



# Infinite-Dimensional Sums-of-Squares for Optimal Control

Eloïse Berthier  
*with* Ulysse Marteau-Ferey, Justin Carpentier,  
Alessandro Rudi, Francis Bach

June 3, 2022

# Introduction

---

We present a representation of non-negative smooth functions in reproducing kernel Hilbert spaces (RKHS), extending the sum-of-squares (SoS) representation of polynomials.

We apply such representations to optimal control problems, leading to a sample-based numerical method.

# Introduction

---

We present a representation of non-negative smooth functions in reproducing kernel Hilbert spaces (RKHS), extending the sum-of-squares (SoS) representation of polynomials. We apply such representations to optimal control problems, leading to a sample-based numerical method.

The preprint is available on arxiv:

<https://arxiv.org/abs/2110.07396>

# Contents

---

- 1 Sums-of-squares Representations for Non-convex Optimization
- 2 Application to Optimal Control Problems

## Non-convex Optimization as a Linear Program

---

We are interested in finding the global minimum of a possibly non-convex function:

$$f = \min_{x \in \mathbb{R}^p} f(x):$$

This is equivalent to:

$$\begin{aligned} f &= \sup_{c \in \mathbb{R}} c \\ \text{s.t. } &\exists x \in \mathbb{R}^p; f(x) \leq c \end{aligned}$$

## Non-convex Optimization as a Linear Program

---

We are interested in finding the global minimum of a possibly non-convex function:

$$f = \min_{x \in \mathbb{R}^p} f(x):$$

This is equivalent to:

$$\begin{aligned} f &= \sup_{c \in \mathbb{R}} c \\ \text{s.t. } \exists x \in \mathbb{R}^p; & f(x) \leq c \end{aligned}$$

## Non-convex Optimization as a Linear Program

---

We are interested in finding the global minimum of a possibly non-convex function:

$$f = \min_{x \in \mathbb{R}^p} f(x):$$

This is equivalent to:

$$\begin{aligned} f &= \sup_{c \in \mathbb{R}} c \\ \text{s.t. } \exists x \in \mathbb{R}^p; & \boxed{f(x) - c \leq 0.} \end{aligned}$$

*How to handle a dense set of constraints?*

## Idea 1: Subsampling Inequalities

---

$$f = \sup_{c \in \mathbb{R}} c$$
$$\text{s.t. } \exists x \in \mathbb{R}^p; f(x) - c \leq 0.$$

Relax it to:

$$f_n = \sup_{c \in \mathbb{R}} c$$
$$\text{s.t. } \exists i \in \{1, \dots, n\}; f(x_i) - c \leq 0;$$

which is equivalent to...



## Idea 1: Subsampling Inequalities

---

$$f = \sup_{c \in \mathbb{R}} c$$
$$\text{s.t. } \exists x \in \mathbb{R}^p; f(x) - c \leq 0.$$

Relax it to:

$$f_n = \sup_{c \in \mathbb{R}} c$$
$$\text{s.t. } \exists i \in \{1, \dots, n\}; f(x_i) - c \leq 0;$$

which is equivalent to...  $f^* - f_n = \min_j f(x_j)$ .

## Idea 1: Subsampling Inequalities

---

$$f = \sup_{c \in \mathbb{R}} c$$
$$\text{s.t. } \exists x \in \mathbb{R}^p; f(x) - c \leq 0.$$

Relax it to:

$$f_n = \sup_{c \in \mathbb{R}} c$$
$$\text{s.t. } \exists i \in \{1, \dots, n\}; f(x_i) - c \leq 0;$$

which is equivalent to...  $f \approx f_n = \min_i f(x_i)$ .

If  $f$  Lipschitz, we need  $O(n^p)$  samples to approximate  $f$  up to  $\epsilon$ .

If  $f \in \mathcal{C}^s(\mathbb{R}^p)$  is smooth, the lower-bound is  $O(n^{p-s})$  [3].

*Can we do any better?*

## Idea 2: Representing Non-negative Functions

---

$$f = \sup_{c \in \mathbb{R}} c$$

s.t.  $\exists x \in \mathbb{R}^p; f(x) \leq c \leq 0.$

We need a *practical representation* of non-negative functions.

## Idea 2: Representing Non-negative Functions

---

$$f = \sup_{c \in \mathbb{R}} c$$

s.t.  $\exists x \in \mathbb{R}^p; f(x) - c \leq 0.$

We need a *practical representation* of non-negative functions.

Imagine we know how to represent some  $g_k$ , e.g., of the form:

$$g_k(x) = h_k(x)^2$$

Then we can generate non-negative functions as sum-of-squares:

$$g(x) = \sum_{k=1}^m g_k(x)^2$$

## Idea 2: Representing Non-negative Functions

---

$$f = \sup_{c \in \mathbb{R}} c$$
$$\text{s.t. } \exists x \in \mathbb{R}^p; \quad f(x) - c \leq 0.$$

We need a *practical representation* of non-negative functions.

Imagine we know how to represent some  $g_k$ , e.g., of the form:

$$g_k(x) = h_k'(x) A_k^{-1} h_k(x)$$

Then we can generate non-negative functions as sum-of-squares:

$$g(x) = \sum_{k=1}^m g_k(x)^2 = \sum_{k=1}^m h_k'(x) A_k^{-1} h_k(x)$$

where  $A = \sum_{k=1}^m A_k^{-1} h_k(x) h_k(x)'$  has rank less than  $m$ .

## Polynomial Sum-of-Squares (SoS)

---

In dimension 1, all non-negative polynomials are SoS.

This is **not true** in larger dimensions.

Powerful theorems describe cases of tight SoS representations [1].

## Polynomial Sum-of-Squares (SoS)

---

In dimension 1, all non-negative polynomials are SoS.

This is **not true** in larger dimensions.

Powerful theorems describe cases of tight SoS representations [1].

### Theorem (Putinar's Positivstellensatz (simplified))

Let  $(h_k)_k$  a family of polynomials and  $W = \{x \in \mathbb{R}^d \mid h_k(x) \geq 0\}$  a **semi-algebraic set**. Assume that  $\{x \in \mathbb{R}^d \mid h_k(x) \geq 0\}$  is **compact** for some  $k$ . If a polynomial  $f$  is strictly positive on  $W$ , then there exists SoS polynomials  $(s_k)_{k=0}^m$  such that:

$$f = s_0 + \sum_{k=1}^m s_k h_k$$

## The Moment – SoS Hierarchy

---

Let  $f$  a polynomial of degree  $d_0$ . We want to solve:

$$f = \min_{x \in \mathbb{R}^p} f(x) \text{ s.t. } g_k \geq f_1; \dots; m; h_k(x) = 0$$

Lasserre's Hierarchy of semi-definite programs (SDP):

Find  $c; X_k = 0; k = 0; \dots; m$  such that

$$g \geq N_{d_0}^p; f - c \mathbf{1} = \sum_{k=0}^m h C^k; X_k$$

$$g \geq N_{2r}^p; n N_{d_0}^p; 0 = \sum_{k=0}^m h C^k; X_k$$



## The Moment – SoS Hierarchy

Let  $f$  a polynomial of degree  $d_0$ . We want to solve:

$$f = \min_{x \in \mathbb{R}^p} f(x) \text{ s.t. } g_k \geq f_1; \dots; m; h_k(x) = 0$$

Lasserre's Hierarchy of semi-definite programs (SDP):

Find  $c; X_k = 0; k = 0; \dots; m$  such that

$$g \geq N_{d_0}^p; f - c \mathbf{1} = \sum_{k=0}^m h C^k; X_k$$

$$g \geq N_{2r}^p \cap N_{d_0}^p; 0 = \sum_{k=0}^m h C^k; X_k$$

The monomials are indexed by  $N_r^p := \{f \in N^p : |j| \leq r\}$  which has size  $s_r(d) = \binom{p+r}{r}$ . This is exponential in the dimension  $p$ .

## SoS Functions in Reproducing Kernel Hilbert Spaces

---

A positive-definite kernel on  $\mathbb{R}^p$  is a function  $K : \mathbb{R}^p \times \mathbb{R}^p \rightarrow \mathbb{R}$  such that  $\forall n \geq 1; \forall (x_1, \dots, x_n)$ , the matrix  $(K(x_i, x_j))$  is PSD.

It is associated to a Hilbert space  $H$  such that:

$$\begin{aligned} \forall x \in \mathbb{R}^p, \quad \phi(x) &:= K(x, \cdot) \in H; \\ \forall f \in H; x \in \mathbb{R}^p, \quad \langle f, \phi(x) \rangle &= f(x). \end{aligned}$$

## SoS Functions in Reproducing Kernel Hilbert Spaces

A positive-definite kernel on  $\mathbb{R}^p$  is a function  $K : \mathbb{R}^p \times \mathbb{R}^p \rightarrow \mathbb{R}$  such that  $\forall n \geq 1; \forall (x_1, \dots, x_n)$ , the matrix  $(K(x_i, x_j))$  is PSD.

It is associated to a Hilbert space  $H$  such that:

$$\forall x \in \mathbb{R}^p, \phi(x) := K(x, \cdot) \in H;$$
$$\forall f \in H; \forall x \in \mathbb{R}^p, \langle f, \phi(x) \rangle = f(x).$$

Using the reproducing property, a sum-of-squares of functions in  $H$ :

$$\forall x \in \mathbb{R}^p; \quad g(x) = \sum_{k=1}^m h_k(x)^2$$

is such that

$$\forall x \in \mathbb{R}^d; \quad g(x) = \langle h'(x); A^{-1}(x) \rangle;$$

where  $A \in S_+(H)$  is a PSD operator, *possibly infinite-dimensional*.

## SoS Representation of Smooth Functions [4]

We consider as our RKHS  $H$  the Sobolev space  $W_2^s(\mathbb{R}^p)$ , with  $s = p/2 + 3$ , of  $s$ -smooth functions.

### Theorem (informal)

If  $f \in H$  has a unique **isolated** global minimum at  $x^*$  s.t.  $\frac{\partial^2 f}{\partial x^2}(x^*) \succ 0$ , then there exists  $h_1, \dots, h_m \in H$ , with  $m = p + 1$  such that:

$$\exists \lambda; f(x) - f(x^*) = \sum_{k=1}^m \lambda_k h_k(x)^2;$$

Hence  $f - f(x^*)$  is a SoS of (smooth) functions in  $H$ , and:

$$\exists A \in S_+(H) \text{ s.t. } \exists \lambda; f(x) - f(x^*) = \lambda'(x) A \lambda(x); A'(x) \succ 0;$$

*No need for a hierarchy!*

# Non-convex Optimization of Smooth Functions

---

$$f = \sup_{c \in \mathbb{R}} c$$
$$\text{s.t. } \exists x \in \mathbb{R}^p; f(x) - c \leq 0.$$

Using the Theorem, if  $f$  is smooth, this is equivalent to:

$$f = \sup_{c \in \mathbb{R}; A \in \mathcal{S}_+(H)} c$$
$$\text{s.t. } \exists x; f(x) - c = h'(x); A'(x) \leq 0.$$

# Non-convex Optimization of Smooth Functions

---

We can now *subsample equalities*:

$$f_n = \sup_{c \in \mathbb{R}; A \in \mathcal{S}_+(H)} c \quad \text{Tr}(A)$$

s.t.  $\exists i \in \{1, \dots, n\}; f(x_i) \quad c = h'(x_i); A'(x_i)i:$

## Non-convex Optimization of Smooth Functions

---

We can now *subsample equalities*:

$$f_n = \sup_{c \in \mathbb{R}; A \in \mathcal{S}_+(H)} c \quad \text{Tr}(A)$$

$$\text{s.t. } \forall i \in \{1, \dots, n\}; f(x_i) \leq c = h'(x_i); A' (x_i) i:$$

Using the reproducing property, this is equivalent to the SDP:

$$f_n = \sup_{c \in \mathbb{R}; B \succeq 0} c \quad \text{Tr}(B)$$

$$\text{s.t. } \forall i \in \{1, \dots, n\}; f(x_i) \leq c = \Phi_i^T B \Phi_i;$$

where the  $\Phi_i \in \mathbb{R}^n$  are vectors computed from the kernel matrix.

This achieves an *almost optimal rate* of  $O(n^{(s-3)+p+1=2})$  for  $s-3+p=2$ . The lower-bound is  $O(n^{s-p})$ .

In short... (see [3])

---



# Contents

---

- 1 Sums-of-squares Representations for Non-convex Optimization
- 2 Application to Optimal Control Problems

## Optimal Control as a Linear Program

---

The optimal control problem is to find  $V$  such that:

$$V(t_0; x_0) = \inf_{u(\cdot)} \int_{t_0}^T L(t; x(t); u(t)) dt + M(x(T))$$
$$\forall t \in [t_0; T]; \dot{x}(t) = f(t; x(t); u(t)); \quad x(0) = x_0:$$

## Optimal Control as a Linear Program

---

The optimal control problem is to find  $V$  such that:

$$V(t_0; x_0) = \inf_{u(\cdot)} \int_{t_0}^T L(t; x(t); u(t)) dt + M(x(T))$$

$$\forall t \in [t_0; T]; \dot{x}(t) = f(t; x(t); u(t)); \quad x(0) = x_0:$$

Under convexity assumptions, this is equivalent to finding a *maximal subsolution* of the Hamilton-Jacobi-Bellman equation [2]:

$$\sup_{V \in C^1([0; T] \times X)} V(0; x_0) = V(0; x_0)$$

$$\partial_t V(t; x; u) + \frac{\partial V}{\partial x}(t; x) \cdot f(t; x; u) + L(t; x; u) - r V(t; x) \geq 0$$

$$\partial_x V(T; x) = M(x):$$

## Optimal Control as a Linear Program

---

The optimal control problem is to find  $V$  such that:

$$V(t_0; x_0) = \inf_{u(\cdot)} \int_{t_0}^T L(t; x(t); u(t)) dt + M(x(T))$$

$$\forall t \in [t_0; T]; \dot{x}(t) = f(t; x(t); u(t)); \quad x(0) = x_0:$$

Under convexity assumptions, this is equivalent to finding a *maximal subsolution* of the Hamilton-Jacobi-Bellman equation [2]:

$$\sup_{V \in C^1([0; T] \times X)} V(0; x_0) \geq V(0; x_0)$$

$$\forall (t; x; u); \quad \frac{\partial V}{\partial t}(t; x) + L(t; x; u) + \inf_{u} \{ \langle \nabla_x V(t; x), f(t; x; u) \rangle - f(t; x; u) \} \leq 0$$

$$\forall x; \quad V(T; x) = M(x):$$

## A Simple Baseline: Subsampling Inequalities

---

Using a linear parameterization of  $V$ , and simply subsampling inequalities leads to an LP:

$$\begin{aligned} \sup_{\theta \in \mathbb{R}^m} \quad & \frac{1}{n} \sum_{i=1}^n V(\theta; x^{(i)}) \\ \text{s.t.} \quad & H(t^{(i)}; x^{(i)}; u^{(i)}) \leq 0 \end{aligned}$$

This is already a non-trivial numerical method.

*Can we do any better?*

## SoS Representation of the Hamiltonian

### Theorem (informal)

Assume that:

$f$  is **control-affine**:  $f(t;x;u) = g(t;x) + B(t;x)u$ ;

$L$  is strongly convex in  $u$ ;

$L$ ,  $B$  and  $V$  are **sufficiently smooth**;

Then  $H$  is a SoS of  $p$  smooth functions  $(w_j)_{j=1}^p \in C^s(\Omega)$ :

$$\forall (t;x;u) \in \Omega; \quad H(t;x;u) = \sum_{j=1}^p w_j(t;x;u)^2$$

**Limit:** in general  $V$  is not even  $C^1$ . A possible workaround is to add noise to the dynamics to smoothen  $V$ .

## A Practical Algorithm for Smooth Optimal Control

---

Adding regularization terms, we get the following SDP:

$$\sup_{B \succ 0; \gamma} c^T \gamma \quad k \gamma^2 \quad \text{Tr}(B) \quad k \gamma^2 + \alpha \log \det B + C$$

$$\text{such that } \gamma_i \geq f_1; \dots; f_n; \quad b_i + a_i \gamma = (\Phi_i)^T B \Phi_i + \gamma_i$$

The dual is solved with damped Newton's method, an algorithm with cost  $O(n^3)$  in time and  $O(n^2)$  in space for each iteration.

## Numerical Example

---

We solve a linear quadratic regulator. We know that:

$$H(t; x; u) = (u + F(t)x)^T R(u + F(t)x):$$

We use:  $K((t; x; u); (t^0; x^0; u^0)) = hu; u^0i + hx; x^0ie^{-jt - t^0j}$ :



## Numerical Example

---

We solve a linear quadratic regulator. We know that:

$$H(t; x; u) = (u + F(t)x)^T R(u + F(t)x):$$

We use:  $K((t; x; u); (t^0; x^0; u^0)) = hu; u^0 i + hx; x^0 i e^{-jt - t^0 j}$ :

# Again...

---

## Conclusion

---

We have presented an **extension** of the SoS framework in RKHS.

It leads to **sample-based** numerical methods involving SDPs.

There are **many potential applications**, e.g., to optimal transport, sampling, modelling of probability distributions...

# References

---



J.-B. Lasserre.

*Moments, Positive Polynomials and their Applications*, volume 1.  
World Scientific, 2010.



J.-B. Lasserre, D. Henrion, C. Prieur, and E. Trélat.

Nonlinear optimal control via occupation measures and LMI-relaxations.  
*SIAM J. on Contr. and Optim.*, 47(4):1643–1666, 2008.



E. Novak.

*Deterministic and Stochastic Error Bounds in Numerical Analysis*.  
Springer, 2006.



A. Rudi, U. Marteau-Ferey, and F. Bach.

Finding global minima via kernel approximations.  
Technical Report 2012.11978, arXiv, 2020.