

Kernel-based method for large empirical truncated moment problem

MICHAEL MULTERER

PAUL SCHNEIDER

ROHAN SEN

USI Lugano, Switzerland

YAMC 2022, Arenzano

September 21, 2022

TABLE OF CONTENTS

INTRODUCTION

TRUNCATED MOMENT PROBLEM

EMPIRICAL PROBLEM

METHODOLOGY

FIGURES

INTRODUCTION

MOTIVATION

- Quadrature rules approximate integrals through a small number of *nodes* and *weights* pertaining to a discrete probability measure.
- These nodes parsimoniously describe the important states or *scenarios* that are the best low-dimensional representation of the underlying complicated distribution.
- These scenarios reconcile moment matrices that often feature in many applied situations.
- Our goal is to extract these low-dimensional scenarios and their probabilities from large and high-dimensional datasets.

CONTRIBUTION

- The extant algorithms do not scale well when solving large and high-dimensional problems and also suffer from numerical instability.
- We propose algorithms that are tractable and computationally efficient at the same and are founded on the intersection of the truncated moment problem from probability theory and reproducing kernel Hilbert spaces.
- We propose a novel approach for the extraction of the scenarios and their probabilities for the specialized case of covariance matrices of high-dimensional random variables.
- We also modify Lasserre's algorithm for multivariate Gaussian quadrature that partially remedies its numerical instability and significantly improves its computational complexity.

Notations

- $\mathcal{S}(\mathcal{H})$: space of symmetric matrices on \mathcal{H}
- $\mathbf{A} \succeq 0$: \mathbf{A} is positive semi-definite matrix
- $\mathbf{A} \in \mathcal{S}_+^N$: $\mathbf{A} \in \mathbb{R}^{N \times N}$ is symmetric and positive semi-definite
- $\mathcal{P}_t(\Omega)$: space of multivariate polynomials on Ω of maximum degree t
- $s(t) : \binom{d+t}{t}$
- $\|\cdot\|_F$: Frobenius norm
- $\|\cdot\|_\star$: trace norm

TRUNCATED MOMENT PROBLEM

MOMENT SEQUENCE AND LINEAR FORM

Truncated sequence (in d variables and of degree t):

$$\mathbf{y} = (y_{\alpha}) \text{ where } \alpha \in \overline{\mathbb{N}}_t^d := \left\{ (\alpha_1, \dots, \alpha_d) : \alpha_i \in \overline{\mathbb{N}}, |\alpha_1 + \dots + \alpha_d| \leq t \right\}$$

Truncated moment sequence (tms):

$$y_{\alpha} = \int_{\Omega} \mathbf{x}^{\alpha} d\mu(\mathbf{x}) \text{ where } \mathbf{x}^{\alpha} := x_1^{\alpha_1} \cdots x_d^{\alpha_d}$$

Monomial basis for $\mathcal{P}_t = \text{span} \{ \mathbf{x}^{\alpha} : |\alpha| \leq t \}$:

$$\boldsymbol{\tau}_t(\mathbf{x}) := [1, x_1, \dots, x_d, \dots, x_1^t, \dots, x_d^t] \in \mathbb{R}^{s(t)}$$

Riesz functional: Given $\mathbf{y} = (y_{\alpha})$ define $\mathcal{L}_{\mathbf{y}} \in (\mathbb{R}[\mathbf{x}])^*$ as:

$$\mathcal{L}_{\mathbf{y}}(p) := \sum_{\alpha} p_{\alpha} y_{\alpha} \text{ for } p = \sum_{\alpha} p_{\alpha} \mathbf{x}^{\alpha}$$

TRUNCATED MOMENT PROBLEM

The truncated moment problem:

Given a truncated sequence y , does there exist a representing measure μ and if so, how to obtain it.

A crucial fact:

Every truncated moment sequence has a representing measure that is a convex combination of at most $s(t) = \binom{d+t}{t}$ many Dirac measures.

Relation with quadrature rules:

Finding measures with a small number of atoms is the equivalent to the problem of finding quadrature rules.

MOMENT MATRIX

Moment matrix: Given $\mathbf{y} = (y_\gamma)_{|\gamma| \leq 2t}$, define $M_t(\mathbf{y}) \in \mathbb{R}^{s(t) \times s(t)}$ as:

$$M_t(\mathbf{y})_{\alpha, \beta} := \mathcal{L}_y(\boldsymbol{\tau}_t \boldsymbol{\tau}_t^\top)_{\alpha, \beta} = \mathcal{L}_y(x^{\alpha+\beta}) = y_{\alpha+\beta}$$

Example for $d = t = 2$:

$$M_2(\mathbf{y}) = \begin{bmatrix} y_{00} & y_{10} & y_{01} & y_{20} & y_{11} & y_{02} \\ y_{10} & y_{20} & y_{11} & y_{30} & y_{21} & y_{12} \\ y_{01} & y_{11} & y_{02} & y_{21} & y_{12} & y_{03} \\ y_{20} & y_{30} & y_{21} & y_{40} & y_{31} & y_{22} \\ y_{11} & y_{21} & y_{12} & y_{31} & y_{22} & y_{13} \\ y_{02} & y_{12} & y_{03} & y_{22} & y_{13} & y_{04} \end{bmatrix}$$

For $\tilde{\mathbf{y}} = (\tilde{y}_\alpha)_{\alpha \in \overline{\mathbb{N}}^d}$, we have $M(\tilde{\mathbf{y}})_{\alpha, \beta} = \tilde{y}_{\alpha+\beta}$ for $\alpha, \beta \in \overline{\mathbb{N}}^d$

FLAT EXTENSION

Let X be a symmetric matrix with block form

$$X = \begin{pmatrix} A & B \\ B^\top & C \end{pmatrix}$$

X is called a *flat extension* of A if

$$\text{rank } X = \text{rank } A$$

If X is a flat extension of A , then $X \succeq 0 \iff A \succeq 0$.

Flat extension theorem: For $y = (y_\alpha)_{|\alpha| \leq 2t}$, if $M_t(y)$ is a flat extension of $M_{t-1}(y)$, then there exists a (unique) sequence $\tilde{y} = (\tilde{y}_\alpha)_{\alpha \in \mathbb{N}^d}$ for which $M(\tilde{y})$ is flat extension of $M_t(y)$.

FINITE ATOMIC REPRESENTING MEASURES

Theorem 1: $\tilde{\mathbf{y}}$ has a unique representing measure μ which is r -atomic with

$$\text{supp}(\mu) = \mathcal{V}_{\mathbb{R}}\left(\text{Ker } \mathbf{M}(\tilde{\mathbf{y}})\right) \subseteq \mathbb{R}^d$$

$$\Updownarrow$$

$$\mathbf{M}(\tilde{\mathbf{y}}) \succeq 0 \quad \text{and} \quad \text{rank } \mathbf{M}(\tilde{\mathbf{y}}) = r$$

Theorem 2: \mathbf{y} has a unique representing measure μ which is r -atomic with

$$\text{supp}(\mu) = \mathcal{V}_{\mathbb{R}}\left(\langle \text{Ker } \mathbf{M}_t(\mathbf{y}) \rangle\right) \subseteq \mathbb{R}^d$$

$$\Updownarrow$$

$$\mathbf{M}_t(\mathbf{y}) \succeq 0 \quad \text{and} \quad \text{rank } \mathbf{M}_t(\mathbf{y}) = \text{rank } \mathbf{M}_{t-1}(\mathbf{y}) = r$$

Scenario representation:

$$\mu = \sum_{i=1}^r p_i \delta_{\xi_i} \quad \text{where } \Xi := \{\xi_1, \dots, \xi_r\} = \mathcal{V}_{\mathbb{R}}\left(\langle \text{Ker } \mathbf{M}_t(\mathbf{y}) \rangle\right)$$

SCENARIOS

Moment matrix representation:

$$M_t(\mathbf{y}) = \sum_{i=1}^r p_i \boldsymbol{\tau}_t(\boldsymbol{\xi}_i) \boldsymbol{\tau}_t(\boldsymbol{\xi}_i)^\top$$

Vandermonde form:

$$M_t(\mathbf{y}) = \mathbf{V}_t(\boldsymbol{\Xi}, \boldsymbol{\tau})^\top \mathbf{D} \mathbf{V}_t(\boldsymbol{\Xi}, \boldsymbol{\tau}) \quad \text{with} \quad \mathbf{D} := \text{diag}(p_1, \dots, p_r)$$

Vandermonde matrix:

$$\mathbf{V}_t(\boldsymbol{\Xi}, \boldsymbol{\tau}) = \left[\boldsymbol{\tau}_t(\boldsymbol{\xi}_1), \dots, \boldsymbol{\tau}_t(\boldsymbol{\xi}_r) \right]^\top \in \mathbb{R}^{r \times s(t)}$$

Gaussian Quadrature: *Lasserre's algorithm for finding a finite atomic representing measure coincides with that of constructing a quadrature rule with minimal number of nodes, known as Gaussian quadrature.*

MODIFICATIONS

Flat extension: Obtained via an SDP that is a trace minimization problem.

Numerical rank computation: No guarantee of convergence of the SDP whose size grows exponentially fast, and also necessitates the numerical rank computation of the input moment matrix.

Pivoted Cholesky decomposition: Circumvents the computational cost of the usual Cholesky and can be made rank-revealing, also performs Gaussian elimination in a numerically most favorable way.

MODIFIED LASSERRE'S ALGORITHM

Algorithm 3 (extraction algorithm) :

- **Input:** The moment matrix $\mathbf{M}_t(\mathbf{y})$ with $\text{rank } \mathbf{M}_t(\mathbf{y}) = r$
- **Output:** The r nodes $\Xi = [\xi_1, \dots, \xi_r]$

1: Perform pivoted Cholesky decomposition to get

$$\mathbf{P}^\top \mathbf{M}_t(\mathbf{y}) \mathbf{P} = \mathbf{P}^\top \mathbf{V}_t^\top \mathbf{D} \mathbf{V}_t \mathbf{P} = \tilde{\mathbf{V}}_t^\top \mathbf{D} \tilde{\mathbf{V}}_t = \mathbf{L} \mathbf{L}^\top$$

2: Reduce \mathbf{L} to an echelon form $\tilde{\mathbf{L}}$.

3: Extract from $\tilde{\mathbf{L}}$ the multiplication matrices $\mathbf{N}_i, i = 1, \dots, d$.

4: Compute $\mathbf{N} := \sum_{i=1}^d \rho_i \mathbf{N}_i$ with random convex combination.

5: Compute the Schur decomposition $\mathbf{N} = \mathbf{Q} \mathbf{T} \mathbf{Q}^\top$ with $\mathbf{Q} = [\mathbf{q}_1 \cdots \mathbf{q}_r]$.

6: Extract $\xi_j(i) = \mathbf{q}_j^\top \mathbf{N}_i \mathbf{q}_j, i = 1, \dots, d; j = 1, \dots, r$.

LEAST-SQUARES WEIGHTS

Algorithm 4 (least-squares weights) :

- **Input:** The moment matrix $\mathbf{M}_t(\mathbf{y})$ with $\text{rank } \mathbf{M}_t(\mathbf{y}) = r$
- **Output:** The r probability weights $\mathbf{p} = [p_1, \dots, p_r]$

- 1: Compute the generalized inverse (for instance using SVD) \mathbf{V}^\dagger of $\mathbf{M}_t(\mathbf{y})$
- 2: Compute $\mathbf{m}_y := \text{diag}\left((\mathbf{V}^\dagger)^\top \mathbf{M}_t(\mathbf{y}) \mathbf{V}^\dagger\right)$
- 3: With $\mathbf{p} := \text{diag}(\mathbf{D})$ perform the minimization with :

$$\underset{\mathbf{p} \in \mathbb{R}_+^r}{\text{minimize}} \quad \|\mathbf{m}_y - \mathbf{p}\|_2^2$$

$$\text{subject to: } \mathbf{1}^\top \mathbf{p} = 1$$

EMPIRICAL PROBLEM

EMPIRICAL TRUNCATED MOMENT PROBLEM

Training sample (drawn from μ):

$$\mathcal{X} = \left\{ \tilde{x}_1, \dots, \tilde{x}_N \right\} \subset \Omega \subseteq \mathbb{R}^d \text{ with } \Omega \text{ Hausdorff and locally compact}$$

Empirical probability measure:

$$\hat{\mu} = \frac{1}{N} \sum_{i=1}^N \delta_{\tilde{x}_i}$$

Empirical moment sequence:

$$\hat{y}_{\alpha} = \int_{\Omega} x^{\alpha} d\hat{\mu}(x) \text{ where } |\alpha| \leq 2t$$

Empirical moment matrix:

$$\widehat{M}_t = \frac{1}{N} \sum_{i=1}^N \tau_t(\tilde{x}_i) \tau_t(\tilde{x}_i)^{\top} \in \mathbb{R}^{s(t) \times s(t)}$$

MEASURE COMPRESSION

Target compressed measure:

$$\tilde{\mu} = \sum_{i=1}^r p_i \delta_{\xi_i} \text{ where } \Xi := \{\xi_1, \dots, \xi_r\} \subset \mathcal{X} \text{ with } r \ll N$$

Model moment matrix:

$$\widetilde{M}_t = \sum_{i=1}^r p_i \tau_t(\xi_i) \tau_t(\xi_i)^\top = V_t(\Xi, \tau)^\top \text{diag}(p) V_t(\Xi, \tau) \in \mathbb{R}^{s(t) \times s(t)}$$

Optimization Problem:

$$\underset{\substack{\xi_i \in \mathbb{R}^d, i=1, \dots, r \\ p \in \mathbb{R}^r}}{\operatorname{argmin}} \quad \left\| \widehat{M}_t - V_t(\Xi, \tau)^\top \text{diag}(p) V_t(\Xi, \tau) \right\|_F \quad (1)$$

$$\text{subject to: } \mathbf{1}^\top p = 1$$

Observations

- (1) is a non-convex optimization problem in general (except when $t = 1$) in the nodes ξ_i and the probability weights p_i for $i = 1, \dots, r$.
- The problem of matching the empirical moments with the model moments is in general, under-determined since we have $r \leq s(t) \ll N$.
- The Vandermonde representation of the positive semi-definite $\widetilde{\mathbf{M}}_t$ pertaining to the minimal generating measure is precisely in the optimal form as in the model proposed by Bach, Rudi et. al. within the RKHS framework.
- We develop a relaxation of the optimization problem (1) and reformulate it within the RKHS framework, which makes it convex.

METHODOLOGY

REPRODUCING KERNEL HILBERT SPACE (RKHS)

Let $(\mathcal{H}, \langle \cdot, \cdot \rangle_{\mathcal{H}}) \subseteq \mathbb{R}^{\Omega}$ be a separable Hilbert space of functions with

$$\Omega \subseteq \mathbb{R}^d.$$

Then \exists a unique reproducing kernel $k : \Omega \times \Omega \longrightarrow \mathbb{R}$ such that:

$$\forall x \in \Omega, \quad k_x := k(x, \cdot) \in \mathcal{H}$$

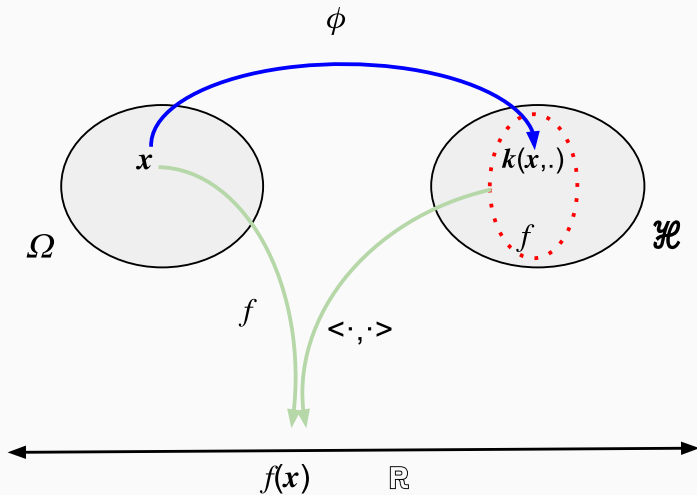
$$\forall f \in \mathcal{H}, \quad f(x) = \langle f, k_x \rangle_{\mathcal{H}} \quad \forall x \in \Omega$$

k is a *symmetric and positive definite kernel* i.e. for any finite

$$\mathcal{X} = \{\tilde{x}_1, \dots, \tilde{x}_N\} \subset \Omega, \quad K := \left[k(\tilde{x}_i, \tilde{x}_j) \right]_{i,j=1}^N \in \mathcal{S}_+^N$$

$$\text{Corollary: } \forall x, \tilde{x} \in \Omega, \quad \langle k_x, k_{\tilde{x}} \rangle_{\mathcal{H}} = \langle \phi(x), \phi(\tilde{x}) \rangle_{\mathcal{H}}$$

INTUITION



FRAMEWORK

Model:

$$f_A(\mathbf{x}) = \phi(\mathbf{x})^\top \mathbf{A} \phi(\mathbf{x}), \quad \mathbf{A} \in \mathcal{S}_+(\mathcal{H}) \quad (2)$$

Objective Function:

$$\inf_{\mathbf{A} \in \mathcal{S}_+(\mathcal{H})} L\left(f_A(\xi_1), \dots, f_A(\xi_r)\right) + \underbrace{\lambda_1 \|\mathbf{A}\|_\star + \lambda_2 \|\mathbf{A}\|_F^2}_{\tilde{\Omega}(\mathbf{A})}, \quad \lambda_2 > 0 \quad (3)$$

Representer Theorem: *Let L be lower semi-continuous, bounded below and convex, and $\tilde{\Omega}(\mathbf{A})$ be as above. Then (3) has a unique minimizer*

$$\mathbf{A}_* = \sum_{i=1}^r \sum_{j=1}^r B_{ij} \phi(\xi_i) \phi(\xi_j)^\top \quad \mathbf{B} \in \mathbb{R}^{r \times r}, \quad \mathbf{B} \succeq 0 \quad (4)$$

SETUP

RKHS:

$$(\mathcal{H}, \langle \cdot, \cdot \rangle_{\mathcal{H}}) := \text{span}\{x^{\alpha} : x \in \Omega, |\alpha| \leq t\} = \mathcal{P}_t(\Omega)$$

$$\text{with } L^2_{\mu}(\Omega) \text{ inner product } \langle f, g \rangle_{\mathcal{H}} = \int_{\Omega} f(x) g(x) d\mu(x)$$

Monomial basis:

$$\tau_t(x) := [1, x_1, \dots, x_d, \dots, x_1^t, \dots, x_d^t]^{\top} \in \mathbb{R}^{s(t)}$$

Gram matrix:

$$[G_t]_{\alpha, \beta} = \langle \tau_t, \tau_t^{\top} \rangle_{\mathcal{H}} = \int_{\Omega} x^{\alpha+\beta} d\mu(x) \quad |\alpha|, |\beta| \leq t$$

Reproducing Kernel:

$$k_t(x, \tilde{x}) = \tau_t(x)^{\top} G_t^{\dagger} \tau_t(\tilde{x}) \quad \forall x, \tilde{x} \in \Omega$$

Empirical Gram matrix:

$$\left[\widehat{\mathbf{G}}_t \right]_{\alpha, \beta} := \int_{\Omega} \tilde{\mathbf{x}}^{\alpha + \beta} d\hat{\mu}(\mathbf{x}) = \frac{1}{N} \sum_{i=1}^N \tilde{\mathbf{x}}_i^{\alpha + \beta} \in \mathbb{R}^{s(t) \times s(t)}$$

Discrete orthonormal basis:

$$\psi_t(\mathbf{x}) := \left(\widehat{\mathbf{G}}_t^{\dagger} \right)^{1/2} \tau_t(\mathbf{x})$$

Kernel matrix:

$$\begin{aligned} \mathbf{K}_t &= \left[k_t(\tilde{\mathbf{x}}_i, \tilde{\mathbf{x}}_j) \right]_{i,j=1}^N = \mathbf{V}_t(\mathcal{X}, \tau) \mathbf{G}_t^{\dagger} \mathbf{V}_t(\mathcal{X}, \tau)^{\top} \\ &= \mathbf{Q}_t(\mathcal{X}, \tau) \mathbf{Q}_t(\mathcal{X}, \tau)^{\top} \end{aligned}$$

$$\mathbf{Q}_t(\mathcal{X}, \tau) = \left[\psi_t(\tilde{\mathbf{x}}_1), \dots, \psi_t(\tilde{\mathbf{x}}_N) \right]^{\top} \in \mathbb{R}^{N \times s(t)}$$

Low-rank approximation

- We would like to extract a suitable subsample of size $r \ll N$ from the original sample such that they correspond to the optimal nodes.
- Computing the spectral decomposition of \mathbf{K}_t can be severely prohibitive as the computational cost is $\mathcal{O}(N^3)$.
- Hence, we will use the diagonally pivoted Cholesky decomposition to select r columns which span the dominant subspace generated by the corresponding kernel functions.
- The said algorithm resorts to a greedy strategy that reduces the trace of the kernel matrix in an iterative manner.
- The computational cost of the pivoted Cholesky is $\mathcal{O}(r^2 N)$.

PIVOTED CHOLESKY DECOMPOSITION

Algorithm 2 (pivoted Cholesky decomposition) :

- **input:** symmetric and positive semi-definite $\mathbf{M} \in \mathbb{R}^{s \times s}$, tolerance $\varepsilon \geq 0$
- **output:** low-rank approximation $\mathbf{M} \approx \mathbf{L}\mathbf{L}^\top$

- 1: **initialization:** set $m := 1$, $\mathbf{d} := \text{diag}(\mathbf{M})$, $\mathbf{L} := []$, $\text{err} := \|\mathbf{d}\|_1$
- 2: **while** $\text{err} > \varepsilon$
- 3: determine $j := \arg \max_{1 \leq i \leq s} d_i$
- 4: compute $\hat{\ell}_m := \mathbf{M}(:, j) - \mathbf{L} * \mathbf{L}^\top(:, j)$
- 5: set $\ell_m := \hat{\ell}_m / \sqrt{d_j}$
- 6: set $\mathbf{L} := [\mathbf{L}, \ell_m]$
- 7: set $\mathbf{d} := \mathbf{d} - \ell_m \odot \ell_m$
- 8: set $\text{err} := \|\mathbf{d}\|_1$
- 9: set $m := m + 1$

FAST EMPIRICAL SCENARIOS

Algorithm 5 (fast empirical scenarios) :

- **Input:** The N samples $\mathcal{X} = \{\tilde{x}_1, \dots, \tilde{x}_N\}$
- **Output:** The r nodes $\Xi = [\xi_1, \dots, \xi_r]$ and probability weights $\mathbf{p} = [p_1, \dots, p_r]$

- 1: Compute the empirical moment matrix $\widehat{\mathbf{M}}_t$ and empirical kernel matrix \mathbf{K}_t
- 2: Perform the pivoted Cholesky decomposition on \mathbf{K}_t to obtain the r nodes $\Xi = [\xi_1, \dots, \xi_r]$ that generate the dominant subspace
- 3: Perform the following optimization problem:

$$\underset{\mathbf{p} \in \mathbb{R}_+^r}{\operatorname{argmin}} \left\| \widehat{\mathbf{M}}_t - \mathbf{V}_t(\Xi, \boldsymbol{\tau})^\top \operatorname{diag}(\mathbf{p}) \mathbf{V}_t(\Xi, \boldsymbol{\tau}) \right\|_F \quad (5)$$

$$\text{subject to: } \mathbf{1}^\top \mathbf{p} = 1$$

to obtain the probabilities $\mathbf{p} = [p_1, \dots, p_r]$

COVARIANCE SCENARIOS

$$M_1(\mathbf{y}) = \mathbf{L} \mathbf{L}^\top.$$

Let \mathbf{H}_v be the Householder reflector, then we have the Vandermonde form:

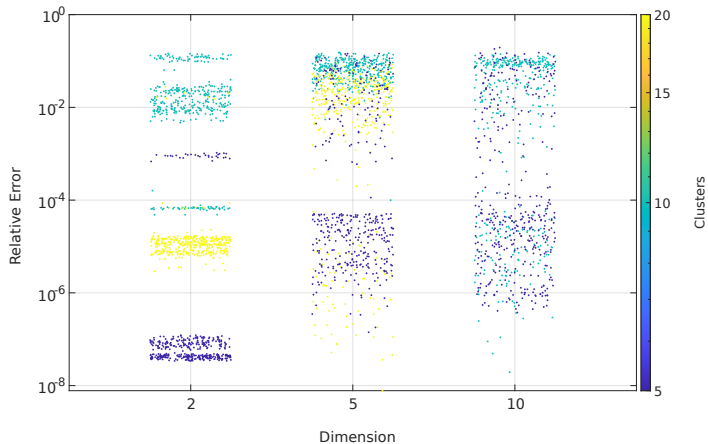
$$\mathbf{M}_1(\mathbf{y}) = \mathbf{L} \mathbf{H}_v^\top \mathbf{H}_v \mathbf{L}^\top = \mathbf{V} \mathbf{D} \mathbf{V}^\top$$

When $t = 1$:

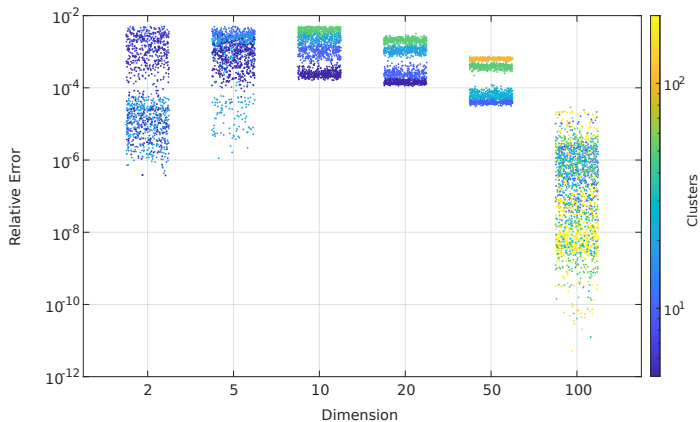
The optimization problem (1) is convex as the Vandermonde matrix is linear in the probability weights and directly solves the interpolation problem

FIGURES

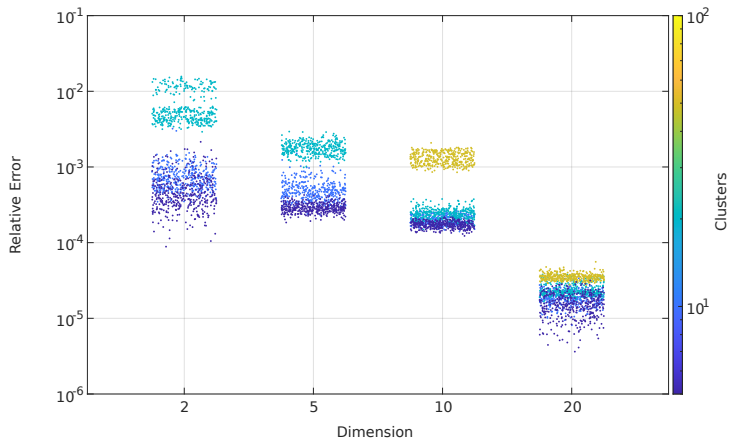
LASSERRE'S GAUSSIAN QUADRATURE



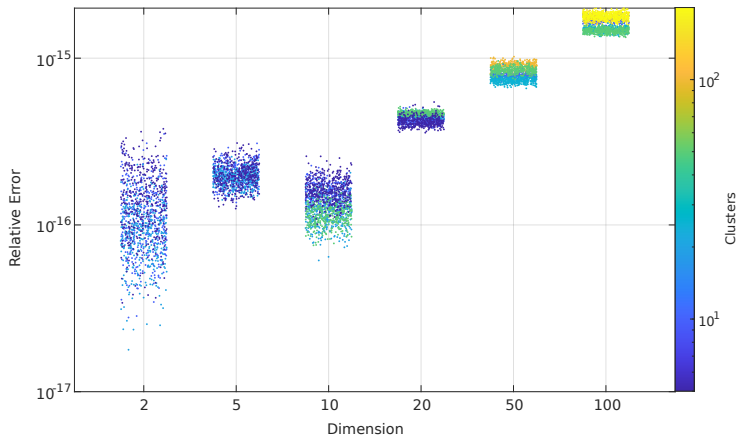
RKHS SCENARIOS (ORDER 1)



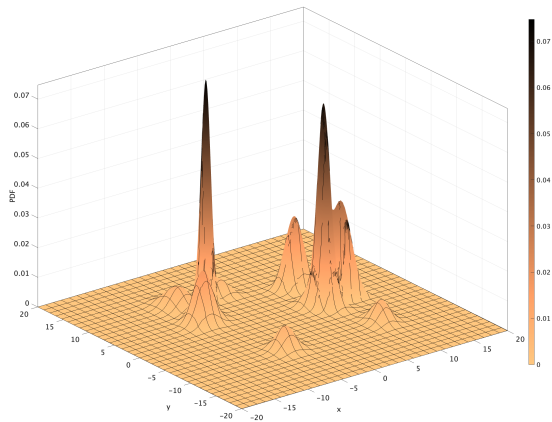
RKHS SCENARIOS (ORDER 2)



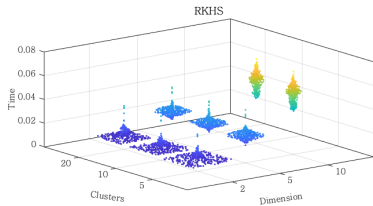
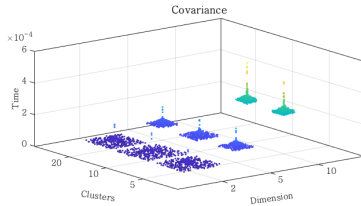
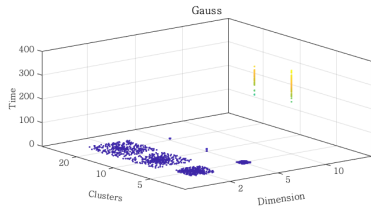
COVARIANCE SCENARIOS



Dimension = 2, Clusters = 10







COMPUTATION TIMES



I would like to thank the
Swiss National Science Foundation
for the grant of this project.

THANK YOU!

REFERENCES

-  A. BERLINET AND C. THOMAS-AGNAN, *Reproducing Kernel Hilbert Spaces In Probability and Statistics*.
-  H. HERBRECHT, M. PETERS, AND R. SCHNEIDER, *On the low-rank approximation by the pivoted cholesky decomposition*, (2011).
-  J. B. LASSERRE, *Moments, Positive Polynomials and Their Applications*.
-  U. MARTEAU-FEREY, F. BACH, AND A. RUDI, *Non-parametric models for non-negative functions*, (2020).