Information and Inference: A Journal of the IMA (2022) 00, 1–30

https://doi.org/10.1093/imaiai/iaac021

# Non-dissipative and structure-preserving emulators via spherical optimization

#### Dihan Dai<sup>†</sup>

Department of Mathematics, and Scientific Computing and Imaging (SCI) Institute, University of Utah, 201 Presidents Cir, 84112 Salt Lake City, USA †Corresponding author. Email: dai@math.utah.edu

#### YEKATERINA EPSHTEYN

Department of Mathematics, University of Utah, 201 Presidents Cir, 84112 Salt Lake City, USA

AND

#### AKIL NARAYAN

Department of Mathematics, and Scientific Computing and Imaging (SCI) Institute, University of Utah, 201 Presidents Cir, 84112 Salt Lake City, USA

[Received on 26 August 2021; revised on 24 April 2022; accepted on 5 July 2022]

Approximating a function with a finite series, e.g., involving polynomials or trigonometric functions, is a critical tool in computing and data analysis. The construction of such approximations via now-standard approaches like least squares or compressive sampling does not ensure that the approximation adheres to certain convex linear structural constraints, such as positivity or monotonicity. Existing approaches that ensure such structure are norm-dissipative and this can have a deleterious impact when applying these approaches, e.g., when numerical solving partial differential equations. We present a new framework that enforces via optimization such structure on approximations and is simultaneously norm-preserving. This results in a conceptually simple convex optimization problem on the sphere, but the feasible set for such problems can be very complex. We establish well-posedness of the optimization problem through results on spherical convexity and design several spherical-projection-based algorithms to numerically compute the solution. Finally, we demonstrate the effectiveness of this approach through several numerical examples.

Keywords: structure-preserving emulators; high-order accuracy; quadratic programming; geodesic convex optimization.

#### 1. Introduction

Approximating an unknown function with a superposition of basis functions (e.g., polynomials or Fourier series) is a widely -used technique in computing and numerical analysis. For example, when solving a system of partial differential equations (PDEs), the class of spectral methods proposes such a superposition ansatz and determines the coefficients through minimization conditions on the PDE residual. Traditionally, fundamental properties of the approximation, such as stability, accuracy and computational efficiency are major considerations for the approximations. However, for certain problems, approximations are required to preserve certain implicit 'structures,' i.e., approximations should inherit certain desirable qualitative features of the original function. Such structure can include positivity, monotonicity, conservation of energy, etc. An approximation that fails to be structure-preserving may lead to numerical instability or even the failure of numerical schemes [36]. From the

broader viewpoint of building predictive emulators from data, this structure can be crucial to generate a meaningful emulator; for example, emulators built to predict population trends should not predict negative values. In this manuscript, we consider building approximations that respect general families of linear homogeneous convex inequality constraints (for which positivity and monotonicity are examples) along with a single quadratic equality constraint (an energy constraint).

Based on the existing framework for linear inequality constraints [35], we impose a new *spherical* constraint, i.e., a quadratic constraint in addition to the linear constraints. While a seemingly benign addition, this extra constraint substantially changes the optimization problem and its properties. With a coordinate vector  $\hat{p}$  provided (e.g., a vector of Fourier coefficients), the formulation we consider in this paper gives rise to an optimization problem of the form,

$$\begin{aligned} & \min_{\widehat{\mathbf{v}} \in \mathbb{R}^N} \|\widehat{\mathbf{v}} - \widehat{\boldsymbol{p}}\|_2^2 \\ \text{s.t. } & g_k(\widehat{\mathbf{v}}, y) \leq 0, \quad \forall y \in \omega_k, k \in [K] \\ & \|\widehat{\mathbf{v}}\|_2 = \|\widehat{\boldsymbol{p}}\|_2, \end{aligned} \tag{1.1}$$

where the linear constraints are given by the y-parameterized scalar-valued functions  $g_k(\cdot,y)$  and the energy-preserving constraint is given by the equality constraint  $\|\widehat{\boldsymbol{v}}\|_2 = \|\widehat{\boldsymbol{p}}\|_2$ . The parameters y can take values from a (possibly uncountably infinite) set  $\omega_k$ , and hence the feasible set can be very complex. Generally, the feasible set in (1.1) is the intersection of a finite collection of homogeneous convex cones and a sphere. The model (1.1) corresponds to a semi-infinite programming (SIP) problem [23]. In several SIP algorithms, a discrete approximation to the domain  $\omega_k$  is constructed (and perhaps refined). For linear constraints corresponding to positivity, this would correspond to requiring positivity at only a finite collection of points on the domain and hence structure is only preserved at a discrete set of points instead of on the whole domain. An additional difficulty is that the feasible set is a subset on the surface of a sphere, which is not a convex set in Euclidean space and hence the approaches from [35] do not apply. Thus, computationally solving (1.1) can be very challenging.

# 1.1 Related problems and approaches

There is existing literature on the study of optimization over ellipsoids (or spheres), which is closely related to the solution of the subproblems in the class of trust-region methods [18, 20, 30, 31]. However, in those approaches, the number of linear constraints is finite, and therefore such approaches are not directly applicable in our setting.

There are existing energy-preserving numerical methods that focus on energy-conservation of a Hamiltonian system, where a differential equation is discretized in a special way so that the energy of the discretized system is preserved, see [5, 10, 21, 29]. However, our focus is to preserve the energy of the approximation to a given function rather than the energy of a differential equation system, which is a different problem. In addition, methods built for differential equations assume very particular types of discretizations; the formulation we investigate in this paper applies to general discretizations.

A number of techniques have been proposed for preserving special kinds of structure for special choices of basis functions. To ensure positivity preservation, one can simply enforce positivity at a finite collection of points in the computational domain. The corresponding feasible set is a convex polytope and there are several algorithms available to computationally solve this problem [6]. Unfortunately, such techniques do not guarantee positivity of the approximation over the entire domain (a generally uncountable Euclidean set). An alternative to constraints over a finite set is to use special mappings.

For example, one can approximate  $\sqrt{f}$  and square the resulting approximation, or approximate  $\log f$ and subsequently exponentiate the approximation in order to guarantee the resulting approximation is positive. However, such mapping functions are not easy to construct for more complicated constraints, and the introduction of such maps can affect accuracy; for example,  $x \mapsto \sqrt{x}$  is not smooth at x = 0. For univariate polynomial approximation, one can take advantage of the special representations of nonnegative polynomials (Lukács theorem) to develop optimizations with a finite number of constraints [8].

In these univariate polynomial approximation methods that leverage Lukács's theorem, Nesterov [26] represented the convex cone of nonnegative polynomials over an interval by the linear images of the cone of positive semidefinite matrices. The corresponding optimization problem can then be solved by semidefinite programming (SDP) [34]. In [27], Nie and Demmel used the Nesterov's idea to determine a rational function satisfying several shape constraints. In [1], Allen and Kirby adapt the cone representation for Bernstein polynomial basis representation. In particular, they considered the best  $L^2$ - approximation and imposed a mass conservation constraint (a linear equality). They also proposed approximation for multivariate polynomials, but the corresponding feasible set is a subset (rather than an exact characterization) of the set of all nonnegative polynomials. We also note that there exists theoretical bounds for structure-preserving univariate polynomial approximations [1, 8, 13]. In special multivariate polynomial settings with a discrete number of constraints, additional estimates are possible [22, 25]. In this work, we are interested in the problem with more generality than the previously mentioned approaches. We allow more general (non-polynomial) basis functions and do not impose any restrictions for the domain or dimension of the space.

Other approaches include using an adaptive construction scheme for certain kinds of constraints [4] or linearly scaling the high-order coefficients of the polynomial to limit the oscillations [36]. Finally, we note that there are existing theoretical investigations for structure-preserving approximation in [2, 3, 14, 28], but these investigations do not translate into algorithms.

Our approach extends the recent technique in [35], which considers building approximations with linear structure, in the sense that the constraints are linear with respect to the approximant. In [35], the authors formalize a model for the structure-preserving problem with linear constraints, which applies to general, nontrivial linear structure. Under a mild condition, the corresponding function approximation problem can be cast to a semi-infinite convex optimization problem in a finite-dimensional Euclidean space with a unique solution. In addition, the work in [35] develops several projection-based algorithms to preserve the desired structures. However, their method, which amounts to filtering the approximation in a nonlinear manner, does not preserve the  $L^2$  norm, and thus is dissipative. The work of this paper preserves the quadratic (energy) norm through a modified formulation of the problem. This slight modification results in nontrivial changes to well-posedness and algorithmic development that we address.

#### Contributions of this paper

In this work, we are interested in providing theory and algorithms to address non-dissipative, structurepreserving function approximation methods of the form (1.1). Using notions of spherical convexity and spherical projections [16, 17], we show that the corresponding function approximation problem can be converted to a spherically convex feasibility problem, and establish uniqueness of the solution under mild conditions, see Theorem 3.1. Based on our theoretical results and by extending algorithms in [35], we propose three algorithms to solve the spherically convex feasibility problem; see sections 4.2, 4.3, and 4.4. Our algorithms do not rely on the discretization of the domain and therefore differ from many existing SIP algorithms [19, 32].

The setup of the general problem is as follows: We first assume that the unconstrained approximation to the unknown function is available, e.g., an unconstrained projection of the unknown function onto a finite-dimensional subspace. The unconstrained approximation is then *post-processed* via our algorithms so that the linear constraints, such as positivity, are satisfied without augmenting or reducing the quadratic energy of the approximant.

This paper is structured as follows. In Section 2, we give the theoretical framework of the structure-preserving function approximation problem as well as formalization of the constraints. In Section 3, we provide a brief overview of spherical geometry and present the uniqueness result of the function approximation problem. In Section 4, we discuss projections on the sphere and develop two algorithms for solving the function approximation problem. Finally, in Section 5, we demonstrate the efficacy of our algorithms with numerical results for polynomial and Fourier series approximations. Our energy-and structure-preserving results show similar rates of convergence as those of the unconstrained approximation as the subspace is refined.

# 2. Setup

Let  $\Omega \subseteq \mathbb{R}^d$  be a spatial domain. Consider the Hilbert space H formed by scalar-valued functions over  $\Omega$  with inner product  $\langle \cdot, \cdot \rangle_H$ ,

$$H = H(\Omega) := \{ f : \Omega \to \mathbb{R} \mid ||f||_H < \infty \}, \qquad ||f||^2 := \langle f, f \rangle_H,$$

A prototypical example is  $H = L^2(\Omega; \mathbb{R})$ . Let  $V \subseteq H$  be an N-dimensional subspace spanned by orthonormal basis functions  $\{v_n\}_{n \in [N]}$ ,

$$V = \operatorname{span}\{v_1, \dots, v_N\}, \qquad \langle v_i, v_k \rangle_H = \delta_{i,k}, \qquad j, k \in [N],$$

where  $\delta_{j,k}$  is Kronecker delta function, and  $[N] := \{1, \dots, N\}$ . Our numerical examples will be restricted to d = 1 or d = 2 on a closed interval or a closed rectangle  $\Omega$ , respectively, but the theoretical framework we develop holds for general choices of d and  $\Omega$ .

We assume throughout this document that V has no common zeros on  $\Omega$ , i.e., that,

$$\forall x \in \Omega \ \exists \ v \in V \text{ such that } v(x) \neq 0.$$

This assumption is true if, for example, V contains constant functions.

# 2.1 The unconstrained problem – linear measurements

We assume availability of an unconstrained function approximation scheme from H onto V using a finite collection of data. In this section, we briefly mention canonical approaches for accomplishing this via linear measurements, but our optimization problem is independent of how this unconstrained approximation is formed.

Let  $u \in H$  be a function about which a finite number of observations  $\{u_m\}_{m \in [M]} := \{\phi_m(u)\}_{m \in [M]} \subset \mathbb{R}$  are available, where  $\phi_1, \cdots, \phi_M$  are M linear functionals on H, and are bounded on V. The functionals can be, e.g.,  $v_m$ -projections  $\langle \cdot, v_m \rangle$  or pointwise evaluations  $\delta_{x_m}(\cdot)$ , where  $\delta_{x_m}$  is the Dirac mass centered

at  $x_m \in \Omega$ . An approximation  $p \in V$  to u is frequently built by enforcing these linear measurements:

Find 
$$p = \sum_{n \in [N]} \widehat{p}_n v_n$$
 satisfying  $A\widehat{p} = b$ , (2.1)

where

$$(A)_{m,n} = \phi_m(v_n), \qquad (m,n) \in [M] \times [N], \qquad b = [\phi_1(u), \cdots, \phi_M(u)]^\top \in \mathbb{R}^M.$$
 (2.2)

The condition M = N is necessary for the problem (2.1) to be well-posed, and so in practice one relaxes (2.1) in appropriate ways depending on whether the system is under-/over-determined. For example, with  $\hat{v}$  the  $v_n$ -coordinates of an element  $v \in V$ , one could relax (2.1) in the following ways:

$$(M > N) \, \widehat{\boldsymbol{p}} = \underset{\widehat{\boldsymbol{v}} \in \mathbb{R}^{N}}{\operatorname{arg \, min}} \, \left\| \boldsymbol{A} \widehat{\boldsymbol{v}} - \boldsymbol{b} \right\|_{2} \qquad \text{(Least squares)}$$

$$(M = N) \, \widehat{\boldsymbol{p}} = \boldsymbol{A}^{-1} \boldsymbol{b} \qquad \text{(Interpolation)}$$

$$(M < N) \, \widehat{\boldsymbol{p}} = \underset{\widehat{\boldsymbol{v}} \in \mathbb{R}^{N}: \, A \widehat{\boldsymbol{v}} = \boldsymbol{b}}{\operatorname{arg \, min}} \, \left\| \widehat{\boldsymbol{v}} \right\|_{1} \qquad \text{(Compressive sampling)}$$

$$(2.3)$$

where  $\|\cdot\|_p$  is the  $\ell^p([N])$  norm on vectors. Theory for well-posedness of each of these problems is mature [9, 12, 15, 33]. The numerical results in this paper utilize the interpolation (M = N) formulation above for simplicity, but this choice is independent of the theory and algorithms developed in this paper. The essential idea is that we assume the ability to construct  $\hat{p}$  that, in the absence of linear inequality or quadratic equality constraints, is considered a good approximation to the original function u based on available data.

#### 2.2 The constraints

In many practical situations, we require not only a solution to (2.1), but instead a solution that also obeys certain physical constraints, such as positivity over  $\Omega$ . The unconstrained approximation (2.1) need not obey any such constraints, even if the original function u does obey them, which may lead to unphysical approximations. We therefore consider the problem of imposing these additional constraints. We consider simultaneously imposing two types of constraints: a (possibly uncountable) set of linear constraints, along with a single quadratic constraint.

- The linear constraints. The constraints we consider in this section are motivated by the following examples of structural desiderata:
- positivity: p(y) > 0 for all  $y \in \Omega$ ,
- monotonicity: p'(y) > 0 for all  $y \in \Omega$ ,
- convexity: p''(y) > 0 for all  $y \in \Omega$ .

As it is pointed out in [35], these constraints can be characterized by families of linear constraints and a unique solution to the linearly constrained problem is guaranteed under some mild assumptions. Linearity in this context refers to linearity of the constraint with respect to the function p in V. In the rest of this subsection, we briefly review the abstract formulation introduced in [35].

Assume that there are K types of linear constraints (e.g., K=2 if we simultaneously impose positivity and monotonicity). For each  $k \in [K]$ , each type of linear constraint is a family defined by the condition,

$$L_k(v, y) \le 0,$$
  $\forall y \in \omega_k,$  (2.4)

where,

- $\omega_k$  is a subset (possibly containing uncountably many elements) of the spatial domain  $\Omega$ ,
- $L_k(\cdot, y)$  is a y-parameterized unit-norm element in the dual space  $V^*$ .

The feasible set of elements  $v \in V$  that satisfy (2.4) for the family-k constraint is given by

$$E_{\nu} := \{ \nu \in V \mid L_{\nu}(\nu, y) \le 0, \ \forall y \in \omega_{\nu} \}.$$
 (2.5)

Positivity, monotonicity, and convexity can be describes by the abstract formulation (2.4). The *linear* feasible set  $E^0$  is the set of all  $v \in V$  that satisfy all K constraints simultaneously, and hence is the intersection of all the  $E_k$ ,

$$E^{0} := \bigcap_{k \in [K]} E_{k} = \bigcap_{k \in [K]} \{ v \in V \mid L_{k}(v, y) \le 0, \ \forall y \in \omega_{k} \}.$$
 (2.6)

Note that  $E^0$  is always non-empty since it contains 0.

Since V is N-dimensional, we can identify the feasible set  $E^0$  in V with a feasible set in the  $v_n$ -coordinate space  $\mathbb{R}^N$ . By the Riesz representation theorem, for any  $L \in V^*$ , there exists a unique Riesz representor  $\ell \in V$  such that,

$$L(v) = \langle v, \ell \rangle_{H}, \quad \forall v \in V.$$

The function  $\ell \in V$  can also be written explicitly using the orthonormal basis  $\{v_n\}_{n \in [N]}$ ,

$$\ell(\cdot) = \sum_{n=1}^{N} \widehat{\ell}_n v_n(\cdot), \qquad \widehat{\ell}_n = \langle \ell, v_n \rangle = L(v_n),$$

and the following relation holds,

$$\|L\|_{V^*} = \|\ell\|_V = \|\widehat{\boldsymbol{\ell}}\|, \qquad \qquad \widehat{\boldsymbol{\ell}} = (\widehat{\ell}_1, \cdots, \widehat{\ell}_N)^{\mathrm{T}}.$$

In what follows we denote the Riesz representor for  $L_k(\cdot,y)$  by  $\ell_k(\cdot,y)$  and the corresponding coordinate vector by  $\widehat{\boldsymbol{\ell}}_k(y) \in \mathbb{R}^N$ . Since  $L_k(\cdot,y)$  is unit-norm, we have

$$||L_k(\cdot, y)||_{V^*} = ||\ell_k(y)||_V = ||\widehat{\ell}_k(y)|| = 1.$$
 (2.7)

Finally, the set  $C_k \subseteq \mathbb{R}^N$  corresponding to the feasible set  $E_k \subseteq V$  is given by,

$$C_k = \bigcap_{\mathbf{y} \in \omega_k} \left\{ \widehat{\mathbf{v}} \in \mathbb{R}^N \,\middle|\, \left\langle \widehat{\mathbf{v}}, \widehat{\boldsymbol{\ell}}_k(\mathbf{y}) \right\rangle \le 0 \right\} =: \bigcap_{\mathbf{y} \in \omega_k} c_k(\mathbf{y}), \qquad k \in [K], \tag{2.8}$$

and the set  $C^0 \subseteq \mathbb{R}^N$  corresponding to  $E^0$  is

$$C^0 = \bigcap_{k \in [K]} C_k. \tag{2.9}$$

It can be verified that all  $C_k$ ,  $k \in [K]$  are closed, convex cones in  $\mathbb{R}^N$  (and that  $E_k$  is a closed convex cone in V) [35] and thus their intersection  $C^0$  is also a convex cone. Note that, although  $C^0$  is simply a convex cone, the geometry of  $C^0$  can be very complicated with infinitely many extreme points since every  $C_k$  is the intersection of infinitely many half-spaces  $c_k(y)$  if the  $\omega_k$  is a set with infinite cardinality, e.g.,  $\omega_k$  is an interval.

REMARK 2.1. In [35], an additional  $r_k$  parameter is introduced to define *affine* convex cones as feasible sets. We specialize here to the homogeneous case  $r_k = 0$ , so that our cones all have vertices at the origin. If  $r_k \neq 0$ , the problem we consider in this paper is not necessarily well-posed; see Example 2.3.

We summarize one example from [35] to demonstrate the notation and how it can be specialized to familiar types of constraints.

Example 2.1. (Positivity) Let  $\Omega = [-1, 1]$  and V be any N-dimensional subspace of  $L^2(\Omega) \cap L^{\infty}(\Omega)$ . We want to impose a positivity-structure for  $v \in V$ :  $v(x) \geq 0$ ,  $\forall x \in \Omega$ . Thus, only K = 1 family is needed and  $\omega_1 = \Omega$ . Fixing  $y \in \omega_1$ , the corresponding unit-norm linear operator is given by

$$L_1(v, y) := -\lambda(y)v(y), \qquad \lambda(y) = \left(\sum_{n=1}^N v_n(y)^2\right)^{-\frac{1}{2}}, \qquad (2.10)$$

where  $\lambda(y)$  is a y-dependent normalization factor. The corresponding y-parameterized Riesz representor  $\ell_1(\cdot, y)$  and its coordinate vector  $\widehat{\ell}_1(y)$  are, respectively,

$$\ell_1(\cdot, y) = -\lambda(y) \sum_{n=1}^N \nu_n(y) \nu_n(\cdot), \qquad \widehat{\ell}_1(y) = \left[ -\lambda(y) \nu_1(y), \cdots, -\lambda(y) \nu_N(y) \right]^\top. \tag{2.11}$$

Once the orthonormal basis is specified,  $\widehat{\ell}_1(y)$  can thus be explicitly identified.

2.2.2 Constraints that are "determining" In order to establish uniqueness of the solution to our quadratic-linear constrained problem, we require an additional condition on the linear constraints  $(L_k, \omega_k)$  defining  $C^0$ .

Definition 2.1. The set of constraints  $(L_k, \omega_k)_{k \in [K]}$  is V-determining if

$$v \in V \text{ and } L_k(v, y) = 0 \ \forall y \in \omega_k \forall k \in [K]$$
  $\Longrightarrow$   $v = 0$ 

We will assume V-determining linear constraints, which amounts to a technical assumption about the geometry of the associated  $\mathbb{R}^N$ -feasible set  $C^0$  that we later exploit. The V-determining condition precludes certain problem setups, but all the practical situations we consider in this paper are V-determining. As a simple example, to enforce positivity for every point in  $\Omega$  as in Example 2.1 we have that  $L_1(v,y)$  is a scaled point evaluation at y. Therefore, the V-determining condition requires that if  $v \in V$  satisfies v(y) = 0 for every  $y \in \Omega$  then v = 0, which is a quite natural condition.

For more intuition, the following lists some additional examples, with  $\Omega = [-1, 1]$ , K = 1, and  $\ell_1$  the normalized point-evaluation operator in Example 2.1,

- If  $V = \text{span}\{x^j | j = 0, \dots, N-1\}$  and  $|\omega_1| \ge N$ , then the linear constraint is V-determining
- If  $V = \text{span}\{x^j | j = 0, ..., N-1\}$  and  $|\omega_1| < N$ , then the linear constraint is *not V*-determining
- If  $V = \text{span}\{x^j | j = 1, ..., N\}$ , and  $|\omega_1| \le N$  with  $0 \in \omega_1$ , then the linear constraint set is *not* V-determining.
- If  $V = \text{span}\{H(x), 1 H(x)\}$ , with H the Heaviside function, and  $\omega_1 = [-1, 0.5]$ , then the linear constraint set is V-determining.
- If  $V = \text{span}\{H(x), 1 H(x)\}$ , with H the Heaviside function, and  $\omega_1 = [-1, -0.5]$ , then the linear constraint set is *not V*-determining.

Note that the V-determining condition is violated only for specialized cases, e.g., either when  $\omega_1$  has finite cardinality less than N, or when V contains very special types of functions. In all the numerical examples we consider, the linear constraints are V-determining.

2.2.3 The quadratic energy constraint In addition to the linear constraints, we further impose a single quadratic norm constraint analogous to an  $L^2$ -energy of the function. To be precise, we impose that our constrained solution must have same norm as the unconstrained solution.

$$||v||_{H} = ||p||_{H},\tag{2.12}$$

where p is the unconstrained solution with coordinate vector  $\hat{p}$  from solving (2.3). The corresponding discretized constraint set in  $\mathbb{R}^N$  is a sphere with radius  $\|\hat{p}\|$ ,

$$C^{H} := \{\widehat{\mathbf{v}} \in \mathbb{R}^{N} | \|\widehat{\mathbf{v}}\|_{\mathbb{R}^{N}} = \|\widehat{\mathbf{p}}\|_{\mathbb{R}^{N}}\}, \tag{2.13}$$

where  $\hat{v}$  is the coordinate vector for v.

We can now state the overall procedure we consider in this paper:

- 1. Given data  $\{u_m\}_{m\in[M]}$ , solve the unconstrained problem (2.3) to obtain the unconstrained solution  $\widehat{p}$ .
- 2. Post-process the unconstrained solution p by solving the constrained problem

$$\boldsymbol{d} = \underset{\widehat{\boldsymbol{v}} \in C}{\arg\min} \frac{1}{2} \|\widehat{\boldsymbol{v}} - \widehat{\boldsymbol{p}}\|^2, \tag{2.14}$$

where  $C = C^H \cap C^0$ . This constrained problem optimizes a quadratic function over a subset C of a sphere in  $\mathbb{R}^N$ .

Our focus is on the theory and algorithms for the second step, post-processing the unconstrained solution to obtain a structure-preserving solution. We show in Theorem 3.1 that the problem above has a unique solution. Furthermore, in Section 4, we are able to naturally extend the existing algorithms proposed in [35] to the new optimization problem. We end this section with two examples that illustrate why some alternative formulations to the two-step procedure above do not necessarily result in wellposed problems.

EXAMPLE 2.2. (An alternative formulation with nonunique solutions) One possible alternative to our framework proposed above is to instead consider the following constrained problem

$$d = \arg\min_{\widehat{\mathbf{v}} \in C} \frac{1}{2} ||A\widehat{\mathbf{v}} - \mathbf{b}||^2. \tag{2.15}$$

with A and b as introduced in Section 2.1. This formulation incorporates the constraints and a leastsquares problem simultaneously. However, the solution to this alternative formulation (2.15) is not necessarily unique. The issue lies in the fact that if the singular values of the full-rank matrix A are not all equal to 1, then the problem corresponds to optimization over an ellipsoid, which can yield non-unique solutions.

For example, consider the case when N=M=2. Let

$$A = \begin{bmatrix} 0.4 & 0 \\ 0 & 1 \end{bmatrix}, b = \begin{bmatrix} 0 \\ 0.5 \end{bmatrix}, C = \{(\cos t, \sin t) | 0.01 \le t \le \pi - 0.01 \}.$$

The constrained set C is a spherically convex set (Definition 3.4). The loss/cost function associated to (2.15) is

$$cost(t) = 0.4\cos^2 t + (\sin t - 0.5)^2 = 0.6\sin^2 t - \sin t + 0.65,$$

which has two distinct global minima over the feasible set C at  $t = \arcsin(\frac{5}{6})$  and  $t = \pi - \arcsin(\frac{5}{6})$ .

EXAMPLE 2.3. (Nonhomogeneous cones as in Remark 2.1) In Remark 2.1, we mention that in previous work [35] an  $r_k$  parameter is introduced to include more general linear constraints. If  $r_k \neq 0$ , then our optimization problem does not necessarily have unique solutions. Consider the following problem. Let the feasible set

$$C = \{(x, y) | x^2 + y^2 = 1\} \bigcap \{(x, y) | y \le 4x + 2, y \le -4x + 2\}.$$

The feasible set consists of two disjoint arcs (red arcs in Figure 1). If the unconstrained solution lies in the middle of the dark green arc (the black dot), there will be two solutions to (2.14), one from each red arc.

# 3. Solution to the constrained optimization problem

In this section, we will study the solution to the constrained optimization (2.14). Specifically, we will show in Theorem 3.1 that the solution is unique under reasonable conditions.

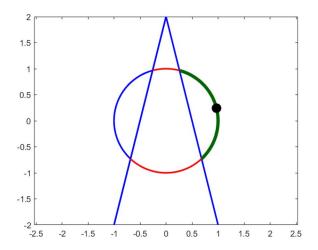


Fig. 1. An illustration to Example 2.3.

#### 3.1 *Spherical geometry*

We introduce some definitions and relevant results for the spherical geometry in this subsection. We refer to [16, 17] for technical details. We largely focus on the unit sphere in this section, i.e.,  $C^H = \mathbb{S}^{N-1}$ . In Section 4, we will use the more general origin-centered sphere of nonzero radius.

The *intrinsic distance* on  $\mathbb{S}^{N-1}$  is defined to be the great circle distance between two points, which corresponds to the angle between the two unit vectors in the ambient space.

DEFINITION 3.1. (Intrinsic distance on the sphere) Given  $u, w \in \mathbb{S}^{N-1}$ , the intrinsic distance between them is

$$d(\mathbf{u}, \mathbf{w}) = \arccos\langle \mathbf{u}, \mathbf{w} \rangle. \tag{3.1}$$

If S is an origin-centered sphere with radius r > 0, then the intrinsic distance between  $u, w \in S$  is

$$d_r(\mathbf{u}, \mathbf{w}) = r \arccos \langle \mathbf{u} / \| \mathbf{u} \|, \mathbf{w} / \| \mathbf{w} \| \rangle. \tag{3.2}$$

Note that the intrinsic distance between two points on a sphere of radius  $r \neq 1$  is given by the intrinsic distance between the unit-normalized points, scaled by the radius.

DEFINITION 3.2. (Geodesics on a sphere) A geodesic on the unit sphere  $\mathbb{S}^{N-1}$  is a *great circle*, i.e, the intersection curve of the sphere and a hyperplane in  $\mathbb{R}^N$  through the origin. The unique arclength-parameterized geodesic segment from  $\boldsymbol{u}$  to  $\boldsymbol{w}$ , where  $\boldsymbol{u}, \boldsymbol{w} \in \mathbb{S}^{N-1}$  and  $\boldsymbol{u} \neq \pm \boldsymbol{w}$ , is given by

$$\gamma_{uw}(t) = \csc d(u, w) [u \sin(d(u, w) - t) + w \sin t], \qquad t \in [0, d(u, w)].$$
 (3.3)

The (non-unique) geodesic segments joining u and -u, starting at u with velocity v satisfying ||v|| = 1 at u, is given by

$$\gamma_{u(-u)} := \cos(t)u + \sin(t)v, \qquad t \in [0, \pi].$$
 (3.4)

For a general sphere S centered at the origin with radius r, the geodesic segments can be defined via rescaling (3.3) or (3.4).

DEFINITION 3.3. (Exponential Mapping) The exponential mapping at u is defined to be

$$exp_{\boldsymbol{u}}: T_{\boldsymbol{u}} \mathbb{S}^{N-1} \to \mathbb{S}^{N-1}, \qquad \qquad \boldsymbol{v} \to \boldsymbol{u} \cos(\|\boldsymbol{v}\|) + \frac{\boldsymbol{v}}{\|\boldsymbol{v}\|} \sin(\|\boldsymbol{v}\|),$$
 (3.5)

which maps an element on the tangent plane  $T_u \mathbb{S}^{N-1}$  at u to the endpoint of the geodesic segment of length ||v|| starting at u in the direction of v. In (3.3), the geodesic segment can be expressed as  $\gamma_{uw}(t) = exp_u(t\gamma'_{uw}(0))$ . For a general sphere S centered at the origin with radius r, the geodesic as well as the exponential mapping can be defined via rescaling (3.5).

DEFINITION 3.4. (Spherically convex set) A subset  $C \subseteq \mathbb{S}^{N-1}$  is said to be *spherically convex* if for any  $s, t \in C$ , all the geodesic segments joining s and t are contained in C.

PROPOSITION 3.1. ([17], Proposition 2) Let  $C \subseteq \mathbb{S}^{N-1}$ . C is a spherically convex set if and only if the cone

$$K_C = \{ z s | s \in C, z \in [0, +\infty) \}$$
 (3.6)

is convex (in Euclidean sense) and pointed, i.e.,  $K_C \cap (-K_C) = \{0\}$ .

If we choose  $C = C^H \cap C^0$  as is done in our problem articulated in Section 2.2.3, then  $K_C = C^0$ , the convex cone corresponding to only the linear equalities. Thus, in our framework spherical convexity of  $C^H \cap C^0$  is determined by whether or not  $C^0$  is a pointed cone.

Proposition 3.2. A closed hemisphere is *not* spherically convex.

*Proof.* Noticing the existence of the antipodal points on a closed hemisphere, then there is a nontrivial v such that  $v, -v \in K_C$ . Therefore  $K_C \cap (-K_C)$  contains at least one nontrivial point, and Proposition 3.1 yields the conclusion.

One major utility of convex sets on the sphere is the ability to perform projections.

DEFINITION 3.5. (Spherical projection onto a closed convex set) Let  $C \subset \mathbb{S}^{N-1}$  be a spherically convex, closed set. The projection of  $z \in \mathbb{S}^{N-1}$  onto C is defined to be:

$$\mathscr{P}_C^s(z) = \left\{ t \in \mathbb{S}^{N-1} | d(t, z) \le d(t, r), \forall r \in C \right\}, \tag{3.7}$$

i.e., the nearest intrinsic distance projection.

The definition above does not immediately reveal uniqueness or computability for this type of projection, but the following proposition proved in [17] shows the relation between spherical projection

onto a closed spherically convex set and the Euclidean projection onto the convex cone spanned by the spherical convex set.

PROPOSITION 3.3. ([17], Proposition 8) Let  $C \subseteq \mathbb{S}^{N-1}$  be a spherical convex set. Take  $z \in \mathbb{S}^{N-1}$ . Let  $u = \mathcal{P}_{K_C}(z)$ , be the Euclidean projection of z onto  $K_C$ , the latter of which is defined in (3.6). If  $u \neq 0$ , then the spherical projection of z onto C is unique, and is given by,

$$\mathscr{P}_C^s(z) = \frac{u}{\|u\|} = \exp_z v,$$

where

$$v = \left(-z \cot \theta + \frac{u}{\|u\|} \csc \theta\right) \theta, \qquad \theta = d(z, u/\|u\|).$$

# 3.2 Uniqueness of the solution to (2.14)

In this subsection, we present a uniqueness theorem for the solution to (2.14). Our first step is to introduce the formulation (3.8) below that is equivalent to (2.14).

LEMMA 3.1. The constrained optimization problem (2.14) is equivalent to finding the spherical projection of  $\hat{p}$  onto the feasible set C,

$$\mathbf{d} = \underset{\widehat{\mathbf{v}} \in C}{\arg\min} d(\widehat{\mathbf{v}}, \widehat{\mathbf{p}}). \tag{3.8}$$

Proof. The proof is direct,

$$\begin{split} \underset{\widehat{\boldsymbol{v}} \in C}{\arg\min} & \frac{1}{2} \|\widehat{\boldsymbol{v}} - \widehat{\boldsymbol{p}}\|^2 = \underset{\widehat{\boldsymbol{v}} \in C}{\arg\min} \left( \frac{1}{2} (\|\widehat{\boldsymbol{v}}\|^2 + \|\widehat{\boldsymbol{p}}\|^2) - \langle \widehat{\boldsymbol{v}}, \widehat{\boldsymbol{p}} \rangle \right) \\ &= \underset{\widehat{\boldsymbol{v}} \in C}{\arg\max} \langle \widehat{\boldsymbol{v}}, \widehat{\boldsymbol{p}} \rangle, \\ &= \underset{\widehat{\boldsymbol{v}} \in C}{\arg\min} d(\widehat{\boldsymbol{v}}, \widehat{\boldsymbol{p}}). \end{split} \tag{3.9}$$

Using Lemma 3.1, we are able to show the following uniqueness theorem.

THEOREM 3.1. Assume the following hold for the constraint set  $C = C^H \cap C^0$  defined by (2.8)-(2.9) and (2.13):

- (a) The set of constraints  $(L_k, \omega_k)_{k \in [K]}$  are V-determining in the sense of Definition 2.1.
- (b) The Euclidean projection onto the linearly-constrained set satisfies  $\mathscr{P}_{C^0}\widehat{p} \neq 0$ .
- (c) The set  $C^H$  in (2.13) is  $\mathbb{S}^{N-1}$ , i.e.,  $\|\widehat{p}\| = 1$ .

Then the solution to (3.8) (or equivalently, (2.14)) is unique.

*Proof.* We first show that C is a closed spherically convex set. A subsequent application of Proposition 3.3 will prove the result.

Since  $c_k(y)$  are closed half-spaces,  $C^H \cap c_k(y)$  are closed hemispheres and

$$C = C^H \bigcap C^0 = \bigcap_{y \in \omega_k, k \in [K]} (C^H \cap c_k(y))$$

is closed. On the other hand, from Proposition 3.1, the set C is spherically convex if and only if the cone  $K_C$  (see also Proposition 3.1 for the definition) is convex and pointed. Direct calculation shows that  $K_C = C^0$ , which has been shown to be a closed convex set in [35].

Dofina

$$W := C^0 \bigcap \{-C^0\} = \left\{ \nu | \left\langle \widehat{\ell}_k(y), \nu \right\rangle = 0, \forall y \in \omega_k, k \in [K] \right\}. \tag{3.10}$$

Take  $x \in W \subseteq C^0$ . By assumption (a) and Definition 2.1, the only element v of V satisfying  $L_k(v, y) = 0$  for every  $y \in \omega_k$  and  $k \in [K]$  is v = 0. Thus,  $W = \{0\}$  and therefore  $K_C = C^0$  is pointed.

By Proposition 3.3 set C is a spherical convex set and the solution to (3.8) (or equivalently, (2.14)) is unique.

COROLLARY 3.1. The conclusions in Theorem 3.1 hold with loosening assumption (c) to  $\|\hat{p}\| > 0$ .

The proof is direct since all arguments hold unchanged via scaling by  $\|\hat{p}\| > 0$ .

# 3.3 Approximation errors from quadratic constraints

The optimization problem we consider, (2.14) or equivalently (3.8), contains both linear and quadratic in/equality constraints in general multivariate and non-polynomial scenarios. As discussed in Section 1.1, there are many existing approaches in more specialized cases that consider only linear inequality constraints for polynomials (such as positivity). Such specializations sometimes yield approximation errors of the linearly constrained problem relative to unconstrained best approximation errors. (See, e.g., [1, 8, 13].)

Since our framework imposes an additional quadratic constraint in addition to linear constraints, a natural question is how much introduction of this extra constraint affects approximation errors. We provide a simple result below, indicating that the error committed by our linear-quadratic constraints is comparable to the error committed by only the linearly constrained problem. Therefore, any theory establishing error of the linearly constrained solution can immediately be ported to our linear-quadratically constrained case.

Precisely, we recall our notation from Section 2.1 where u is a given function in a Hilbert space H and p is some prescribed approximation to u from the finite-dimensional space V. Recalling that  $E^0 \subset V$  is the closed convex set corresponding to our linear inequality constraints defined in (2.6), then the projections  $\mathscr{P}_{E^0}(u)$  and  $\mathscr{P}_{E^0}(p)$  onto  $E^0 \subseteq V$  are unique (and the latter is equivalent to the vector projection  $\mathscr{P}_{C^0}(\widehat{p})$  in  $\mathbb{R}^N$ ). A simple application of triangle inequality yields

$$\|u - \mathscr{P}_{E^0}(p)\|_H \le \|u - p\|_H + \|p - \mathscr{P}_{E^0}(p)\|_H,$$
 (3.11)

providing an upper bound on the error committed with *only* linear constraints. In particular, the first error term above is due to imposition of p as the initial unconstrained approximation, and the second term is due solely to the optimization/projection. In our linear-quadratically constrained case, if  $\hat{p}^*$  is

any vector solution to the optimization (2.14) corresponding to the function  $p^* \in V$ , then we would hope that  $||u-p^*||_H$  has error comparable to (3.11). The following lemma establishes precisely this statement.

LEMMA 3.2. Let  $u \in H$  and  $p \in V$  be given, where p has expansion coefficients  $\widehat{p}$ . Let a linear-quadratic constraint set  $C = C^H \cap C^0$  be given, where  $C^0$  is the linear constraint set, and let  $p^* \in V$  denote the function associated to a(ny) solution  $\widehat{p}^*$  of the linear-quadratic constrained problem (2.14) or equivalently (3.8). Then we have the following estimates,

$$\|u - p^*\|_H \le \|u - p\|_H + 2 \|p - \mathcal{P}_{E^0}(p)\|_H$$
 (3.12)

$$\|u - p^*\|_H \le 5 \|p - p^{\dagger}\|_H + 2\sqrt{2} \|u - \mathcal{P}_{E^0}(u)\|_H,$$
 (3.13)

where  $p^{\dagger}$  is the *H*-best approximation from *V*, i.e.,

$$p^{\dagger} = \operatorname*{arg\,min}_{p \in V} \|u - p\|_{H}.$$

Before proving this result we emphasize the practicality of these statements: (3.12) provides an upper bound comparable to (3.11), indicating that the linear-quadratically constrained solution  $p^*$  has error similar to an approximation with only linear constraints. The estimate (3.13) is a statement in terms of best approximations, showing that the error in  $p^*$  splits into (i) a component  $\|p-p^{\dagger}\|_H$  that is due to the quality of the initial unconstrained approximation provided, and (ii) the best approximation error from the linearly constrained space  $E^0 \subseteq V$ . In other words, modulo the quality of p relative to  $p^{\dagger}$ , the linear-quadratic constrained approximation  $p^*$  commits an error scaling precisely like the error in the best possible linearly-constrained approximation.

*Proof. of Lemma* 3.2 We first prove (3.12). Using the triangle inequality, we first divide the error into one part solely due to imposition of p as the starting unconstrained approximation, and a second part due to the linear-quadratic optimization p:

$$\|u - p^*\|_H \le \|u - p\|_H + \|p - p^*\|_H$$
 (3.14a)

We compute the last term using our notation in  $\mathbb{R}^N$ , i.e.,

$$\|p - p^*\|_H = \|\widehat{\boldsymbol{p}} - \widehat{\boldsymbol{p}}^*\|. \tag{3.14b}$$

To bound this norm, first assume that  $\mathscr{P}_{C^0}(\widehat{p})$ , the solution to the linearly-constrained problem, is non-zero. Then define,

$$q := \frac{\|\widehat{p}\|}{\|\mathscr{P}_{C^0}(\widehat{p})\|} \mathscr{P}_{C^0}(\widehat{p}) \in C.$$

A direct computation shows,

$$\begin{aligned} \|\widehat{\boldsymbol{p}} - \widehat{\boldsymbol{p}}^*\| &\overset{(2.14)}{\leq} \|\widehat{\boldsymbol{p}} - \boldsymbol{q}\| \\ &\leq \|\widehat{\boldsymbol{p}} - \mathcal{P}_{C^0}(\widehat{\boldsymbol{p}})\| + \|\mathcal{P}_{C^0}(\widehat{\boldsymbol{p}}) - \boldsymbol{q}\| \\ &= \|\widehat{\boldsymbol{p}} - \mathcal{P}_{C^0}(\widehat{\boldsymbol{p}})\| + \|\mathcal{P}_{C^0}(\widehat{\boldsymbol{p}})\| - \|\widehat{\boldsymbol{p}}\| \| \\ &\leq 2 \|\widehat{\boldsymbol{p}} - \mathcal{P}_{C^0}(\widehat{\boldsymbol{p}})\|, \end{aligned}$$

where the last inequality uses the reverse triangle inequality. If we have  $\mathscr{P}_{C^0}(\widehat{p}) = \mathbf{0}$ , then defining q instead as any feasible point in C and using the same computations along with  $\|q\| = \|\widehat{p} - \mathscr{P}_{C^0}(\widehat{p})\|$  proves the same inequality. Therefore,

$$\|\widehat{\boldsymbol{p}} - \widehat{\boldsymbol{p}}^*\| \le 2 \|\widehat{\boldsymbol{p}} - \mathscr{P}_{C^0}(\widehat{\boldsymbol{p}})\| \tag{3.15}$$

is true in general. Combining this with (3.14) proves (3.12).

Finally, since  $u - p^{\dagger}$  is orthogonal to V then using the Pythagorean theorem can be used to begin the following inequality chain:

$$\|u-p^*\|_{H} = \sqrt{\|u-p^{\dagger}\|_{H}^{2} + \|p^{\dagger} - p^{*}\|_{H}^{2}}$$

$$\leq \|u-p^{\dagger}\|_{H} + \|p^{\dagger} - p^{*}\|_{H}$$

$$\leq \|p-p^{\dagger}\|_{H} + \|u-p^{\dagger}\|_{H} + \|p-p^{*}\|_{H}$$

$$\leq \|p-p^{\dagger}\|_{H} + \|u-p^{\dagger}\|_{H} + 2\|p-\mathcal{P}_{E^{0}}(p)\|_{H}$$

$$\leq \|p-p^{\dagger}\|_{H} + \|u-p^{\dagger}\|_{H} + 2(\|p-p^{\dagger}\|_{H} + \|p^{\dagger} - \mathcal{P}_{E^{0}}(p^{\dagger})\|_{H} + \|\mathcal{P}_{E^{0}}(p^{\dagger}) - \mathcal{P}_{E^{0}}(p)\|_{H})$$

$$= 3\|p-p^{\dagger}\|_{H} + 2\|\mathcal{P}_{E^{0}}(p^{\dagger}) - \mathcal{P}_{E^{0}}(p)\|_{H} + (\|u-p^{\dagger}\|_{H} + 2\|p^{\dagger} - \mathcal{P}_{E^{0}}(p^{\dagger})\|_{H})$$

$$\stackrel{(*)}{\leq} 5\|p-p^{\dagger}\|_{H} + (\|u-p^{\dagger}\|_{H} + 2\|p^{\dagger} - \mathcal{P}_{E^{0}}(p^{\dagger})\|_{H})$$

$$\leq 5\|p-p^{\dagger}\|_{H} + 2(\|u-p^{\dagger}\|_{H} + \|p^{\dagger} - \mathcal{P}_{E^{0}}(p^{\dagger})\|_{H})$$

$$\stackrel{(**)}{\leq} 5\|p-p^{\dagger}\|_{H} + 2\sqrt{2}\sqrt{\|u-p^{\dagger}\|_{H}^{2} + \|p^{\dagger} - \mathcal{P}_{E^{0}}(p^{\dagger})\|_{H}}$$

$$= 5\|p-p^{\dagger}\|_{H} + 2\sqrt{2}\|u-\mathcal{P}_{E^{0}}(u)\|_{H}, \qquad (3.16)$$

where the inequality (\*\*) uses  $|a| + |b| \le \sqrt{2}\sqrt{a^2 + b^2}$ , and the inequality (\*) uses the non-expansive property of projections onto closed convex sets,

$$\left\|\mathscr{P}_{E^0}(q)-\mathscr{P}_{E^0}(r)\right\|_H\leq \|q-r\|_H\,,$$

D. DAI ET AL.

see, e.g., [11, Theorem 3]. The final equality (3.16) is the fact that,

$$\left\| u - \mathscr{P}_{E^0}(u) \right\|_H^2 = \min_{q \in E^0} \left\| u - q \right\|_H^2 = \min_{q \in E^0} \left\| u - p^\dagger \right\|_H^2 + \left\| p^\dagger - q \right\|_H^2 = \left\| u - p^\dagger \right\|_H^2 + \left\| p^\dagger - \mathscr{P}_{E^0}(p^\dagger) \right\|_H^2.$$

Note that Lemma 3.2 does not require uniqueness of the optimization (3.8). In particular, we do *not* require assumptions (a)-(c) from Theorem 3.1. The corollary (3.13) shows that if the unconstrained input p to our linear-quadratic optimization problem is a good approximation to u (in particular if it is the best approximation  $p^{\dagger}$ ), then the optimization error committed is at most  $2\sqrt{2}$  times that of the best possible linearly constrained approximation to u.

# 4. Algorithm: spherical projections

Having established the well-posedness of the problem (3.8), we proceed to discuss algorithms for solving the problem. In particular, we extend some procedures from [35] to the spherical optimization problem (3.8). We will shift our focus back to a general sphere centered at the origin with radius  $r \neq 0$ , equipped with the intrinsic distance  $d_r(\cdot, \cdot)$  (Equation 3.2). Here,  $r = \|\widehat{p}\| \neq 0$  is the norm of the unconstrained solution.

## 4.1 *Spherical projection onto a closed hemisphere*

To start, we first compute the spherical projection of a point on the sphere onto a closed hemisphere  $c_k(y) \cap C^H$ , which later serves as an ingredient of our main algorithms. Note that Proposition 3.3 is not directly applicable since a closed hemisphere is not spherically convex The proof in this section is elementary but we provide it in order to make our work self-contained.

Theorem 4.1. Let  $\widehat{p}$  be given, and fix (k,y). If  $\widehat{\ell}_k(y)$  is *not* parallel to  $\widehat{p}$ , then the spherical projection of  $\widehat{p}$  onto the closed hemisphere  $C^H \cap c_k(y)$  is unique, i.e., the solution  $c_k(\widehat{p};y)$  to

$$\boldsymbol{c}_{k}(\widehat{\boldsymbol{p}}; y) = \underset{\widehat{\boldsymbol{v}} \in C^{H} \cap \mathcal{C}_{k}(y)}{\arg\min} d_{r}\left(\widehat{\boldsymbol{v}}, \widehat{\boldsymbol{p}}\right), \tag{4.1}$$

is unique, and is given by

$$\boldsymbol{c}_{k}(\widehat{\boldsymbol{p}};\boldsymbol{y}) = \left\{ \begin{array}{c} \widehat{\boldsymbol{p}}, \ \widehat{\boldsymbol{p}} \in C^{H} \cap c_{k}(\boldsymbol{y}), \\ \frac{\mathscr{P}_{l}\widehat{\boldsymbol{p}}}{\|\mathscr{P}_{l}\widehat{\boldsymbol{p}}\|} \|\widehat{\boldsymbol{p}}\|, \ \widehat{\boldsymbol{p}} \not\in C^{H} \cap c_{k}(\boldsymbol{y}), \end{array} \right. \tag{4.2}$$

where  $\mathscr{P}_L$  is the Euclidean projection operator onto the subspace L, with the latter defined as,

$$L := \partial c_k(y) = \left\{ s | \left\{ \widehat{\ell}_k(y), s \right\} = 0 \right\}.$$

*Proof.* For simplicity, we will suppress  $k, \widehat{p}$ , and y notationally in the proof, i.e.,  $c := c_k(\widehat{p}; y)$ , and  $\widehat{\ell} := \widehat{\ell}_k(y)$ . Since  $C^H \cap c_k(y)$ , is non-empty and compact, there is at least one solution to (4.1).

Following similar computations to the proof to Lemma 3.1, it can be shown that

$$\mathbf{c} = \underset{\widehat{\mathbf{v}} \in C^H \cap c_k(y)}{\operatorname{arg}} \max_{\langle \widehat{\mathbf{v}}, \widehat{\mathbf{p}} \rangle}, \tag{4.3}$$

is equivalent to (4.1). If  $\widehat{p} \in C^H \cap c_k(y)$ , then  $c = \widehat{p}$  is the unique solution to (4.3) by the Cauchy-Schwarz inequality, which verifies part of (4.2). Thus, the remainder of the proof assumes  $\widehat{p}$  is not in the feasible set. Let  $\widehat{c}$  be any solution to (4.3). Since  $\widehat{c}$  lies in  $c_k(y)$  and since  $\widehat{p}$  lies in  $C^H$  but is not feasible, then we have

$$\langle \widehat{\boldsymbol{c}}, \widehat{\boldsymbol{\ell}} \rangle \leq 0, \qquad \langle \widehat{\boldsymbol{p}}, \widehat{\boldsymbol{\ell}} \rangle > 0.$$

By the above inequalities, any solution  $\hat{c}$  to (4.3) satisfies,

$$\begin{split} \langle \widehat{\boldsymbol{c}}, \widehat{\boldsymbol{p}} \rangle &= \left\langle \mathscr{P}_L \widehat{\boldsymbol{c}}, \mathscr{P}_L \widehat{\boldsymbol{p}} \right\rangle + \left\langle (I - \mathscr{P}_L) \widehat{\boldsymbol{c}}, (I - \mathscr{P}_L) \widehat{\boldsymbol{p}} \right\rangle = \left\langle \mathscr{P}_L \widehat{\boldsymbol{c}}, \mathscr{P}_L \widehat{\boldsymbol{p}} \right\rangle + \left\langle \widehat{\boldsymbol{c}}, \widehat{\boldsymbol{\ell}} \right\rangle \left\langle \widehat{\boldsymbol{p}}, \widehat{\boldsymbol{\ell}} \right\rangle, \\ &\stackrel{(i)}{\leq} \left\langle \mathscr{P}_I \widehat{\boldsymbol{c}}, \mathscr{P}_I \widehat{\boldsymbol{p}} \right\rangle \stackrel{(ii)}{\leq} \left\| \mathscr{P}_I \widehat{\boldsymbol{c}} \right\| \left\| \mathscr{P}_I \widehat{\boldsymbol{p}} \right\| \stackrel{(iii)}{\leq} \left\| \widehat{\boldsymbol{c}} \right\| \left\| \mathscr{P}_I \widehat{\boldsymbol{p}} \right\| = \left\| \widehat{\boldsymbol{p}} \right\| \left\| \mathscr{P}_I \widehat{\boldsymbol{p}} \right\| \end{split}$$

The choice  $\hat{c} = c$  in (4.2) is the unique solution that achieves equality in (i), (ii), and (iii) above. To see this, first note that c is feasible since it lies in both  $C^H$  and  $c_k(y)$ , and is well-defined since  $\hat{p}$  is not parallel to  $\hat{\ell}_k(y)$  and hence  $\mathscr{P}_L\hat{p} \neq 0$ . Equality in (i) and (iii) can be established by noting that  $c \in L$ , so that  $\langle c, \hat{\ell} \rangle = 0$  and  $\mathscr{P}_L c = c$ . Equality in (ii) is achieved if and only if  $\mathscr{P}_L c = c$  has the same direction as  $\mathscr{P}_L\hat{p}$ , which the choice (4.2) satisfies. This also shows that c is the only vector that achieves this equality, and hence (4.3) (equivalently, (4.1)) has a unique solution (4.2).

REMARK 4.1. The solution (4.2) implies that when  $\widehat{p}$  is not feasible and is not parallel to  $\widehat{\ell}_k(y)$ , the solution to (4.1) can be computed by first computing a Euclidean projection onto the hyperplane L, and the simply rescaling this projection to have norm  $\|\widehat{p}\|$ . We exploit this fact in algorithms.

#### 4.2 A greedy approach

We first introduce a new notation for the spherical projection

$$\boldsymbol{c}_{k}(\widehat{\boldsymbol{p}}; \mathbf{y}) := \mathscr{P}^{s}_{c_{k}(\mathbf{y})}\widehat{\boldsymbol{p}}, \tag{4.4}$$

which by Theorem 4.1 is well-defined for every  $\widehat{p}$  that is not a multiple of  $\widehat{\ell}_k(y)$ .

A greedy procedure, in the spirit of the greedy algorithm of [35], iteratively updates  $\hat{p}$  by repeatedly identifying most-violated constraints. Defining  $\hat{p}^0 = \hat{p}$ , and using  $\hat{p}^j$  to denote the iterate at step j, we seek to compute,

$$\widehat{\boldsymbol{p}}^{j+1} = \mathscr{P}^{s}_{c_{k^*}(y^*)}\widehat{\boldsymbol{p}}^{j}, \qquad (k^*, y^*) := \underset{k \in [K], y \in \omega_k}{\arg\max} \ d_r(\widehat{\boldsymbol{p}}^{j}, c_k(y)), \tag{4.5}$$

for  $j \geq 1$ . Lemma 4.1 first allows us to conclude that the set of (k, y) such that  $d_r(\widehat{\boldsymbol{p}}, c_k(y) \cap C^H) > 0$  is equal to the set of (k, y) such that  $\widehat{\boldsymbol{p}} \notin c_k(y)$ .

D. DAI ET AL.

LEMMA 4.1. Let  $\hat{p}$  be the solution to the current iteration, then

$$d_r(\widehat{\boldsymbol{p}}, c_k(y) \cap C^H) > 0 \Leftrightarrow \widehat{\boldsymbol{p}} \notin c_k(y), \tag{4.6}$$

where  $r = \|\widehat{\boldsymbol{p}}\|$ .

*Proof.* Let  $dist(\cdot, \cdot)$  be the Euclidean distance function, then

$$\operatorname{dist}(\widehat{\boldsymbol{p}}, c_k(y)) = \min_{\boldsymbol{s} \in c_k(y)} \|\widehat{\boldsymbol{p}} - \boldsymbol{s}\|_2 = \|\widehat{\boldsymbol{p}} - \mathcal{P}_L \widehat{\boldsymbol{p}}\|_2 = \|\widehat{\boldsymbol{p}}\| \sin \theta_k(y), \tag{4.7}$$

where  $L:=\partial c_k(y)$  is the boundary of the half-space  $c_k(y)$ , and  $\theta_k(y)=\arccos\left\langle\widehat{\boldsymbol{p}}/r,\boldsymbol{c}_k(\widehat{\boldsymbol{p}};y)/r\right\rangle$  is the angle between  $\widehat{\boldsymbol{p}}$  and its spherical projection  $\boldsymbol{c}_k(\widehat{\boldsymbol{p}};y)$  onto the half-space  $c_k(y)$ . The last equality in (4.7) is true since  $\mathscr{P}_L\widehat{\boldsymbol{p}}$  and  $\boldsymbol{c}_k(\widehat{\boldsymbol{p}};y)$  have the same direction (Theorem 4.1). The angle  $\theta_k(y)$  is the angle between the vector  $\widehat{\boldsymbol{p}}$  and the plane L, thus  $\theta_k(y)\in[0,\pi/2]$ .

For any (k, y), a direct calculation using the Pythagorean theorem and the definition of the intrinsic distance  $d_r(\cdot, \cdot)$  yields

$$\begin{split} \widehat{\pmb{p}} \not\in c_k(y) &\Leftrightarrow \operatorname{dist}(\widehat{\pmb{p}}, c_k(y)) > 0, \\ &\Leftrightarrow \|\widehat{\pmb{p}}\| \sin \theta_k(y) > 0, \\ &\Leftrightarrow \theta_k(y) > 0, \\ &\Leftrightarrow d_r(\widehat{\pmb{p}}, c_k(y) \cap C^H) > 0, \end{split}$$

The proof to Lemma 4.1 motivates the following relation

$$\begin{split} (k^*, y^*) &= \underset{k \in [K], y \in \omega_k}{\arg\max} \ d_r(\widehat{\boldsymbol{p}}, c_k(y)) \\ &= \underset{k \in [K], y \in \omega_k}{\arg\max} \ (\operatorname{dist}(\widehat{\boldsymbol{p}}, c_k(y))) \\ &= \underset{y \in \omega_k, k \in [K]}{\arg\min} \ \operatorname{sdist}(\widehat{\boldsymbol{p}}, c_k(y)), \end{split} \tag{4.8}$$

where the Euclidean signed distance between  $\hat{p}$  and the Euclidean half-space  $c_k(y)$  can be computed by (see [35])

$$sdist(\widehat{\boldsymbol{p}}, c_k(y)) = -\langle \widehat{\boldsymbol{\ell}}_k(y), \widehat{\boldsymbol{p}} \rangle. \tag{4.9}$$

Equations (4.8)-(4.9) imply that, to determine the parameters for the geodesically farthest hemisphere, we only need to determine the parameters for the Euclidean-farthest hyperplane, which is a much easier computational task.

Algorithm 1 summarizes the above iterative procedure of computing the solution to (3.8) through greedy spherical projection.

**Algorithm 1:** Iterative greedy spherical projection to compute the solution to (4.1).

**Require:** the matrix A and the observational vector b, tolerance parameter  $\delta \geq 0$ 

**Require:** constraints  $\{(\widehat{\ell}_k(y), \omega_k)\}_{k \in [K]}$ 

**Ensure:** greedy solution  $\hat{p}$ .

1: Compute the unconstrained solution  $\hat{p}$ , e.g. via solving (2.3).

2: **while** sdist( $\widehat{p}$ ,  $c_k(y)$ )  $\leq -\delta$  for some  $k \in [K]$ ,  $y \in \omega_k$  **do** 

3: compute  $(y^*, k^*)$  via (4.8)

4: update  $\hat{p}$  via (4.5)

5: end while

## 4.3 An averaging approach

The greedy procedure above can lead to oscillatory behavior of the iteration trajectory. To mitigate this behavior, we introduce an averaging projection approach to suppress potential oscillatory behavior of iterates in the previous greedy approach. The notion of an average position of a collection of points on the sphere is defined by the *Karcher mean*, which is a natural extension of the Euclidean weighted average.

DEFINITION 4.1. (Karcher mean) Let  $S \subset \mathbb{R}^N$  be the sphere centered at the origin with radius r. Let  $q_1, \dots, q_J$  be J points lying on S associated with nonnegative convex weights  $w_1, \dots, w_J \in [0, 1]$ . The *Karcher mean* is given by the solution to

$$q = \underset{x \in S}{\operatorname{arg\,min}} \left( \frac{1}{2} \sum_{i=1}^{n} w_i \cdot d_r^2(x, q_i) \right) := \underset{x \in S}{\operatorname{arg\,min}} f(x). \tag{4.10}$$

The Karcher mean as defined above is unique under mild assumptions.

THEOREM 4.2. ([7], Theorem 1) With S the radius-r origin-centered sphere in  $\mathbb{R}^N$ , suppose that given points  $q_1, \dots, q_J$  all lie in a closed hemisphere  $\mathscr{H} \subset S$ , with at least one point  $q_j$  in the interior of  $\mathscr{H}$  with  $w_j > 0$ . Then f(x) defined in (4.10) has a single critical point q in the interior of  $\mathscr{H}$ , and this point q is the global minimum of f, hence the unique Karcher mean.

The definition of the Karcher mean in (4.10) can be extended to a collection of infinitely many points by integration, and our averaged projection algorithm is based on this generalized Karcher mean. Let  $\hat{p}^0 = \hat{p}$  be the first iterate in the algorithm. To compute the next iterate, the averaging algorithm first identifies all the parameters (y,k) for which the associated linear constraints are violated,

$$\omega_{k^{-}}^{j} = \{ y \in \omega_{k} | \operatorname{sdist}(\widehat{\boldsymbol{p}}^{j}, c_{k}(y)) < 0 \}. \tag{4.11}$$

Instead of projecting onto the most violated constraint (as in the previous section) we seek a point that minimizes the Karcher mean objective over all violated constraints. With  $r = \|\widehat{p}\|$ , the averaged position

D. DAI *ET AL*.

 $\hat{p}^{j+1}$  at the next iteration is given by

$$\widehat{\boldsymbol{p}}^{j+1} = \underset{\boldsymbol{x} \in C^H}{\arg\min} \left( \frac{1}{2} \sum_{k \in \mathcal{D}_j^-} \int_{\omega_{k^-}^j} d_r^2(\boldsymbol{x}, \boldsymbol{c}_k(\widehat{\boldsymbol{p}}^j; \boldsymbol{y})) w_k(\widehat{\boldsymbol{p}}^j, \boldsymbol{y}) d\boldsymbol{y} \right), \tag{4.12}$$

where  $\mathcal{D}_{j}^{-} = \{k \in [K] | \omega_{k^{-}}^{j} \neq \varnothing\}$  is the set of the indexes where the corresponding constraints are violated at jth iteration, and the weight,

$$w_k(\widehat{\boldsymbol{p}}^j; y) := \left( \frac{d_r^2(\widehat{\boldsymbol{p}}^j, \boldsymbol{c}_k(\widehat{\boldsymbol{p}}^j; y))}{\sum_{\ell \in \mathcal{D}_j^-} \int_{\omega_{\ell^-}^j} d_r^2(\widehat{\boldsymbol{p}}^j, \boldsymbol{c}_\ell(\widehat{\boldsymbol{p}}^j; z)) dz} \right)$$
(4.13)

is introduced for each spherical projection  $c_k(y)$  in order to prioritize updates that mitigate the impact of the more violated constrained sets. In practice, we approximate the integral (4.12) via quadrature with positive weights, i.e.,

$$\widehat{\boldsymbol{p}}^{j+1} = \underset{\boldsymbol{x} \in C^H}{\arg\min} \left( \frac{1}{2} \sum_{k \in \mathcal{D}_j^-} \sum_{q=1}^{Q_k} w_{k,q} d_r^2(\boldsymbol{x}, \boldsymbol{c}_k(\widehat{\boldsymbol{p}}^j; y_{k,q})) \right), \tag{4.14}$$

where  $Q_k$  is the number of the quadrature points associated with the kth constraint, and  $\{w_{k,q}\}_{k \in [K], q \in [Q_k]}$  are the product of the (positive) quadrature weight with an approximations to the weight (4.13) at the quadrature points  $\{y_{k,q}\}_{k \in [K], q \in [Q_k]}$ .

All the discussion above is provided in the context of assuming that the hemisphere condition in Theorem 4.2 holds. The following lemma shows that all the candidate spherical projection  $c_k(\widehat{p}^j; y_{k,q})$  indeed lie on the same hemisphere.

LEMMA 4.2. The spherical projections  $c_k(\widehat{p}; y)$  defined by (4.2) with  $\widehat{\ell} = \widehat{\ell}_k(y), k \in [K]$  are always on the hemisphere

$$\mathcal{H} = \{ \boldsymbol{x} \middle| \|\boldsymbol{x}\| = \|\widehat{\boldsymbol{p}}\|, \langle \widehat{\boldsymbol{p}}, \boldsymbol{x} \rangle \ge 0 \}, \tag{4.15}$$

for any y and k.

*Proof.* Since  $\|c_k(\widehat{p}; y)\| = \|\widehat{p}\|$ , we only need to verify the second condition in (4.15). Using (4.2), a direct computation yields,

$$\langle \widehat{\boldsymbol{p}}, \boldsymbol{c}_k(\widehat{\boldsymbol{p}}; y) \rangle = \|\widehat{\boldsymbol{p}}\| \|\mathscr{P}_{c_k(y)}\widehat{\boldsymbol{p}}\| \ge 0.$$

All the above is almost sufficient to guarantee that the algorithm described by (4.14) has a unique solution. The last obstacle we have yet to overcome is to ensure that the points  $c_k(\widehat{p}^i; y_{k,q})$  are uniquely defined. To ensure this, we make the following assumption.

Assumption 4.1.  $\hat{\ell}_k(y)$  is not parallel to  $\hat{p}^j$  for all (k, y) pairs for  $y \in \omega_k^-$ .

Assumption 4.1 is necessary to ensure unique existence of the spherical projections  $c_k(\vec{p}'; y_{k,q})$ . Although we cannot yet theoretically justify of Assumption 4.1, in all of our numerical experiments, Assumption 4.1 holds. Under this assumption, we can prove uniqueness of the update (4.14).

Proposition 4.1. Under Assumption 4.1, the solution  $\hat{p}^{j+1}$  to (4.14) is unique.

*Proof.* Under Assumption 4.1 and Lemma 4.2, the candidate points  $c_k(\widehat{p}^j; y_{k,q})$  are all in the interior of  $\mathscr{H}$ . Therefore, from Theorem 4.2, the solution  $\widehat{p}^{j+1}$  to (4.14) is unique and furthermore lies in the interior of  $\mathscr{H}$ .

Algorithms that compute the average spherical projection by solving the optimization problem (4.14) can be adapted from [7, Algorithms A1 or A2) or [24, Equation 10).

## 4.4 A hybrid approach

We propose a final algorithm: the averaging algorithm of the previous section results in less oscillatory iterate trajectories, but moves relatively slowly. The algorithm in this section combines the ideas of the greedy and averaging approach. First we denote the greedy update (4.5) at the *j*th iteration by  $\widehat{p}_g^{j+1}$  and the average update (4.12) by  $\widehat{p}_a^{j+1}$ . The hybrid update we propose moves in the direction of the averaged update  $\widehat{p}_a^{j+1}$ , but with a distance defined by the greedy update  $\widehat{p}_g^{j+1}$ . Specifically, at the *j*th iteration,

- i. Compute the geodesic projection  $\hat{p}_g^{j+1}$  and the average projection  $\hat{p}_a^{j+1}$  via (4.5) and (4.14), respectively.
- ii. If  $d_r(\widehat{p}^j, \widehat{p}_g^{j+1})/r < 10^{-6}$ , i.e., the most violated constraint is very close to the current update, we simply perform the greedy update, setting  $\widehat{p}^{j+1} = \widehat{p}_g^{j+1}$ .
- iii. Otherwise, we compute  $\widehat{p}^{j+1}$  by moving  $\widehat{p}^{j}$  along the unique geodesic from  $\widehat{p}^{j}$  to  $\widehat{p}_{a}^{j+1}$  by a distance given by the intrinsic distance between  $\widehat{p}^{j}$  and  $\widehat{p}_{g}^{j+1}$ . Let  $\widetilde{p}^{j}$ ,  $\widetilde{p}_{a}^{j}$ , and  $\widetilde{p}_{g}^{j}$  be the unitnorm versions of  $\widehat{p}^{j}$ ,  $\widehat{p}_{a}^{j+1}$ , and  $\widehat{p}_{g}^{j+1}$ , respectively. Then the update we propose is

$$\widetilde{\boldsymbol{p}}^{j+1} = \exp_{\widetilde{\boldsymbol{p}}^j} \left( d(\widetilde{\boldsymbol{p}}^j, \widetilde{\boldsymbol{p}}_g^{j+1}) \boldsymbol{v} \right),$$

$$\widehat{\boldsymbol{p}}^{j+1} = \|\widehat{\boldsymbol{p}}^j\| \widetilde{\boldsymbol{p}}^{j+1}, \tag{4.16}$$

where v is the unit-speed velocity at the base point  $\tilde{p}^j$  of the geodesic segment leading to  $\tilde{p}_a^{j+1}$ , i.e.,  $v = \gamma'_{\tilde{p}'\tilde{p}_a^{j+1}}(0)$ .

We use the greedy update in step 2 above since in this case we are relatively close to the solution, and so typically the greedy procedure converges very quickly.

# 5. Numerical experiments

Throughout this section, we take M=N observations, and the observation functionals  $\{\phi_n\}_{n\in[N]}$  are chosen to be the projection functionals  $\phi_n(\cdot):=\langle\cdot,v_n\rangle$  onto the given subspace V. We denote the

unknown function by u, the H-best projection onto V by v, the norm-constrained solution (the solution to (2.14)) by  $v_{\rm NC}$ , and the linearly constrained solution by  $v_{\rm LC}$ . I.e.,  $v_{\rm LC}$  is the solution to (2.14) but with only the linear inequality constraints,

$$v_{\text{LC}} := \sum_{n \in [N]} w_n v_n, \qquad \qquad \mathbf{w} := \operatorname*{arg\,min}_{\widehat{\mathbf{v}} \in C^0} \frac{1}{2} \|\widehat{\mathbf{v}} - \widehat{\mathbf{p}}\|^2.$$

For other choices of observation functionals, e.g., pointwise observations (collocation-based approximations), our theory and algorithms can be generalized naturally.

For our univariate examples, we consider the Sobolev spaces on a general interval [a, b] as our Hilbert spaces H,

$$H^{q}([a,b]) := \left\{ u : [a,b] \to \mathbb{R} \middle| \|u\|_{H^{q}}^{2} < \infty \right\}, \qquad \|u\|_{H^{q}}^{2} := \sum_{j=0}^{q} \int_{a}^{b} \left[ u^{(j)}(x) \right] dx, \tag{5.1}$$

and choose the subspace V according to the choice of the pair (a, b),

if 
$$(a,b) = (-1,1)$$
, then  $V = V^{\text{poly}} := \text{span}\left\{ \{x^n\}_{n=0}^{N-1} \right\}$ ,  
if  $(a,b) = (0,\pi)$ , then  $V = V^{\cos} := \text{span}\left\{ \{\cos nx\}_{n=0}^{N-1} \right\}$ . (5.2)

We will test our algorithms for  $H^0(=L^2)$ ,  $H^1$ , and  $H^2$  using the linear constraint sets,

- (Positivity)  $U_0 := \{u \in H | u(x) \ge 0 \ \forall x \in [a, b]\}$
- (Monotonicity)  $U_1 := \{u \in H | u'(x) \ge 0 \ \forall x \in [a, b]\}$
- (Convexity)  $U_2 := \{ u \in H | u''(x) \ge 0 \ \forall x \in [a, b] \}$

Although our theoretical result in Theorem 3.1 does not guarantee the uniqueness of the solution when a boundedness constraint imposed, we still test our algorithms with imposing the constraint,

• (Boundedness)  $G_0 := \{u \in H | u(x) \le 1 \ \forall x \in [a, b]\},\$ 

for some of our tests. When a boundedness constraint is imposed, the hyperplane is an affine plane. In this case, we first project the current iteration to the affine hyperplane, then rescale the point with respect to the vertex  $\mathbf{r}_0$  of the cone (the projection of the origin onto the affine plane) to the sphere, i.e., the projection (4.4) is replaced by

$$\boldsymbol{c} = \mathscr{P}_H^s \widehat{\boldsymbol{p}} := \boldsymbol{r}_0 + \frac{\sqrt{\|\widehat{\boldsymbol{p}}\|^2 - \|\boldsymbol{r}_0\|^2}}{\sqrt{\|\mathscr{P}_H \widehat{\boldsymbol{p}}\|^2 - \|\boldsymbol{r}_0\|^2}} (\mathscr{P}_H \widehat{\boldsymbol{p}} - \boldsymbol{r}_0).$$

We will also introduce a metric to measure the change between the constrained solutions and the unconstrained solution:

$$\eta_* = \frac{\|v - v_*\|_H}{\|v - u\|_H},\tag{5.3}$$

where asterisk "\*" on the subscript of v can be either "LC" (the linearly constrained approximation using the dissipative formulation in [35]) or 'NC' (the non-dissipative procedure in this article). Since v - u is H-orthogonal to V, the Pythagorean theorem implies,

$$||v_* - u||_H = \sqrt{1 + \eta^2} ||v - u||_H^2.$$

The quantity  $\sqrt{1+\eta^2}$  can therefore be used to measure the error in a constrained solution relative to error in the unconstrained solution (which in this case is the *H*-best approximation from *V*). Values of  $\eta$  that are O(1) indicate that the error committed by the constrained solution is comparable to that of the unconstrained solution. It is also interesting to measure the difference between the linearly constrained solution and the norm-constrained solution  $\|v_{LC} - v_{NC}\|$ . In all of our experiments, the norm-constrained solution  $v_{NC}$  differs only slightly from the linearly-constrained solution  $v_{LC}$ .

Algorithm 1 is the greedy algorithm, but it is also the template for the average algorithm. To apply the average algorithm, one only needs to replace the update of  $\hat{p}$  with (4.14). In line 2 of Algorithm 1,  $\delta$  is set to be  $10^{-10}$ . In addition, we restrict the maximum number of iterations to be 10,000 to avoid infinite loops.

**Time complexity**: A *single* step of the greedy algorithm requires determining  $(y^*, k^*)$  associated with the farthest hyperplane (Algorithm 1), whose complexity is basis-dependent. For univariate polynomial approximation, the complexity is  $O(N^2)$ , where N is the dimension of the subspace V [35]. A *single* step of the average algorithm requires computing the integrals and the Karcher mean ((4.14)). The cost of computing the integrals consists of two parts,

- determining the negative region  $\mathcal{D}_j^-$ : this part requires finding the zeros of the unconstrained approximation. For univariate polynomial approximation, the total complexity is  $O(KN^2)$ .
- applying the quadrature rule: let  $Q = \max_{k \in [K]} Q_k$ . Since the number of negative regions is at most O(N) and computing  $d_r(\mathbf{x}, \mathbf{c}_k(\widehat{\mathbf{p}}^j; y_{k,q}))$  is of O(N), the total cost is  $O(QN^2)$ .

For the computation of Karcher mean, we use the algorithm of locally linear rate of convergence ([7], Algorithm A1). Let  $\mu$  be its rate of convergence and  $\epsilon$  be a desired tolerance level, the complexity is  $O(\log_{\mu} \epsilon)$  and is independent of the choice of the basis functions.

The time complexity for finding a feasible solution to the original optimization problem remains unknown since we do not know how many steps the algorithms take to achieved a desired tolerance level. The algorithms converge quickly in some examples but slowly on others (e.g., Table 2). In our numerical test, the greedy algorithm performs better than the average algorithm.

#### 5.1 Performance comparison of algorithms

In this section, we present the comparison of solutions from the linearly-constrained optimization in [35] and the norm-constrained optimization proposed in our work. We consider the case (a, b) = [-1, 1] and degree-(N-1) polynomial approximations with  $V = V^{\text{poly}}$ . The test functions are chosen as a step function and its antiderivatives:

$$u_{j+1}(x) := c_{j+1} \int_{-1}^{x} u_j(t)dt, \qquad u_0(x) = \begin{cases} 0, & x \le 0, \\ 1, & x > 0, \end{cases}$$
 (5.4)

Table 1 Comparison of number of iterations (I) and relative errors  $(\eta)$  on the test function  $u = u_2$  with positivity-constraint imposed for different values of N and different algorithms. The ambient Hilbert space is  $H = L^2([-1, 1])$ .

	N=6		N = 31	
	I	η	I	η
Greedy	17	1.1479	23	0.9859
Average	87	1.1483	220	0.9856
Hybrid	15	1.1494	22	0.9885

Table 2 Comparison of the minimum values of  $v_{NC}$ ,  $v'_{NC}$ , and  $v''_{NC}$  using average approach at 10,000 iterations. 'Converge' indicates the corresponding procedure finds a feasible solution before the maximum number of iterations is reached.

	N = 6			N = 31		
	$v_{NC}$	$v_{NC}'$	v" <sub>NC</sub>	$v_{NC}$	$v'_{NC}$	$v_{NC}^{\prime\prime}$
$H^0$	-1.05e-6	9.35e-5	-1.86e-4	-2.18e-7	-2.35e-3	-3.06e-2
$H^1$	converge	converge	converge	-1.33e-8	-3.84e-7	-1.51e-3
$H^2$	converge	converge	converge	converge	converge	converge

where  $c_{j+1}$  is a normalized constant that ensures the  $u_{j+1}(1) = 1$ . We provide a summary of the performance of our three proposed approaches for norm-preserving optimization Table 1. We observe that, compared to greedy approach, there is a slight decrease in the number of iterations. Both greedy approach and the hybrid approach are much faster than the average approach. The relative error of the three proposed approaches are comparable.

#### 5.2 Polynomial space approximation example

In this section, we continue to consider the approximation using (a,b) = (-1,1) and  $V = V^{\text{poly}}$ . In our first experiment, we test the capability of our algorithms for approximating the step function  $u_0(x)$  in (5.4) and the similarity between  $v_{\text{LC}}$  and  $v_{\text{NC}}$ . We compute the approximation for N = 6 and N = 31 Figure 2 and consider three choices of linear constraint sets  $E^0$  introduced in (2.6):

- (i) (positivity)  $E^0 = U_0$ ,
- (ii) (positivity and monotonicity)  $E^0 = U_0 \cap U_1$ , and
- (iii) (positivity, monotonicity, boundedness)  $E^0 = U_0 \cap U_1 \cap G_0$ .

We note that, for (i) and (ii), the norm-constrained solutions are simply slight adjustment of the linearly constrained solution, both visually and quantitatively. The discrepancy is more obvious in case (iii). Increasing the order in the polynomial approximation further decreases the discrepancy between  $v_{\rm LC}$  and  $v_{\rm NC}$ . We also note that, the Gibbs'-type oscillations presented in the left column of Figure 2 can be alleviated by enforcing the monotonicity and the boundedness constraint. All computed  $\eta_{\rm NC}$  values are order 1, which shows that our norm-preserving approximations are comparable to the H-best approximation.

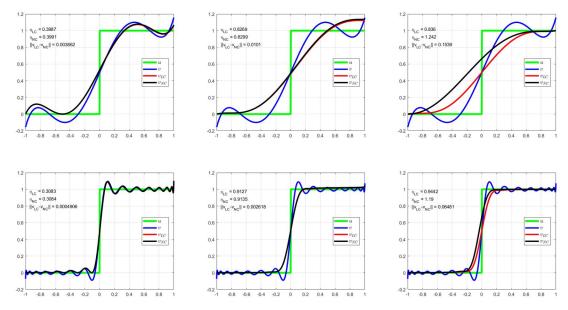


Fig. 2. Greedy algorithm results: comparison of different methods: degree 5 polynomial positivity-preserving approximation to the step function for different constraints and different polynomial spaces. Left: constraint  $U_0$ . Middle: constraint  $U_0 \cap U_1$ . Right:  $U_0 \cap U_1 \cap G_0$ . Top:  $N = \dim V = 6$ . Bottom:  $N = \dim V = 31$ .

In the second experiment of this section, we investigate how the choice of ambient Hilbert space H affects the accuracy of the approximation. We approximate the function  $u_2(x)$  with linear constraint set  $E^0 = U_0 \cap U_1 \cap U_2$  for N = 6 and N = 31 on different Hilbert spaces  $H = H^0, H^1$ , and  $H^2$ . We observe relatively large values of both  $\eta_{NC}$  and  $\eta_{LC}$ , but increasing the regularity of the Hilbert space and/or increasing the order of the polynomial can reduce these relative errors. Similar to the previous test, the discrepancy between  $v_{LC}$  and  $v_{NC}$  decreases as the order of polynomial order increases. It increases as the complexity of the linear constraint set increases. Nevertheless, both approximations are qualitatively good for N = 31. The results are shown in Figure 2.

Quantitatively, we find that the minimum values of  $v_{NC}$ ,  $v_{NC}'$ , and  $v_{NC}''$  converge slowly for some examples with less regular Hilbert space  $H = H^0$ ,  $H^1$ . Among our three proposed approach, the greedy approach and the hybrid approach performs slightly better than the average approach. For different choices of the ambient Hilbert spaces, we report in Table 2 the minimum values of  $v_{NC}$ ,  $v_{NC}'$  and  $v_{NC}''$  using average approach at 10,000 iterations. The 'converge' in Table 2 indicates that the procedures achieve the desired tolerance levels. We note that, by increasing the regularity of the ambient Hilbert space, our procedures can identify a feasible solution much faster. We emphasize that the simplicity of this example belies the complexity and difficulty of the geometry of the problem, which is evidenced by algorithms requiring more iterations to complete. In the remaining examples of this paper, all the algorithms identify an element of the feasible set (to within precision tolerances).

# 5.3 Convergence rate

In this subsection, we compare the rates of convergence between the unconstrained solution and the norm-constrained solution. We consider  $u = u_0(x)$  and  $u = u_2(x)$  using  $V = V^{\text{poly}}$ . The ambient Hilbert

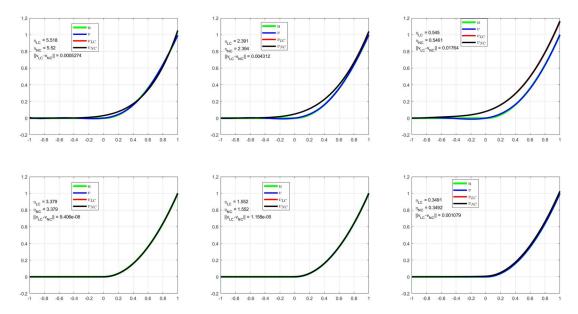


Fig. 3. Comparison of different approximations to  $u_2(x)$  for different ambient Hilbert spaces and different polynomial spaces. The red curves  $(v_{LC})$  are covered by the black curves  $(v_{NC})$ . The constraint is  $U_0 \cap U_1 \cap U_2$ . Left:  $H = H^0$ . Middle:  $H = H^1$ . Right:  $H = H^2$ . Top:  $N = \dim V = 6$ . Bottom:  $N = \dim V = 31$ .

space is  $H = L^2([-1, 1])$ . We compute the rate of convergence on the constrained sets  $U_0$ ,  $U_0 \cap U_1$ , and  $U_0 \cap U_1 \cap G_0$ . We observe from Figure 4 that our norm-constrained solutions have a similar rate of convergence to the unconstrained ( $H = L^2$ -optimal) solution u (even when a boundedness constraint is imposed).

# 5.4 M-shape function using cosine basis

In this section, we will choose  $V = V^{\cos}$  for approximating an M-shape function defined on  $[0, \pi]$ ,

$$u(x) = \begin{cases} -\left(x - \frac{\pi}{8}\right)\left(x - \frac{\pi}{2}\right) & \frac{\pi}{8} \le x < \frac{\pi}{2}, \\ -\left(x - \frac{\pi}{2}\right)\left(x - \frac{7\pi}{8}\right) & \frac{\pi}{2} \le x < \frac{7\pi}{8}, \\ 0 & \text{otherwise,} \end{cases}$$
(5.5)

with positivity constraint  $U_0$  imposed. For a cosine polynomial, the difficult part for applying our algorithm is to determine the y-parameter corresponding to the most violated constraint (or the negative y-region), which requires to find the zeros of a trigonometry polynomials. Fortunately, this difficulty can be resolved by taking advantage of the Chebyshev polynomials. The results are shown in Figure 5.

# 5.5 Two-dimensional cylinder indicator function

In our last example, we consider the approximation to a cylinder

$$u(x,y) = \begin{cases} 1 & \text{if } \sqrt{(x-0.5)^2 + (y-0.5)^2} < 0.5, \\ 0 & \text{otherwise.} \end{cases}$$
 (5.6)

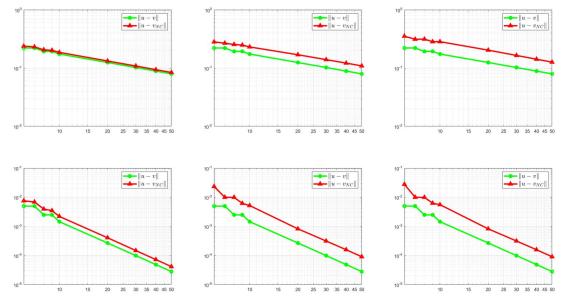


Fig. 4. Rate of convergence. Approximations to  $u = u_0(x)$  (top) and  $u = u_2(x)$  (bottom) with  $U_0$  (left),  $U_0 \cap U_1$  (middle), and  $U_0 \cap U_1 \cap G_0$  (right) imposed. The x-axis indicates the dimension of the polynomial space  $V = V^{\text{poly}}$ . The ambient Hilbert space is  $H = L^2([-1, 1])$ .

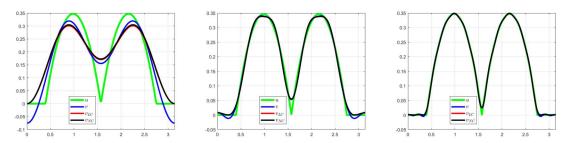


Fig. 5. Comparison of the approximations to (5.5) with different N. Constraint:  $U_0$ , positivity-preserving. From left to right: N = 6, 16, 31. The red curve is covered by the black curve.

The computational domain is  $[-1,1] \times [-1,1]$ , and the polynomial space is the tensor product space  $V^{\mathrm{poly}} \otimes V^{\mathrm{poly}}$ , where N is chosen to be 15. The positivity constraint  $U_0$  is imposed. The computation requires to find the global minimum of a two-dimensional nonconvex function (4.9) (see also (2.10)–(2.11)). We use MATLAB's optimization function fmincon, using the sequential quadratic programming option, and approximate the global minimum by solving the optimization with several randomly initialized starting points. The constraints we set for fmincon are the boundaries for the computational domain.

The results are shown in Figure 6. The numerical results show that our norm-constrained approximation can preserve both the positivity and the norm by 'correcting' the linearly constrained solution. The function is entirely non-negative for both the linearly constrained solution (bottom middle, Figure 6) and the norm-constrained solution (bottom right, Figure 6).

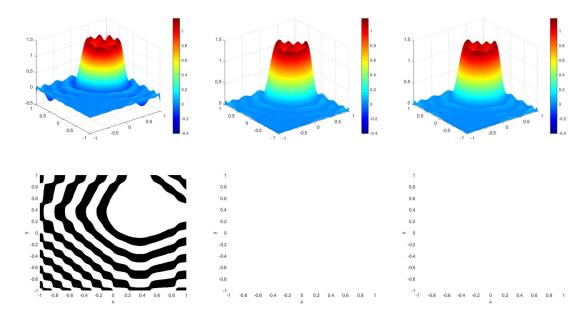


Fig. 6. Comparison of different approximations to (5.6), greedy procedure. Constraint:  $U_0$ , positivity-preserving. Top: mesh plot. Bottom: negative region indicator function, where the black region represent the region where the approximation is negative. Left: unconstrained solution u. Middle: linearly constrained solution  $v_{LC}$ . Right: Norm-constrained solution  $v_{NC}$ .  $\eta_{LC} = 0.1229$ ,  $\eta_{NC} = 0.1230$ ,  $||v_{LC} - v_{NC}|| = 0.0030$ .

## Data availability statement

No new data were generated or analyzed in support of this research.

#### **Funding**

D. Dai and A. Narayan are supported by NSF DMS-1848508 and AFOSR FA9550-20-1-0338.

# **Conflict of interest**

There is no conflict of interest.

#### REFERENCES

- ALLEN, L. & KIRBY, R. C. (2021) Bounds-constrained polynomial approximation using the Bernstein basis. arXiv preprint arXiv:2104.11819.
- 2. Beatson, R. (1982) Restricted range approximation by splines and variational inequalities. *SIAM J. Numer. Anal.*, **19**(2), 372–380.
- 3. Beatson, R. K. (1978) The degree of monotone approximation. *Pacific J. Math.*, **74**(1), 5–14.
- 4. Berzins, M. (2007) Adaptive polynomial interpolation on evenly spaced meshes. SIAM Rev., 49(4), 604–627.
- Besse, C., Descombes, S., Dujardin, G. & Lacroix-Violet, I. (2021) Energy-preserving methods for nonlinear Schrödinger equations. IMA J. Numer. Anal., 41(1), 618–653.
- 6. BOYD, S. P. & VANDENBERGHE, L. (2004) Convex Optimization. Cambridge University Press.
- 7. Buss, S. R. & Fillmore, J. P. (2001) Spherical averages and applications to spherical splines and interpolation. *ACM Transactions on Graphics*, **20**(2), 95–126.

- 8. CAMPOS-PINTO, M., CHARLES, F. & DESPRÉS, B. (2019) Algorithms for positive polynomial approximation. *SIAM J. Numer. Anal.*, **57**(1), 148–172.
- 9. CANDES, E., ROMBERG, J. & TAO, T. (2006) Robust uncertainty principles: Exact signal reconstruction from highly incomplete frequency information. *IEEE Trans. Inform. Theory*, **52**(2), 489–509.
- 10. CELLEDONI, E., McLACHLAN, R. I., McLAREN, D. I., OWREN, B., QUISPEL, G. R. W. & WRIGHT, W. M. (2009) Energy-preserving Runge-Kutta methods. *ESAIM Math. Model. Numer. Anal.*, **43**(4), 645–649.
- 11. Cheney, W. & Goldstein, A. A. (1959) Proximity maps for convex sets. *Proc. Amer. Math. Soc.*, 10(3), 448–450.
- 12. COHEN, A., DAHMEN, W. & DEVORE, R. (2009) Compressed sensing and best k-term approximation. *J. Amer. Math. Soc.*, **22**(1), 211–231.
- DESPRÉS, B. (2017) Polynomials with bounds and numerical approximation. Numer. Algorithms, 76(3), 829– 859.
- 14. DEVORE, R. A. (1974) Degree of Monotone Approximation. in Linear Operators and Approximation II / Lineare Operatoren und Approximation II: Proceedings of the Conference held at the Oberwolfach Mathematical Research Institute, Black Forest, March 30–April 6, 1974 / Abhandlungen zur Tagung im Mathematischen Forschungsinstitut Oberwolfach, Schwarzwald, vom 30. März bis 6. April 1974, ed. by P. L. BUTZER, & B. SZOŐKEFALVI-NAGY, pp. 337–351. Birkhäuser Basel, Basel.
- 15. Donoho, D. (2006) Compressed sensing. *IEEE Trans. Inform. Theory*, **52**(4), 1289–1306.
- 16. Ferreira, O., Iusem, A. & Németh, S. (2014) Concepts and techniques of optimization on the sphere. *Top*, **22**(3), 1148–1170.
- 17. Ferreira, O. P., Iusem, A. N. & Németh, S. Z. (2013) Projections onto convex sets on the sphere. J. Global Optim., 57(3), 663–676.
- 18. GANDER, W., GOLUB, G. H. & VON MATT, U. (1989) A constrained eigenvalue problem. *Linear Algebra Appl.*, **114**, 815–839.
- GOBERNA, M. Á. & LÓPEZ, M. A. (2013) Semi-Infinite Programming: Recent Advances. Springer Science & Business Media.
- 20. HAGER, W. W. (2001) Minimizing a quadratic over a sphere. SIAM J. Optim., 12(1), 188–208.
- 21. HAIRER, E. (2010) Energy-preserving variant of collocation methods. *JNAIAM. J. Numer. Anal. Ind. Appl. Math.*, **5**, 73–84.
- 22. HAUCK, C. & McCLARREN, R. (2010) Positive P<sub>N</sub> closures. SIAM J. Sci. Comput., **32**(5), 2603–2626, Publisher: Society for Industrial and Applied Mathematics.
- 23. HETTICH, R. & KORTANEK, K. O. (1993) Semi-infinite programming: Theory, methods, and applications. *SIAM Rev.*, **35**(3), 380–429.
- 24. Krakowski, K. A., Hüper, K. & Manton, J. H. (2007) On the computation of the karcher mean on spheres and special orthogonal groups. *RoboMat 2007, Workshop on Robotics and Mathematics*. Coimbra, Portugal: Centro Internacional de Matemática, pp. 119–124.
- LAIU, M. P., HAUCK, C. D., McCLARREN, R. G., O'LEARY, D. P. & TITS, A. L. (2016) Positive filtered P<sub>N</sub> moment closures for linear kinetic equations. SIAM J. Numer. Anal., 54(6), 3214–3238.
- 26. Nesterov, Y. (2000) Squared Functional Systems and Optimization. ( H. Frenk, K. Roos, T. Terlaky & S. Zhang eds). US, Boston, MA: Springer, pp. 405–440.
- 27. NIE, J. & DEMMEL, J. W. (2006) Shape Optimization of Transfer Functions. *Multiscale Optimization Methods and Applications*. (W. W. HAGER, S.-J. HUANG, P. M. PARDALOS & O. A. PROKOPYEV eds). US, Boston, MA: Springer, pp. 313–326.
- 28. Nochetto, R. & Wahlbin, L. (2002) Positivity preserving finite element approximation. *Math. Comp.*, **71**(240), 1405–1419.
- 29. QUISPEL, G. & McLaren, D. I. (2008) A new class of energy-preserving numerical integration methods. *J. Phys. A*, **41**(4), 045206.
- 30. RENDL, F. & WOLKOWICZ, H. (1997) A semidefinite framework for trust region subproblems with applications to large scale minimization. *Math. Programming*, **77**(1), 273–299.

- 31. Sorensen, D. C. (1997) Minimization of a large-scale quadratic function subject to a spherical constraint. *SIAM J. Optim.*, 7(1), 141–161.
- 32. Stein, O. (2012) How to solve a semi-infinite optimization problem. *European J. Oper. Res.*, **223**(2), 312–320.
- 33. Trefethen, L. N. (2012) Approximation Theory and Approximation Practice. SIAM.
- 34. VANDENBERGHE, L. & BOYD, S. (1996) Semidefinite programming. SIAM Rev., 38(1), 49–95.
- 35. ZALA, V., KIRBY, M. & NARAYAN, A. (2020) Structure-preserving function approximation via convex optimization. SIAM J. Sci. Comput., 42(5), A3006–A3029.
- 36. ZHANG, X. & SHU, C.-W. (2011) Maximum-principle-satisfying and positivity-preserving high-order schemes for conservation laws: survey and new developments. *Proceedings of the Royal Society A: Mathematical, Physical and Engineering Sciences*, **467**(2134), 2752–2776.