcal{H}_2 -optimal ROMs across the parameter space, so a good sample of local IRKA
 316 subspaces may not exist. In Figure SM1, we show some results for $k = 14$, where we
 317 use a sample of 10 points for our model. For the three segments of the parameter space
 318 where the error curve of local IRKA is relatively continuous, our method is able to
 319 maintain the error level, but overall the error curve is discontinuous and the ROMs
 320 can be unstable. This situation gets worse as k increases in this example.

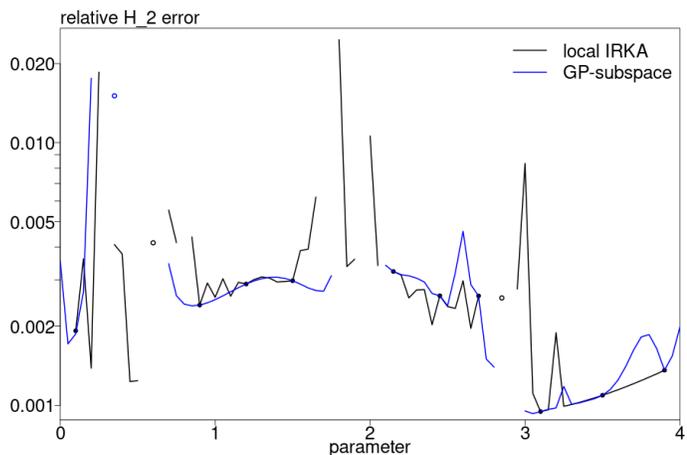


FIG. SM1. Relative \mathcal{H}_2 error for the microthruster. $k = 14$. Training data are shown as solid points. Disconnected test data are shown as hollow points.

- 322 [SM1] U. BAUR, P. BENNER, B. HAASDONK, C. HIMPE, I. MARTINI, AND M. OHLBERGER, *Chapter 9:*
 323 *Comparison of methods for parametric model order reduction of time-dependent problems,*
 324 *in Model Reduction and Approximation: Theory and Algorithms*, P. Benner, M. Ohlberger,
 325 A. Cohen, and K. Willcox, eds., SIAM, 2017, pp. 377–407.
- 326 [SM2] J. R. MAGNUS AND H. NEUDECKER, *Matrix Differential Calculus with Applications in Statistics*
 327 *and Econometrics*, Wiley, 2019.
- 328 [SM3] G. STRANG, *Linear Algebra and Learning from Data*, Wellesley-Cambridge Press, 2019.