

# A Faster Interior-Point Method for Sum-Of-Squares Optimization

Shunhua Jiang ✉

Columbia University, New York, NY, USA

Bento Natura ✉

London School of Economics, UK

Omri Weinstein ✉

The Hebrew University, Jerusalem, Israel

Columbia University, New York, NY, USA

---

## Abstract

We present a faster interior-point method for optimizing sum-of-squares (SOS) polynomials, which are a central tool in polynomial optimization and capture convex programming in the Lasserre hierarchy. Let  $p = \sum_i q_i^2$  be an  $n$ -variate SOS polynomial of degree  $2d$ . Denoting by  $L := \binom{n+d}{d}$  and  $U := \binom{n+2d}{2d}$  the dimensions of the vector spaces in which  $q_i$ 's and  $p$  live respectively, our algorithm runs in time  $\tilde{O}(LU^{1.87})$ . This is polynomially faster than state-of-art SOS and semidefinite programming solvers [16, 15, 27], which achieve runtime  $\tilde{O}(L^{0.5} \min\{U^{2.37}, L^{4.24}\})$ .

The centerpiece of our algorithm is a dynamic data structure for maintaining the inverse of the Hessian of the SOS barrier function under the *polynomial interpolant basis* [27], which efficiently extends to multivariate SOS optimization, and requires maintaining spectral approximations to low-rank perturbations of *elementwise (Hadamard) products*. This is the main challenge and departure from recent IPM breakthroughs using inverse-maintenance, where low-rank updates to the slack matrix readily imply the same for the Hessian matrix.

**2012 ACM Subject Classification** Mathematics of computing → Continuous functions; Mathematics of computing → Convex optimization; Mathematics of computing → Semidefinite programming; Mathematics of computing → Stochastic control and optimization

**Keywords and phrases** Interior Point Methods, Sum-of-squares Optimization, Dynamic Matrix Inverse

**Digital Object Identifier** 10.4230/LIPIcs.ICALP.2022.79

**Category** Track A: Algorithms, Complexity and Games

**Related Version** *Full Version*: <https://arxiv.org/abs/2202.08489>

**Funding** *Shunhua Jiang*: Supported by NSF CAREER award CCF-1844887.

*Bento Natura*: This project has received funding from the European Research Council (ERC) under the European Union's Horizon 2020 research and innovation programme (grant agreement no. 757481–ScaleOpt).

*Omri Weinstein*: Supported by NSF CAREER award CCF-1844887 and ISF grant #3011005535.

**Acknowledgements** The second author would like to thank Vissarion Fisikopoulos and Elias Tsigaridas for introducing him from a practical perspective to Sum-of-Squares Optimization under the interpolant basis.

## 1 Introduction

Polynomial optimization is a fundamental problem in many areas of applied mathematics, operations research, and theoretical computer science, including combinatorial optimization [5, 36, 4], statistical estimation [13, 14], experimental design [26], control theory [12], signal



© Shunhua Jiang, Bento Natura, and Omri Weinstein;  
licensed under Creative Commons License CC-BY 4.0

49th International Colloquium on Automata, Languages, and Programming (ICALP 2022).

Editors: Mikołaj Bojańczyk, Emanuela Merelli, and David P. Woodruff;

Article No. 79; pp. 79:1–79:20



Leibniz International Proceedings in Informatics

Schloss Dagstuhl – Leibniz-Zentrum für Informatik, Dagstuhl Publishing, Germany



processing [31], power systems engineering [11], discrete geometry [2, 3] and computational algebraic geometry [20]. In the most basic formulation, we are given a collection of  $k$  real  $n$ -variate polynomials  $f_1, \dots, f_k$  and an objective function  $f : \mathbb{R}^n \rightarrow \mathbb{R}$ , and the goal is to minimize  $f$  over the set  $\mathcal{S} := \{t \in \mathbb{R}^n \mid \forall i \in \{1, \dots, k\} : f_i(t) \geq 0\}$ , that is, to find

$$\inf_{t \in \mathbb{R}^n} \{f(t) \mid t \in \mathcal{S}\}, \quad (1)$$

which is equivalent to checking polynomial nonnegativity  $\sup_{c \in \mathbb{R}} \{c \mid f(t) - c \geq 0, \forall t \in \mathcal{S}\}$ . This is then equivalent to computing  $\sup_{c \in \mathbb{R}} \{c \mid f - c \in \mathcal{K}(\mathcal{S})\}$ , where  $\mathcal{K}(\mathcal{S})$  denotes the convex cone of all polynomials of degree at most  $\deg(f)$  that are non-negative on the set  $\mathcal{S}$ . This is an instance of the more general conic programming:

$$\min_{x \in \mathbb{R}^N} c^\top x \text{ s.t. } Ax = b, x \in \mathcal{K}, \quad (2)$$

where  $\mathcal{K} \subset \mathbb{R}^N$  is some convex cone<sup>1</sup>. The conic optimization problem over the cone  $\mathcal{K}(\mathcal{S})$  is intractable in general because there is no simple characterization of  $\mathcal{K}(\mathcal{S})$ . Nevertheless, there always exists an increasing family of convex cones of *weighted sum-of-squares* polynomials that converges to any such cone  $\mathcal{K}(\mathcal{S})$ .

We first introduce the notion of *sum-of-squares* (SOS) polynomials: Denoting by  $\mathcal{V}_{n,d}$  the vector space of all  $n$ -variate polynomials of (total) degree at most  $d$ , a polynomial  $p \in \mathcal{V}_{n,2d}$  is said to be *sum-of-squares* (SOS) if it can be written as a finite sum of square polynomials, i.e., there exist  $q_1, \dots, q_\ell$  such that  $p = \sum_{i=1}^{\ell} q_i^2$ . The set  $\Sigma_{n,2d}$  of SOS polynomials of degree at most  $2d$  is a (proper) cone contained in  $\mathcal{V}_{n,2d}$ , of dimension  $U := \dim(\mathcal{V}_{n,2d}) = \binom{n+2d}{2d}$ , as the vector space  $\mathcal{V}_{n,2d}$  is isomorphic to  $\mathbb{R}^U$ . If  $p$  can be written as  $p = \sum_{i=1}^k f_i s_i$  for  $s_1 \in \Sigma_{n,2d_1}, \dots, \Sigma_{n,2d_k}$  and  $k$  nonzero polynomials  $\mathbf{f} := (f_1, \dots, f_k)$ , then it is said to be *weighted sum-of-squares* (WSOS).

*Putinar's Positivstellensatz* [29] states that under mild conditions, any polynomial  $p$  that is non-negative on  $\mathcal{S}$  can be written as a WSOS polynomial  $\sum_{i=1}^k f_i s_i$ , albeit with (potentially) *unbounded degree*  $s_i$ 's. In WSOS optimization we consider sum-of-squares polynomials  $s_i$  with bounded degree, so the hierarchy of WSOS optimization with increasing degree (known as the Lasserre hierarchy) can be viewed as a tool for approximating general polynomial optimization. For more details of this approximation scheme for polynomial optimization, we refer the readers to the textbooks [19, 7].

This paper concerns algorithms for *(W)SOS optimization*, which is the conic optimization program (2) where the underlying cone  $\mathcal{K}$  is the (W)SOS cone:

$$\min_{x \in \mathbb{R}^U} c^\top x \text{ s.t. } Ax = b, x \in \Sigma_{n,2d}, \quad (3)$$

where  $x \in \mathcal{V}_{n,2d}$  is the vector of coefficients which encodes the polynomial. Henceforth, we focus on the case where  $\mathcal{K} = \Sigma_{n,2d}$  is the SOS cone. See our full version for how to extend our algorithm for SOS optimization to WSOS.

The computational complexity of solving Problem 3 naturally depends on the dimensions

$$L := \dim(\mathcal{V}_{n,d}) = \binom{n+d}{d}, \quad U := \dim(\mathcal{V}_{n,2d}) = \binom{n+2d}{2d} \quad (4)$$

of the underlying vector spaces (Note that  $L \leq U \leq L^2$ ). We now turn to explain the previous approaches for SOS optimization solvers.

<sup>1</sup> A subset  $\mathcal{K} \subset \mathbb{R}^N$  is a convex cone if  $\forall x, y \in \mathcal{K}$  and  $\alpha, \beta \in \mathbb{R}_+$ ,  $\alpha x + \beta y \in \mathcal{K}$ .

## SOS Optimization as SDPs

A fundamental fact is that the *dual* SOS cone is a *slice of the SDP cone* [24]. More formally, for any *fixed bases*  $\mathbf{p} = (p_1, p_2, \dots, p_L)$  and  $\mathbf{q} = (q_1, q_2, \dots, q_U)$  to  $\mathcal{V}_{n,d}$  and  $\mathcal{V}_{n,2d}$  respectively, there exists a unique linear mapping  $\Lambda : \mathbb{R}^U \rightarrow \mathbb{R}^{L \times L}$  satisfying

$$\Lambda(\mathbf{q}(t)) = \mathbf{p}(t)\mathbf{p}(t)^\top, \quad \forall t \in \mathbb{R}^n. \quad (5)$$

Here we define  $\mathbf{p}(t) = (p_1(t), p_2(t), \dots, p_L(t))^\top$  and  $\mathbf{q}(t) = (q_1(t), q_2(t), \dots, q_U(t))^\top$ . An equivalent way to view the definition of  $\Lambda$  in (5) is as follows: For polynomials  $p_i, p_j \in \mathbf{p}$  there are unique coefficients  $\lambda_{iju}$  such that  $p_i p_j = \sum_{u \in U} \lambda_{iju} q_u$ . These  $\lambda_{iju}$  define the mapping  $\Lambda$  unambiguously.

This in turn implies that a polynomial  $s \in \mathcal{V}_{n,2d}$  (we view  $s$  as a vector in  $\mathbb{R}^U$  that corresponds to its coefficients over the basis  $\mathbf{q}$ ) is in the dual SOS cone  $\Sigma_{n,2d}^*$  if and only if  $\Lambda(s)$  is a *positive semidefinite* (PSD) matrix (proved by [24], see Theorem 10 for details). As [27] recently observed, the choice of the bases  $\mathbf{p}, \mathbf{q}$  crucially affects the complexity of the optimization problem, more on this below.

Equation (5) implies the well-known fact that optimization over SOS polynomials (3) can be reduced to semidefinite programming

$$\min_{X \succeq 0} \{ \langle C, X \rangle \mid \text{tr}(A_i X) = b_i, \forall i \in [m] \}, \quad (\text{SDP})$$

and can thus be solved using off-the-shelf SDP solvers. However, despite recent breakthroughs on the runtime of general SDP solvers via *interior-point methods* (IPMs) [16, 15], this SDP reformulation of (3) does not scale well for moderately large degrees, i.e., whenever  $U \ll L^2$  in (4). This is because the SDP reformulation always incurs a factor of at least  $L^2$ , even when  $U \ll L^2$ , as this is the SDP variable size (the PSD matrix  $X$  has size  $L \times L$ ). Indeed, for the current fast-matrix-multiplication (FMM) exponent  $\omega \approx 2.37$  [21, 1], the running time of state-of-the-art SDP solvers [16, 15] for SOS optimization (Problem 3) is<sup>2</sup>

$$\tilde{O}(L^{0.5} \cdot \min\{UL^2 + U^{2.37}, L^{4.24}\}). \quad (6)$$

An alternative approach is to solve Problem 3 directly by designing an *ad-hoc* IPM for the dual SOS cone, avoiding the blowup in the SDP reformulation. This was exactly the motivation of [27]. Like all aforementioned SDP solvers, [27]’s SOS solver is based on IPMs [25], which iteratively minimize the original objective function plus a barrier function via Newton steps. When applied to the SOS Problem (3), *the choice of the specific bases  $\mathbf{p}, \mathbf{q}$  crucially affects the structure of the (Hessian of the) barrier function  $F(s) = F(\Lambda(s))$* , and hence the cost-per-iteration of the IPM. As such, choosing a “good” and efficient basis is key to a fast algorithm for (3). One of the main contribution of [27] is an efficient basis for the SOS cone, which efficiently scales to multivariate SOS, yielding an IPM whose total runtime is

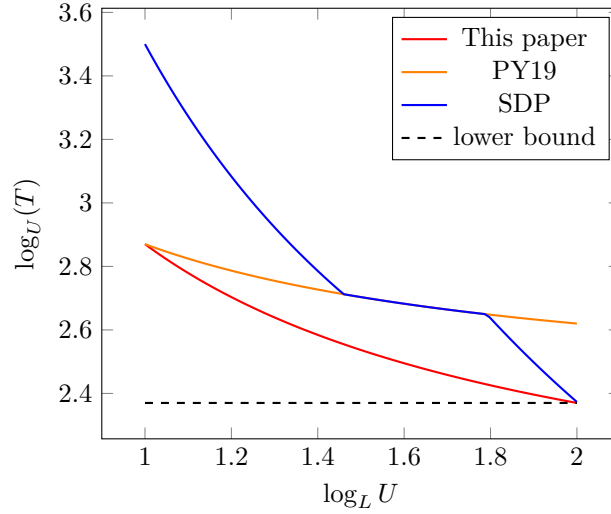
$$\tilde{O}(L^{0.5}U^\omega) \approx \tilde{O}(L^{0.5}U^{2.37}). \quad (7)$$

Our main result is a polynomially faster IPM for Problem 3:

► **Theorem 1** (Main Result, Informal version of Theorem 16). *With current FMM exponent, there is an algorithm for solving Problem (3), whose total running time is  $\tilde{O}(LU^{1.87})$ .*

Indeed, this runtime is polynomially faster than (7) and (6), as shown in Figure 1. We now turn to elaborate on the technical approach for proving Theorem 1.

<sup>2</sup> We use  $\tilde{O}(\cdot)$  to hide  $U^{o(1)}$  and  $\log(1/\delta)$  factors.



■ **Figure 1** Overview of current running times of recent solvers for SOS. The lower bound stems from solving a linear system in  $U$  variables, i.e.,  $T = \Omega(U^\omega)$  where  $\omega \approx 2.37$ .

### Faster IPMs via Inverse-Maintenance

Interior-Point Methods (IPMs [18, 30]) are a powerful class of second-order optimization algorithms for convex optimization, which essentially reduce a conic optimization problem (2) to solving *a sequence of slowly-changing linear systems* (via Newton steps). Since their discovery in the mid 80’s, IPMs have emerged as the “gold-standard” of convex optimization, as they are known to converge fast in *both theory and practice* [35]. The main computational cost of IPMs is computing, in each iteration, the inverse of the Hessian of the underlying *barrier function*  $F(s) = F(\Lambda(s))$ , which naively costs at least  $U^\omega$  time per iteration for the SOS optimization problem [27]. A recent influential line of work [9, 16], inspired by [37]’s seminal work, has demonstrated that *dynamically maintaining* the inverse of the Hessian matrix under *low-rank* updates using clever *data structures*, can lead to much cheaper *cost-per-iteration*. All of these results rely on a careful combination of dynamic data structures with the geometry (e.g., spectral approximation) of the underlying optimization method and barrier function. Similar techniques have been extended to other optimization problems as well [23, 38, 17, 33, 34]. This paper extends this line of work to SOS optimization.

### Our Techniques

Following the framework of [27], we also choose the *polynomial interpolant basis* representation, which corresponds to a linear operator  $\Lambda : \mathbb{R}^U \rightarrow \mathbb{R}^{L \times L}$  is  $\Lambda(s) = P^\top \text{diag}(s)P$ , where  $P \in \mathbb{R}^{U \times L}$  is the matrix whose entries are the evaluation of the Lagrange interpolation polynomials, through some unisolvent<sup>3</sup> set of points in  $\mathcal{V}_{n,d}$  (see Section 3 for a formal definition). This basis induces the aforementioned convenient form of  $\Lambda$ , and generalizes to the multivariate case. The Hessian of the barrier function  $F(s) = -\log \det(\Lambda(s))$  is given by

$$H(s) = (P(P^\top \text{diag}(s)P)^{-1}P^\top)^{\circ 2} \in \mathbb{R}^{U \times U},$$

<sup>3</sup> Any set of points in  $\mathbb{R}^n$  for which the evaluation of a polynomial in  $\mathcal{V}_{n,d}$  on these points uniquely defines the polynomial.

where  $A \circ B$  denotes the element-wise (Hadamard) product of two matrices. The main bottleneck of each iteration of IPMs is to compute the Hessian inverse  $H(x)^{-1}$  of the Newton step, which naively takes  $O(U^\omega)$  time.

In IPM theory, it has long been known that it suffices to compute a spectral approximation of the Hessian. We follow the “lazy update” framework in recent developments of LP and SDP solvers [9, 16], which batches together low-rank updates to  $M := P(P^\top \text{diag}(s)P)^{-1}P^\top$ , where  $\text{rk}(M) = L$ . In each iteration, we can compute a spectral approximation  $M^{\text{new}} = M + UV^\top$ , where  $U, V$  are low rank matrices with size  $U \times r$  where  $r \ll U$  is chosen to optimize the runtime. Since  $\widetilde{M} \approx M$  implies that  $\widetilde{M}^{\circ 2} \approx M^{\circ 2}$ , this also gives a spectral approximation of the Hessian.

The main challenge here, compared to previous LP and SDP solvers [37, 22, 9, 16, 15], is that low-rank updates to  $M$  do not readily translate to a low-rank update to  $(M^{\circ 2})^{-1}$ , since Hadamard-products can *increase* the rank  $\text{rk}(A \circ B) \leq \text{rk}(A) \cdot \text{rk}(B)$ , in contrast to standard matrix multiplication which does not increase the rank  $\text{rk}(AB) \leq \max\{\text{rk}(A), \text{rk}(B)\}$ . This means that we cannot directly apply Woodbury’s identity to efficiently update the inverse of the Hessian, which is the common approach in all aforementioned works. Instead, we employ the following property which relates rank-one Hadamard-product perturbations to standard matrix products

$$M \circ (u \cdot v^\top) = \text{diag}(u) \cdot M \cdot \text{diag}(v),$$

which means that we can translate the rank- $r$  update of  $M$  into a rank- $Lr$  update of  $M^{\circ 2}$  for  $r \leq L$ . With some further calculations, applying Woodbury’s identity on the resulting matrix, implies that we can compute  $((M^{\text{new}})^{\circ 2})^{-1}$  in time

$$O(\mathcal{T}_{\text{mat}}(U, U, Lr)),$$

which is never worse than  $\mathcal{T}_{\text{mat}}(U, U, U) = U^\omega$  as long as  $r \leq U/L$ . Modifying the amortization tools of [16] and [15], combined with basic spectral theory for Hadamard products, we show that our amortized cost per iteration is bounded by

$$O(U^2 + U^{\omega-1/2} \cdot L^{1/2}),$$

which becomes  $O(U^2 + U^{1.87}L^{0.5})$  if we plug in the current matrix multiplication exponent.

## 2 Preliminaries

In this section we provide the definitions and the tools that we will use. For any integer  $n > 0$ , we define  $[n] = \{1, 2, \dots, n\}$ . We use  $\mathbb{R}_+$  and  $\mathbb{R}_{\geq 0}$  to denote the set of positive and non-negative real numbers respectively. We use  $0_n, 1_n \in \mathbb{R}^n$  to denote the all-zero and all-one vectors of size  $n$ .

Given a vector  $v \in \mathbb{R}^n$ , for any  $m \leq n$ , we use  $v_{[1:m]} \in \mathbb{R}^m$  to denote the first  $m$  entries of  $v$ . For a vector  $v \in \mathbb{R}^n$ , we use  $\text{diag}(v) \in \mathbb{R}^{n \times n}$  to denote the diagonal matrix whose diagonal entries are  $v$ . For a square matrix  $A \in \mathbb{R}^{n \times n}$ , we use  $\text{diag}(A) \in \mathbb{R}^n$  to denote the vector of the diagonal entries of  $A$ . We use  $\text{rk}(A)$  to denote the rank of a matrix  $A$ . We use  $\ker(A)$  and  $\text{Im}(A)$  to denote the kernel space and the column space of  $A$ .

We say a matrix  $A \in \mathbb{R}^{n \times n}$  is PSD (denoted as  $A \succeq 0$ ) if  $A$  is symmetric and  $x^\top Ax \geq 0$  for all  $x \in \mathbb{R}^n$ . We use  $\mathbb{S}^{n \times n}$  to denote the set of PSD matrices of size  $n \times n$ . The spectral norm of a matrix  $A \in \mathbb{R}^{n \times d}$  is defined as  $\|A\|_2 = \max_{x \in \mathbb{R}^d, \|x\|_2=1} \|Ax\|_2$ . The Frobenius norm of  $A$  is defined as  $\|A\|_F = \sqrt{\sum_{i \in [n]} \sum_{j \in [d]} A_{i,j}^2}$ . For any PSD matrix  $M \in \mathbb{S}^{n \times n}$ , we define the  $M$ -norm as  $\|x\|_M = \sqrt{x^\top Mx}$ ,  $\forall x \in \mathbb{R}^n$ .

## 79:6 A Faster Interior-Point Method for Sum-Of-Squares Optimization

We use  $\mathcal{T}_{\text{mat}}(a, b, c)$  to denote the time to multiply two matrices of sizes  $a \times b$  and  $b \times c$ . A basic fact of fast matrix multiplication is that  $\mathcal{T}_{\text{mat}}(a, b, c) = \mathcal{T}_{\text{mat}}(b, c, a) = \mathcal{T}_{\text{mat}}(c, a, b)$  (see e.g. [6]), and we will use these three terms interchangeably.

► **Fact 2** (Woodbury identity). *Let  $A \in \mathbb{R}^{n \times n}$ ,  $C \in \mathbb{R}^{k \times k}$ ,  $U \in \mathbb{R}^{n \times k}$ ,  $V \in \mathbb{R}^{k \times n}$  where  $A$  and  $C$  are invertible, then*

$$(A + UCV)^{-1} = A^{-1} - A^{-1}U(C^{-1} + VA^{-1}U)^{-1}VA^{-1}.$$

► **Definition 3** (Hadamard product). *For any two matrices  $A, B \in \mathbb{R}^{m \times n}$ , the Hadamard product  $A \circ B$  is defined as*

$$(A \circ B)_{i,j} = A_{i,j} \cdot B_{i,j}, \quad \forall i \in [m], j \in [n].$$

We also use  $A^{\circ 2}$  to denote  $A \circ A$ .

The Hadamard product has the following properties (the proofs are straightforward).

► **Fact 4** (Properties of Hadamard product). *For matrices  $A, B \in \mathbb{R}^{m \times n}$ , and vectors  $x \in \mathbb{R}^m$ ,  $y \in \mathbb{R}^n$ , we have the following properties.*

1.  $x^\top (A \circ B)y = \text{tr}[\text{diag}(x)A \text{diag}(y)B^\top]$ ,
2.  $A \circ (x \cdot y^\top) = \text{diag}(x) \cdot A \cdot \text{diag}(y)$ .

► **Definition 5** (Spectral approximation). *For any two symmetric matrices  $A, \tilde{A} \in \mathbb{R}^{n \times n}$ , any parameter  $\epsilon \in (0, 1)$ , we say  $\tilde{A}$  and  $A$  are  $\epsilon$ -spectral approximation of each other, denoted as  $\tilde{A} \approx_\epsilon A$ , if we have*

$$e^{-\epsilon} \cdot x^\top Ax \leq x^\top \tilde{A}x \leq e^\epsilon \cdot x^\top Ax, \quad \forall x \in \mathbb{R}^n.$$

Spectral approximation has the following properties.

► **Fact 6** (Properties of spectral approximation). *For any two PSD matrices  $A, \tilde{A} \in \mathbb{R}^{n \times n}$ , any parameter  $\epsilon \in (0, 1)$ , if  $\tilde{A} \approx_\epsilon A$ , then we have*

1.  $B^\top AB \approx_\epsilon B^\top \tilde{A}B$ , for any matrix  $B \in \mathbb{R}^{n \times n}$ .
2. If both  $A$  and  $\tilde{A}$  are invertible, then  $A^{-1} \approx_\epsilon \tilde{A}^{-1}$ .
3.  $e^{-\epsilon} \text{tr}[A] \leq \text{tr}[\tilde{A}] \leq e^\epsilon \text{tr}[A]$ .
4.  $\tilde{A}^{\circ 2} \approx_{2\epsilon} A^{\circ 2}$ .

**Proof.** The proofs of the first three claims are straightforward. We only prove the last claim.

For any vector  $x \in \mathbb{R}^n$ , we have

$$\begin{aligned} x^\top A^{\circ 2}x &= \text{tr}[\text{diag}(x)A \text{diag}(x)A] \\ &= \text{tr}[A^{1/2} \text{diag}(x)A \text{diag}(x)A^{1/2}] \\ &\leq e^\epsilon \cdot \text{tr}[A^{1/2} \text{diag}(x)\tilde{A} \text{diag}(x)A^{1/2}] \\ &= e^\epsilon \cdot \text{tr}[\tilde{A}^{1/2} \text{diag}(x)A \text{diag}(x)\tilde{A}^{1/2}] \\ &\leq e^{2\epsilon} \cdot \text{tr}[\tilde{A}^{1/2} \text{diag}(x)\tilde{A} \text{diag}(x)\tilde{A}^{1/2}] = e^{2\epsilon} \cdot x^\top \tilde{A}^{\circ 2}x \end{aligned}$$

where the first step follows from Fact 4, the second and the fourth steps follow from the trace invariance under cyclic permutations and the fact that  $A^{1/2}$  exists when  $A$  is PSD, the third and the fifth steps follow from Part 3 of this fact.

Similarly we can prove  $x^\top A^{\circ 2}x \geq e^{-2\epsilon} \cdot x^\top \tilde{A}^{\circ 2}x$ . Thus we have  $A^{\circ 2} \approx_{2\epsilon} \tilde{A}^{\circ 2}$ . ◀

### 3 Background of sum-of-squares optimization

In this section we provide the background of sum-of-squares optimization. We refer the readers to [28, 27] for more details.

► **Definition 7** (Polynomial space). *We use  $\mathcal{V}_{n,d}$  to denote the set of  $n$ -variate polynomials over the reals of degree at most  $d$ , where the degree means the total degree, i.e., the degree of  $x_1^{d_1} \cdots x_n^{d_n}$  is  $\sum_{i=1}^n d_i$ .*

► **Definition 8** (Degree of polynomial space). *We define  $L := \dim(\mathcal{V}_{n,d}) = \binom{n+d}{n}$  and  $U := \dim(\mathcal{V}_{n,2d}) = \binom{n+2d}{n}$ .*

After fixing a basis  $(p_1, p_2, \dots, p_L)$  of  $\mathcal{V}_{n,d}$ , there exists a one-to-one correspondence between any polynomial  $p = \sum_{i=1}^L x_i \cdot p_i \in \mathcal{V}_{n,d}$  and the vector  $[x_1, x_2, \dots, x_L] \in \mathbb{R}^L$ . From now on when the basis is clear from context, we will use  $\mathcal{V}_{n,d}$  and  $\mathbb{R}^L$  interchangeably, and similarly  $\mathcal{V}_{n,2d}$  and  $\mathbb{R}^U$  interchangeably.

► **Definition 9** (SOS polynomials). *A polynomial  $p \in \mathcal{V}_{n,2d}$  is said to be a sum-of-squares (SOS) polynomial if  $p$  can be written as a sum of squares of polynomials, i.e.  $p = \sum_{i=1}^M q_i^2$  for some  $M \in \mathbb{N}$  and polynomials  $q_1, q_2, \dots, q_M \in \mathcal{V}_{n,d}$ .*

*We use  $\Sigma_{n,2d}$  to denote the set of  $n$ -variate SOS polynomials of degree at most  $2d$ .*

The set  $\Sigma_{n,2d}$  is a closed convex and pointed cone in  $\mathcal{V}_{n,2d}$  with non-empty interior (Theorem 17.1 of [24]). The SOS optimization problem requires the variable  $x \in \mathbb{R}^U$  to be in the SOS cone, and it is a special case of conic programming. Given a constraint matrix  $A \in \mathbb{R}^{m \times U}$  where  $m \leq U$ , and  $b \in \mathbb{R}^m$  and  $c \in \mathbb{R}^U$ , the SOS optimization can be written in the following primal-dual formulation:

$$\begin{array}{ll} \text{Primal:} & \min \langle c, x \rangle \\ \text{s.t.} & Ax = b \\ & x \in \Sigma_{n,2d}, \end{array} \quad \begin{array}{ll} \text{Dual:} & \max \langle y, b \rangle \\ \text{s.t.} & A^\top y + s = c \\ & s \in \Sigma_{n,2d}^*. \end{array} \quad (\text{SOS})$$

Here  $\Sigma_{n,2d}^* := \{s \in \mathbb{R}^U \mid s^\top x \geq 0, \forall x \in \Sigma_{n,2d}\}$  denotes the dual cone of  $\Sigma_{n,2d}$ .

Nesterov in [24] noted that the dual SOS cone allows the following characterization.

► **Theorem 10** (Dual cone characterization, Theorem 17.1 of [24]). *For any ordered bases  $\mathbf{p} = (p_1, \dots, p_L)$  and  $\mathbf{q} = (q_1, \dots, q_U)$  of  $\mathcal{V}_{n,d}$  and  $\mathcal{V}_{n,2d}$ , let  $\Lambda : \mathbb{R}^U \rightarrow \mathbb{R}^{L \times L}$  be the unique linear mapping satisfying  $\Lambda(\mathbf{q}) = \mathbf{p}\mathbf{p}^\top$ .<sup>4</sup> Then the dual cone  $\Sigma_{n,2d}^*$  admits the characterization under the bases  $\mathbf{p}$  and  $\mathbf{q}$ :*

$$\Sigma_{n,2d}^* = \{s \in \mathbb{R}^U \mid \Lambda(s) \succeq 0\}. \quad (8)$$

As barrier functions for the cone of positive semidefinite matrices are well-known, this also gives rise to a barrier function for the dual SOS cone. With the standard log-det barrier for the semidefinite cone, the following function  $F : \Sigma_{n,2d}^* \rightarrow \mathbb{R}$  is a barrier function for  $\Sigma_{n,2d}^*$ :

$$F(s) = -\log \det(\Lambda(s)).$$

Furthermore, the barrier parameter  $\nu_F$  of  $F(s)$  is bounded by the barrier parameter  $L$  of the original log-det barrier function ([24]).

<sup>4</sup> This equation means  $\forall t \in \mathbb{R}^n, \Lambda([q_1(t), \dots, q_U(t)]^\top) = [p_1(t), \dots, p_L(t)]^\top \cdot [p_1(t), \dots, p_L(t)]$ .



### Interpolant basis

The barrier function depends on the choice of the basis for both  $\mathcal{V}_{n,d}$  and  $\mathcal{V}_{n,2d}$ , as the linear map  $\Lambda$  depends on these two bases. We follow the approach of [27] and focus on the so-called *interpolant* bases, which generalises well to multivariate polynomials and is numerically stable.

For the vector space  $\mathcal{V}_{n,2d}$ , consider a set of *unisolvent* points  $\mathcal{T} = \{t_1, t_2, \dots, t_U\} \subseteq \mathbb{R}^n$ , which is a set points such that every polynomial in  $\mathcal{V}_{n,2d}$  is uniquely determined by its values on the points in  $\mathcal{T}$ . For univariate polynomials any set of  $U$  points suffices, but this does not hold anymore for the multivariate case. To also ensure numerical stability, the so called (approximate) Fekete points can be used as unisolvent points [32, 8].

The interpolant basis is defined as follows. Let us fix a set of unisolvent points  $\mathcal{T} = \{t_1, t_2, \dots, t_U\} \subseteq \mathbb{R}^n$ . Now every  $t_u \in \mathcal{T}$  implies a Lagrange polynomial  $q_u$  which is the unique polynomial that satisfies  $q_u(t_u) = 1$  and  $q_u(t_v) = 0$  for all  $t_v \neq t_u \in \mathcal{T}$ . The Lagrange polynomials form a basis  $\mathbf{q} = (q_1, \dots, q_U)$  of  $\mathcal{V}_{n,2d}$ . Choose any basis  $\mathbf{p} = (p_1, \dots, p_L)$  of  $\mathcal{V}_{n,d}$ . Define the matrix  $P \in \mathbb{R}^{U \times L}$  as

$$P_{u,\ell} = p_\ell(t_u), \quad \forall u \in [U], \ell \in [L].$$

By the definition of the Lagrange polynomials,  $p_i p_j = \sum_{u=1}^U p_i(t_u) p_j(t_u) q_u$ , so we have  $\mathbf{p}\mathbf{p}^\top = P^\top \text{diag}(\mathbf{q})P$ . Thus under the bases  $\mathbf{p}$  and  $\mathbf{q}$ , the linear map  $\Lambda : \mathbb{R}^U \rightarrow \mathbb{R}^{L \times L}$  takes on the following convenient form:

$$\Lambda(s) = P^\top \text{diag}(s)P. \quad (9)$$

For more details on how to pick the unisolvent points and how to construct  $P$ , we refer the readers to [27] and the references therein.

## 4 Algorithm

Since in this paper we focus on the theoretical running time of the algorithm, for simplicity we use the barrier method (see e.g. [30, Chapter 2]) instead of the more sophisticated Skajaa–Ye Algorithm used by [27].

The dual formulation of (SOS) is equivalent to the following optimization problem

$$\min -b^\top y \quad \text{s.t.} \quad y \in \overline{D}_F,$$

where with an abuse of the notation we define  $F : \mathbb{R}^m \rightarrow \mathbb{R}_+$  to be the barrier function

$$F(y) = -\log \det(\Lambda(c - A^\top y)) \quad (10)$$

for  $c \in \mathbb{R}^U$ , and  $A \in \mathbb{R}^{m \times U}$ , and  $\Lambda(s) = P^\top \text{diag}(s)P$  is the linear operator defined in Eq (9).  $D_F \subseteq \mathbb{R}^m$  is the domain of  $F$ , and  $\overline{D}_F$  is the closure of  $D_F$ .

The barrier parameter of the barrier function  $F$  is  $\nu_F = L$ . The gradient and the Hessian of the barrier function  $F$  are (define  $s := c - A^\top y$ ):

$$\begin{aligned} g(y) &= A \cdot \text{diag} \left( P(P^\top \text{diag}(s)P)^{-1} P^\top \right), \\ H(y) &= A \cdot \left( P(P^\top \text{diag}(s)P)^{-1} P^\top \right)^{\circ 2} \cdot A^\top. \end{aligned}$$

For any  $\eta > 0$ , define a function  $F_\eta : \mathbb{R}^m \rightarrow \mathbb{R}$ :

$$F_\eta(y) = -\eta \cdot b^\top y + F(y).$$



The gradient and the Hessian of  $F_\eta(y)$  are:

$$\begin{aligned} g_\eta(y) &= -\eta \cdot b + A \cdot \text{diag} \left( P(P^\top \text{diag}(s)P)^{-1} P^\top \right), \\ H_\eta(y) &= A \cdot \left( P(P^\top \text{diag}(s)P)^{-1} P^\top \right)^{\circ 2} \cdot A^\top. \end{aligned}$$

Note that  $H_\eta(y) = H(y)$  for any  $\eta$ .

In each iteration the barrier method increases  $\eta$  by a factor of  $1 + \Theta(\frac{1}{\sqrt{L}})$ , and it performs a Newton step

$$y \leftarrow y - H_\eta(y)^{-1} \cdot g_\eta(y).$$

By standard IPM theory it suffices to use a spectral approximation of the Hessian matrix in the Newton step. For more details see e.g. [30].

The main technical part of our algorithm is to efficiently maintain a matrix  $N$  that is the spectral approximation of the inverse of the Hessian matrix. To do this, we maintain another matrix  $\tilde{S}$  that is a spectral approximation of  $S := P^\top \text{diag}(s)P$ , and we use the subroutine `LOWRANKUPDATE` (Algorithm 3, Lemma 14) to update  $\tilde{S}$ . After  $\tilde{S}$  is updated, we use another subroutine `UPDATEHESSIANINV` (Algorithm 2, Lemma 11) to update  $N$ . A complete description of our algorithm can be found in Algorithm 1.

## 5 Updating Hessian inverse efficiently

In this section we prove how to update the Hessian inverse efficiently. We present the algorithm `UPDATEHESSIANINV` in Algorithm 2.

► **Lemma 11** (Hessian inverse update). *In the algorithm `UPDATEHESSIANINV` (Algorithm 2), the inputs are the maintained matrices  $T, N$  and the updates  $V_1, V_2 \in \mathbb{R}^{L \times r}$  where  $r$  satisfies  $Lr \leq U$ . The inputs satisfy that for some  $\tilde{S} \in \mathbb{S}^{L \times L}$ ,*

$$\begin{aligned} T &= \tilde{S}^{-1} \in \mathbb{R}^{L \times L}, \\ N &= (A \cdot (P\tilde{S}^{-1}P^\top)^{\circ 2} \cdot A^\top)^{-1} \in \mathbb{R}^{m \times m}, \end{aligned}$$

Let  $\tilde{S}^{\text{new}} = \tilde{S} + V_1 V_2^\top$ . The algorithm outputs two matrices  $T^{\text{new}}, N^{\text{new}}$  such that

$$\begin{aligned} T^{\text{new}} &= (\tilde{S}^{\text{new}})^{-1} \in \mathbb{R}^{L \times L}, \\ N^{\text{new}} &= (A \cdot (P(\tilde{S}^{\text{new}})^{-1}P^\top)^{\circ 2} \cdot A^\top)^{-1} \in \mathbb{R}^{m \times m}. \end{aligned}$$

Furthermore, the algorithm takes  $O(\mathcal{T}_{\text{mat}}(U, U, Lr))$  time.

**Proof.** We first prove the correctness by analyzing each step of the algorithm.

**Step 1.** Compute  $\bar{V}_1, \bar{V}_2 \in \mathbb{R}^{L \times r}$  and  $T^{\text{new}} \in \mathbb{R}^{L \times L}$ .

$$\begin{aligned} T^{\text{new}} &= T + \bar{V}_1 \cdot \bar{V}_2^\top \\ &= T - TV_1 \cdot (I + V_2^\top TV_1)^{-1} \cdot V_2^\top T^\top \\ &= (\tilde{S} + V_1 V_2^\top)^{-1} = (\tilde{S}^{\text{new}})^{-1}, \end{aligned}$$

where the first two steps follow from algorithm description, the third step follows from the Woodbury identity (Fact 2) and  $T = \tilde{S}^{-1}$ .

Thus  $T^{\text{new}}$  satisfies the requirement of the output.

■ **Algorithm 1** Main SOS algorithm.

---

**Parameters:**  $\delta \in (0, 1)$ ,  $\epsilon_N \in (0, 0.05)$ ,  $\alpha = \frac{\epsilon_N}{20\sqrt{L}}$ ,  $t = 40\epsilon_N^{-1}\sqrt{L}\log(L/\delta)$ .

**Input** :  $A \in \mathbb{R}^{m \times U}$ ,  $b \in \mathbb{R}^m$ ,  $c \in \mathbb{R}^U$

**Output** : A near feasible and optimal solution.

- 1 Construct  $P \in \mathbb{R}^{U \times L}$  of the interpolant basis. Convert  $A, b, c$  to the interpolant basis.
- 2 Use Lemma 22 to obtain a modified dual SOS optimization problem which has an initial solution  $(y, s) \in \mathbb{R}^m \times \mathbb{R}^U$  that is optimal for  $F_\eta$ , where  $\eta = 1$ .
- 3  $\tilde{S} \leftarrow S \leftarrow P^\top \text{diag}(s)P$ ; //  $\tilde{S}, S \in \mathbb{R}^{L \times L}$
- 4  $T \leftarrow S^{-1}$ ; //  $T \in \mathbb{R}^{L \times L}$
- 5  $N \leftarrow (A(PTP^\top)^{\circ 2}A^\top)^{-1}$ ; //  $N \in \mathbb{R}^{m \times m}$
- 6  $g \leftarrow -\eta \cdot b + A \cdot \text{diag}(P(P^\top \text{diag}(s)P)^{-1}P^\top)$ ; //  $g \in \mathbb{R}^m$
- 7 **for**  $i = 1, 2, \dots, t$  **do**
- 8  $\delta_y \leftarrow -N \cdot g$ ; //  $\delta_y \in \mathbb{R}^m$
- 9  $y^{\text{new}} \leftarrow y + \delta_y$ ; //  $y^{\text{new}} \in \mathbb{R}^m$
- 10  $s^{\text{new}} \leftarrow c - A^\top y^{\text{new}}$ ; //  $s^{\text{new}} \in \mathbb{R}^U$
- 11  $\eta^{\text{new}} \leftarrow \eta \cdot (1 + \alpha)$ ;
- 12  $S^{\text{new}} \leftarrow P^\top \text{diag}(s^{\text{new}})P$ ; //  $S^{\text{new}} \in \mathbb{R}^{L \times L}$
- 13  $\tilde{S}^{\text{new}}, V_1, V_2 \leftarrow \text{LOWRANKUPDATE}(S^{\text{new}}, \tilde{S})$ ;
- 14 // Lemma 14,  $\tilde{S}^{\text{new}} \in \mathbb{R}^{L \times L}$ ,  $V_1, V_2 \in \mathbb{R}^{L \times r_i}$  or  $V_1 = V_2 = \text{null}$
- 15 **if**  $V_1 = V_2 = \text{null}$  **then**
- 16  $T^{\text{new}} \leftarrow (\tilde{S}^{\text{new}})^{-1}$ ; //  $T^{\text{new}} \in \mathbb{R}^{L \times L}$
- 17  $N^{\text{new}} \leftarrow (A \cdot (PT^{\text{new}}P^\top)^{\circ 2} \cdot A^\top)^{-1}$ ; //  $N^{\text{new}} \in \mathbb{R}^{m \times m}$
- 18 **else**
- 19  $T^{\text{new}}, N^{\text{new}} \leftarrow \text{UPDATEHESSIANINV}(T, N, V_1, V_2)$ ;
- 20 // Lemma 11,  $T^{\text{new}} \in \mathbb{R}^{L \times L}$ ,  $N^{\text{new}} \in \mathbb{R}^{m \times m}$
- 21  $g^{\text{new}} \leftarrow -\eta^{\text{new}} \cdot b + A \cdot \text{diag}(P(P^\top \text{diag}(s^{\text{new}})P)^{-1}P^\top)$ ; //  $g^{\text{new}} \in \mathbb{R}^m$
- 22  $(\eta, y, s, \tilde{S}, T, N, g) \leftarrow (\eta^{\text{new}}, y^{\text{new}}, s^{\text{new}}, \tilde{S}^{\text{new}}, T^{\text{new}}, N^{\text{new}}, g^{\text{new}})$ ;
- 23 **return**  $(y, s)$

---

**Step 2.** Compute  $Y', Z' \in \mathbb{R}^{U \times (L+r)}$  and  $Y, Z \in \mathbb{R}^{U \times (L+r)r}$ . We prove that  $Y$  and  $Z$  satisfy  $(PTP^\top)^{\circ 2} + Y \cdot Z^\top = (PT^{\text{new}}P^\top)^{\circ 2}$ :

$$\begin{aligned}
(PTP^\top)^{\circ 2} + Y \cdot Z^\top &= (PTP^\top)^{\circ 2} + \sum_{i=1}^r \text{diag}(u_i) \cdot (Y' \cdot (Z')^\top) \cdot \text{diag}(v_i) \\
&= (PTP^\top)^{\circ 2} + (Y' \cdot (Z')^\top) \circ \left( \sum_{i=1}^r u_i \cdot v_i^\top \right) \\
&= (PTP^\top)^{\circ 2} + (2PTP^\top + (P\bar{V}_1) \cdot (P\bar{V}_2)^\top) \circ ((P\bar{V}_1) \cdot (P\bar{V}_2)^\top) \\
&= (PTP^\top + (P\bar{V}_1) \cdot (P\bar{V}_2)^\top)^{\circ 2} \\
&= (PT^{\text{new}}P^\top)^{\circ 2},
\end{aligned}$$

where the first step follows from the algorithm description of  $Y$  and  $Z$ , the second step follows from Part 2 of Fact 4 that  $\text{diag}(x) \cdot A \cdot \text{diag}(y) = A \circ (x \cdot y^\top)$ , the third step follows from  $Y' \cdot (Z')^\top = 2PTP^\top + (P\bar{V}_1) \cdot (P\bar{V}_2)^\top$  and  $(P\bar{V}_1) \cdot (P\bar{V}_2)^\top = \sum_{i=1}^r u_i \cdot v_i^\top$  (see algorithm description of  $Y'$  and  $Z'$ ), the last step follows from  $T^{\text{new}} = T + \bar{V}_1 \cdot \bar{V}_2^\top$ .

■ **Algorithm 2** UPDATEHESSIANINV.

---

**Input** :  $T \in \mathbb{R}^{L \times L}$ ,  $N \in \mathbb{R}^{m \times m}$ ,  $V_1, V_2 \in \mathbb{R}^{L \times r}$   
**Output** :  $T^{\text{new}} \in \mathbb{R}^{L \times L}$ ,  $N^{\text{new}} \in \mathbb{R}^{m \times m}$

---

```

1 // Step 1
2  $\bar{V}_1 \leftarrow -TV_1 \cdot (I + V_2^\top TV_1)^{-1}$ ; //  $\bar{V}_1 \in \mathbb{R}^{L \times r}$ 
3  $\bar{V}_2 \leftarrow TV_2$ ; //  $\bar{V}_2 \in \mathbb{R}^{L \times r}$ 
4  $T^{\text{new}} \leftarrow T + \bar{V}_1 \cdot \bar{V}_2^\top$ ; //  $T^{\text{new}} \in \mathbb{R}^{L \times L}$ 
5 // Step 2
6  $Y' \leftarrow [2PT, P\bar{V}_1]$ ; //  $Y' \in \mathbb{R}^{U \times (L+r)}$ 
7  $Z' \leftarrow [P, P\bar{V}_2]$ ; //  $Z' \in \mathbb{R}^{U \times (L+r)}$ 
8  $Y \leftarrow [\text{diag}(u_1)Y', \dots, \text{diag}(u_r)Y']$ ,  $u_i$  is the  $i$ -th column of  $P\bar{V}_1$ ; //  $Y \in \mathbb{R}^{U \times (L+r)r}$ 
9  $Z \leftarrow [\text{diag}(v_1)Z', \dots, \text{diag}(v_r)Z']$ ,  $v_i$  is the  $i$ -th column of  $P\bar{V}_2$ ; //  $Z \in \mathbb{R}^{U \times (L+r)r}$ 
10 // Step 3
11  $N^{\text{new}} \leftarrow N - N \cdot (AY) \cdot (I + (AZ)^\top N(AY))^{-1} \cdot (AZ)^\top \cdot N$ ; //  $N^{\text{new}} \in \mathbb{R}^{m \times m}$ 
12 return  $T^{\text{new}}, N^{\text{new}}$ 

```

---

**Step 3.** Compute  $N^{\text{new}} \in \mathbb{R}^{m \times m}$ .

$$\begin{aligned}
N^{\text{new}} &= N - N \cdot (AY) \cdot (I + (AZ)^\top N(AY))^{-1} \cdot (AZ)^\top \cdot N \\
&= (A \cdot (PTP^\top)^{\circ 2} \cdot A^\top + (AY) \cdot (AZ)^\top)^{-1} \\
&= (A \cdot (PT^{\text{new}}P^\top)^{\circ 2} \cdot A^\top)^{-1},
\end{aligned}$$

where the first step follows from the algorithm description of  $N^{\text{new}}$ , the second step follows from  $N = (A \cdot (PTP^\top)^{\circ 2} \cdot A^\top)^{-1}$  and the Woodbury identity (Fact 2), and the last step follows from  $(PT^{\text{new}}P^\top)^{\circ 2} = (PTP^\top)^{\circ 2} + Y \cdot Z^\top$ .

Thus  $N^{\text{new}}$  satisfies the requirement of the output.

**Time complexity.** It is easy to see that the most time-consuming step is to compute  $N^{\text{new}}$  on Line 11, and in total this step takes  $O(\mathcal{T}_{\text{mat}}(m, U, Lr) + \mathcal{T}_{\text{mat}}(m, m, Lr) + (Lr)^\omega)$  time.

Since  $Lr \leq U$  and  $m \leq U$ , overall this algorithm takes at most  $O(\mathcal{T}_{\text{mat}}(U, U, Lr))$  time. ◀

## 6 Correctness

### 6.1 Standard results from IPM theory

We use the following two results of the barrier method that hold for any cone with a barrier function. The proofs are standard, (see e.g., [30, Section 2.4]). For completeness we include a proof in the full version.

► **Lemma 12** (Invariance of Newton step, [30]). *Consider the following optimization problem:  $\min -b^\top y$  s.t.  $y \in \bar{D}_F$ , where  $F : \mathbb{R}^m \rightarrow \mathbb{R}_+$  is a barrier function with barrier parameter  $\nu_F$ ,  $D_F \subseteq \mathbb{R}^m$  is the domain of  $F$ , and  $\bar{D}_F$  is the closure of  $D_F$ . For any  $\eta \geq 1$ , define  $F_\eta(y) = -\eta b^\top y + F(y)$ . Let  $g_\eta(y) \in \mathbb{R}^m$  and  $H(y) \in \mathbb{R}^{m \times m}$  denote the gradient and the Hessian of  $F_\eta$  at  $y$ .*

*Let  $0 < \epsilon_N \leq 0.05$  be a parameter. If a feasible solution  $y \in D_F$ , a parameter  $\eta > 0$ , and a positive definite matrix  $\tilde{H} \in \mathbb{S}^{n \times n}$  satisfy the following:*

$$\|g_\eta(y)\|_{H(y)^{-1}} \leq \epsilon_N, \quad \tilde{H} \approx_{0.02} H(y).$$

*Then  $\eta^{\text{new}} = \eta \cdot (1 + \frac{\epsilon_N}{20\sqrt{\nu_F}})$ ,  $y^{\text{new}} = y + \delta_y$  where  $\delta_y = -\tilde{H}^{-1}g_{\eta^{\text{new}}}(y)$  satisfy  $y^{\text{new}} \in D_F$  and*

$$\|\delta_y\|_{H(y)} \leq 2\epsilon_N, \quad \|g_{\eta^{\text{new}}}(y^{\text{new}})\|_{H(y^{\text{new}})^{-1}} \leq \epsilon_N.$$

► **Lemma 13** (Approximate optimality, [30]). *Consider the following optimization problem:  $\min -b^\top y$  s.t.  $y \in \overline{D}_F$ , where  $F : \mathbb{R}^m \rightarrow \mathbb{R}_+$  is a barrier function with barrier parameter  $\nu_F$ ,  $D_F \subseteq \mathbb{R}^m$  is the domain of  $F$ , and  $\overline{D}_F$  is the closure of  $D_F$ . Let OPT be the optimal objective value of this optimization problem. For any  $\eta \geq 1$ , define  $F_\eta(y) = -\eta b^\top y + F(y)$ . Let  $g_\eta(y) \in \mathbb{R}^m$  and  $H(y) \in \mathbb{R}^{m \times m}$  denote the gradient and the Hessian of  $F_\eta$  at  $y$ .*

*Let  $0 < \epsilon_N \leq 0.05$ . If a feasible solution  $y \in D_F$  satisfies  $\|g_\eta(y)\|_{H(y)^{-1}} \leq \epsilon_N$ , then we have  $-b^\top y \leq \text{OPT} + \frac{\nu_F}{\eta} \cdot (1 + 2\epsilon_N)$ .*

## 6.2 Low rank update

We use the following low rank update procedure of [16] and [15], which we modify by using a cutoff when  $r \geq U/L$ . The proof of the following lemma can be found in [15, Theorem 10.8].

■ **Algorithm 3** LOWRANKUPDATE of [16].

---

**Parameters:** A real number  $\epsilon_S < 0.01$ .

**Input** : New exact matrix  $S^{\text{new}} \in \mathbb{R}^{L \times L}$ , old approximate matrix  $\tilde{S} \in \mathbb{R}^{L \times L}$ ,

**Output** : New approximate matrix  $\tilde{S}^{\text{new}} \in \mathbb{R}^{L \times L}$ , update matrices  $V_1, V_2 \in \mathbb{R}^{L \times r}$ .

- 1  $Z_{\text{mid}} \leftarrow (S^{\text{new}})^{-1/2} \tilde{S} (S^{\text{new}})^{-1/2} - I$
- 2 Compute spectral decomposition  $Z_{\text{mid}} = X \text{diag}(\lambda) X^\top$
- 3 Let  $\pi : [L] \rightarrow [L]$  be a sorting permutation such that  $|\lambda_{\pi(i)}| \geq |\lambda_{\pi(i+1)}|$ .
- 4 **if**  $|\lambda_{\pi(1)}| \leq \epsilon_S$  **then**
- 5      $\tilde{S}_{\text{new}} \leftarrow \tilde{S}$ ;
- 6     **return**  $(\tilde{S}^{\text{new}}, 0, 0)$
- 7 **else**
- 8      $r \leftarrow 1$ ;
- 9     **while**  $2r \leq U/L$  and  $(|\lambda_{\pi(2r)}| > \epsilon_S$  or  $|\lambda_{\pi(2r)}| > (1 - 1/\log L)|\lambda_{\pi(r)}|)$  **do**
- 10          $r \leftarrow r + 1$ ;
- 11      $r \leftarrow 2r$ ;
- 12     **if**  $r \geq U/L$  **then**
- 13          $\tilde{S}^{\text{new}} \leftarrow S^{\text{new}}$  ;   // Here we deviate from [16]
- 14         **return**  $(\tilde{S}^{\text{new}}, \text{null}, \text{null})$
- 15     **else**
- 16          $\lambda_{\pi(i)}^{\text{new}} \leftarrow \begin{cases} 0 & \text{if } i = 1, 2, \dots, r \\ \lambda_{\pi(i)} & \text{else} \end{cases}$
- 17          $\Omega \leftarrow \text{supp}(\lambda^{\text{new}} - \lambda)$  ;   //  $|\Omega| = r$
- 18          $V_1 \leftarrow ((S^{\text{new}})^{1/2} \cdot X \cdot \text{diag}(\lambda^{\text{new}} - \lambda))_{:, \Omega}$  ;                     //  $V_1 \in \mathbb{R}^{L \times r}$
- 19          $V_2 \leftarrow ((S^{\text{new}})^{1/2} \cdot X)_{:, \Omega}$  ;   //  $V_2 \in \mathbb{R}^{L \times r}$
- 20          $\tilde{S}^{\text{new}} \leftarrow \tilde{S} + (S^{\text{new}})^{1/2} X \text{diag}(\lambda^{\text{new}} - \lambda) X^\top (S^{\text{new}})^{1/2}$  ;
- 21         //  $\tilde{S}^{\text{new}} = \tilde{S} + V_1 V_2^\top \in \mathbb{R}^{L \times L}$
- 21         **return**  $(\tilde{S}^{\text{new}}, V_1, V_2)$ ;

---

► **Lemma 14** (Low rank update). *The algorithm LOWRANKUPDATE (Algorithm 3) has the following properties:*

(i) *The output matrix  $\tilde{S}^{\text{new}} = \tilde{S} + V_1 V_2^\top$  is a spectral approximation of the input matrix:*

$$\tilde{S}^{\text{new}} \approx_{\epsilon_S} S^{\text{new}}.$$

- (ii) Consider a total of  $t$  iterations of `LOWRANKUPDATE`. Initially  $\tilde{S}^{(0)} = S^{(0)}$ , and we use  $(S^{(i)}, \tilde{S}^{(i-1)})$  and  $(\tilde{S}^{(i)}, V_1^{(i)}, V_2^{(i)})$  to denote the input and the output of the  $i$ -th iteration. We define the rank  $r_i$  to be the rank of  $V_1^{(i)}$  if  $V_1^{(i)} \neq \mathbf{null}$ , and otherwise we define  $r_i = U/L$ .

If the input exact matrices  $S^{(0)}, S^{(1)}, \dots, S^{(t)} \in \mathbb{R}^{L \times L}$  satisfy

$$\|(S^{(i-1)})^{-1/2} S^{(i)} (S^{(i-1)})^{-1/2} - I\|_F \leq 0.02, \quad \forall i \in [t]. \quad (11)$$

Then for any non-increasing sequence  $g \in \mathbb{R}_+^L$ , the ranks  $r_i$  satisfy

$$\sum_{i=1}^t r_i \cdot g_{r_i} \leq O(t \cdot \|g\|_2 \cdot \log L).$$

Furthermore, the algorithm `LOWRANKUPDATE` takes  $O(L^\omega)$  time.

### 6.3 Slowly moving guarantee

In SOS optimization, the matrix  $S = P^\top \text{diag}(s)P$  corresponds to the slack matrix of the SDP. The following lemma proves similar to SDP, in SOS the matrix  $S$  is changing slowly. Using this lemma we will prove that the requirement Eq. (11) of Lemma 14 is satisfied, which means we can approximate the change to the slack by a low-rank matrix.

► **Lemma 15** (Slowly moving guarantee). *Let  $c \in \mathbb{R}^U$  and  $A \in \mathbb{R}^{m \times U}$  be the input to the optimization problem. Let  $P \in \mathbb{R}^{U \times L}$  be the matrix of the interpolant basis.*

*For any  $y \in \mathbb{R}^m$  and  $y^{\text{new}} = y + \delta_y \in \mathbb{R}^m$ , let  $s = c - A^\top y \in \mathbb{R}^U$  and  $S = P^\top \text{diag}(s)P \in \mathbb{R}^{L \times L}$ . Similarly define  $s^{\text{new}}$  and  $S^{\text{new}}$  from  $y^{\text{new}}$ . Let  $H(y) = A \cdot (P(P^\top \text{diag}(s)P)^{-1} P^\top)^{\circ 2}$ .  $A^\top \in \mathbb{R}^{m \times m}$ . If  $s, s^{\text{new}} \in \Sigma_{n,2d}^*$ , then  $S$  and  $S^{\text{new}}$  are both PSD, and we have*

$$\|S^{-1/2} S^{\text{new}} S^{-1/2} - I\|_F = \|\delta_y\|_{H(y)}.$$

**Proof.** Note that if  $s, s^{\text{new}} \in \Sigma_{n,2d}^*$ , then by the dual cone characterization (Theorem 10)  $S$  and  $S^{\text{new}}$  are both PSD.

For convenience we define  $M = PS^{-1}P^\top \in \mathbb{R}^{U \times U}$ . Note that  $H(y) = A \cdot M^{\circ 2} \cdot A^\top$ . We also define  $\delta_s = s^{\text{new}} - s = -A^\top \delta_y$ .  $\forall u \in [U]$ , we use  $p_u \in \mathbb{R}^L$  to denote the  $u$ -th row of  $P$ .

$$\begin{aligned} \|S^{-1/2} S^{\text{new}} S^{-1/2} - I\|_F^2 &= \|S^{-1/2} (S^{\text{new}} - S) S^{-1/2}\|_F^2 \\ &= \|S^{-1/2} (P \text{diag}(\delta_s) P^\top) S^{-1/2}\|_F^2 \\ &= \text{tr} (S^{-1} (P^\top \text{diag}(\delta_s) P) S^{-1} (P^\top \text{diag}(\delta_s) P)) \\ &= \text{tr} (S^{-1} (\sum_{u \in U} (\delta_s)_u \cdot p_u p_u^\top) S^{-1} (\sum_{v \in U} (\delta_s)_v \cdot p_v p_v^\top)) \\ &= \sum_{u,v \in U} (\delta_s)_u (\delta_s)_v \cdot \text{tr} (S^{-1} p_u p_u^\top S^{-1} p_v p_v^\top) \\ &= \sum_{u,v \in U} (\delta_s)_u (\delta_s)_v \cdot (p_v^\top S^{-1} p_u)^2 \\ &= \sum_{u,v \in U} (\delta_s)_u (\delta_s)_v \cdot M_{uv}^2 = \|\delta_s\|_{M^{\circ 2}}^2, \end{aligned} \quad (12)$$

where the third step follows from  $\|A\|_F^2 = \text{tr}(A^\top A)$  and the cyclic property of trace, and the sixth step again follows from the cyclic property of trace.

Since  $\delta_s = -A^\top \delta_y$ , we have

$$\|\delta_s\|_{M^{\circ 2}}^2 = \delta_s^\top M^{\circ 2} \delta_s = \delta_y^\top A M^{\circ 2} A^\top \delta_y = \delta_y^\top H(y) \delta_y = \|\delta_y\|_{H(y)}^2. \quad (13)$$

Combining Eq. (12) and (13) we get the bound in the lemma statement. ◀

## 6.4 Proof of correctness

Finally we are ready to prove the correctness of Algorithm 1.

► **Theorem 16** (Correctness of Algorithm 1). *Consider the following optimization problem with  $A \in \mathbb{R}^{m \times U}$ ,  $b \in \mathbb{R}^m$ , and  $c \in \mathbb{R}^U$ :*

$$\begin{array}{ll} \text{Primal:} & \min \langle c, x \rangle \\ & \text{s.t. } Ax = b \\ & x \in \Sigma_{n,2d}, \end{array} \quad \begin{array}{ll} \text{Dual:} & \max \langle y, b \rangle \\ & \text{s.t. } A^\top y + s = c \\ & s \in \Sigma_{n,2d}^*. \end{array}$$

Let  $\text{OPT}$  denote the optimal objective value of this optimization problem. Assume Slater's condition and that any primal feasible  $x \in \Sigma_{n,2d}$  satisfies  $\|x\|_1 \leq R$ .

Then for any error parameters  $\delta \in (0, 1)$ ,  $\epsilon_S \leq 0.01$ , and  $\epsilon_N \leq 0.05$ , Algorithm 1 outputs  $x \in \Sigma_{n,2d}$  that satisfies

$$\langle c, x \rangle \leq \text{OPT} + \delta \cdot R \|c\|_\infty \quad \text{and} \quad \|Ax - b\|_1 \leq 8\delta L \cdot (LR\|A\|_\infty + \|b\|_1).$$

**Proof.** We consider the optimization problem  $\min -b^\top y$  s.t.  $y \in \bar{D}_F$ , where  $F: \mathbb{R}^m \rightarrow \mathbb{R}_+$  is the barrier function defined in Eq. (10), and  $\bar{D}_F$  is the closure of the domain of  $F$ . The barrier parameter of  $F$  is  $\nu_F = L$ . This optimization problem is equivalent to the dual formulation and its optimal value is  $-\text{OPT}$ . For any  $\eta$ , let  $F_\eta(y) = -\eta b^\top y + F(y)$ .

In the beginning Algorithm 1 first uses Lemma 22 to convert the optimization problem to another form which has an initial feasible solution  $y$  that is close to the optimal solution of  $F_\eta$  with  $\eta = 1$ . The initial  $y$  satisfies  $\|g_\eta(y)\|_{H(y)^{-1}} \leq \epsilon_N$  by Lemma 22. Initially we also have  $\tilde{S} = S = P^\top \text{diag}(s)P$  (Line 3 in Algorithm 1).

Next we prove the correctness of Algorithm 1 inductively. At each iteration, we assume the following induction hypothesis is satisfied: (1)  $\|g_\eta(y)\|_{H(y)^{-1}} \leq \epsilon_N$ , (2)  $\tilde{S} \approx_{\epsilon_S} S$ . We aim to prove that the updated  $y^{\text{new}}$ ,  $\eta^{\text{new}}$ ,  $S^{\text{new}}$ , and  $\tilde{S}^{\text{new}}$  still satisfy these two conditions.

In Lemma 11 we have proved that in Algorithm 1 we always maintain  $N = (A \cdot (P\tilde{S}^{-1}P^\top)^{\circ 2} \cdot A^\top)^{-1}$ . Let  $\tilde{H} = N^{-1}$ , we have

$$\tilde{H} = A \cdot (P\tilde{S}^{-1}P^\top)^{\circ 2} \cdot A^\top \approx_{2\epsilon_S} A \cdot (PS^{-1}P^\top)^{\circ 2} \cdot A^\top = H(y),$$

where in the second step we use the induction hypothesis that  $\tilde{S} \approx_{\epsilon_S} S$ , and by Fact 6 we have  $\tilde{S}^{-1} \approx_{\epsilon_S} S^{-1}$ , and hence  $P\tilde{S}^{-1}P^\top \approx_{\epsilon_S} PS^{-1}P^\top$ , and hence  $(P\tilde{S}^{-1}P^\top)^{\circ 2} \approx_{2\epsilon_S} (PS^{-1}P^\top)^{\circ 2}$ .

The new vector  $y^{\text{new}}$  is computed as  $y^{\text{new}} = y + \delta_y$  where  $\delta_y = -\tilde{H}^{-1}g_\eta(y)$  (Line 8 and 9 of Algorithm 1). And  $\eta$  is updated to  $\eta^{\text{new}} = \eta \cdot (1 + \frac{\epsilon_N}{20\sqrt{L}})$  (Line 11 of Algorithm 1). Since  $\|g_\eta(y)\|_{H(y)^{-1}} \leq \epsilon_N$ , and  $\tilde{H} \approx_{2\epsilon_S} H(y)$  where  $2\epsilon_S \leq 0.02$  by its definition in Algorithm 3, the requirements of Lemma 12 are satisfied, so we have

$$\|g_{\eta^{\text{new}}}(y^{\text{new}})\|_{H(y^{\text{new}})^{-1}} \leq \epsilon_N, \quad \text{and} \quad \|\delta_y\|_{H(y)} \leq 2\epsilon_N.$$

This proves the first induction hypothesis.

Then using Lemma 15 and since  $\epsilon_N \leq 0.01$  by its definition in Algorithm 1, we have

$$\|S^{-1/2}(S^{\text{new}})S^{-1/2} - I\|_F \leq \|\delta_y\|_{H(y)} \leq 2\epsilon_N \leq 0.02.$$

Thus the input matrix  $S^{\text{new}}$  to `LOWRANKUPDATE` satisfies the requirement of Eq. (11) of Lemma 14, and we have that  $\tilde{S}^{\text{new}} \approx_{\epsilon_S} S^{\text{new}}$ . This proves the second induction hypothesis.

Finally, we know that after  $t = 40\epsilon_N^{-1}\sqrt{L}\log(L/\delta)$  iterations,  $\eta$  becomes  $(1 + \frac{\epsilon_N}{20\sqrt{L}})^t \geq 2L/\delta^2$ , so using Lemma 13, we have

$$b^\top y \geq \text{OPT} - \frac{\nu_F}{\eta} \cdot (1 + 2\epsilon_N) \geq \text{OPT} - \delta^2.$$

Thus the initialization lemma (Lemma 22) ensures that we have a solution  $x \in \Sigma_{n,2d}$  to the original primal optimization problem which satisfies

$$\langle c, x \rangle \leq \text{OPT} + \delta \cdot R\|c\|_\infty, \quad \|Ax - b\|_1 \leq 8\delta L \cdot (LR\|A\|_\infty + \|b\|_1). \quad \blacktriangleleft$$

## 7 Time complexity

### 7.1 Worst case time

We first bound the worst case running time of Algorithm 1. The running time of the  $i$ -th iteration depends on the updated rank  $r_i$  of LOWRANKUPDATE, which is defined to be the size of  $V_1^{(i)}$  if  $V_1^{(i)} \neq \text{null}$ , and  $U/L$  otherwise (see Lemma 14).

► **Lemma 17** (Worst case time of Algorithm 1). *In Algorithm 1, the initialization time is  $O(U^\omega)$ , and the running time in the  $i$ -th iteration is  $O(\mathcal{T}_{\text{mat}}(U, U, \min\{Lr_i, U\}))$ .*

**Proof.**

**Initialization time.** The most time-consuming step of initialization is Line 5, where computing  $N = (A(PTP^\top)^{\circ 2}A^\top)^{-1}$  takes  $O(\mathcal{T}_{\text{mat}}(U, U, L) + \mathcal{T}_{\text{mat}}(U, U, m))$  time. This is bounded by  $O(U^\omega)$  since  $L, m \leq U$ .

**Time per iteration.** In each iteration the most time-consuming steps are (1) computing  $S^{\text{new}}$  and calling LOWRANKUPDATE on Line 12-14, (2) executing the if-clause on Line 15-20, and (3) computing  $g^{\text{new}}$  on Line 21.

1. Computing  $S^{\text{new}}$  on Line 12 takes  $\mathcal{T}_{\text{mat}}(U, L, L)$  time. Calling LOWRANKUPDATE on Line 14 takes  $O(L^\omega)$  time by Lemma 14.
2. In the if-clause on Line 15-20, if  $V_1 = V_2 = \text{null}$ , then  $Lr_i \geq U$ , and we compute  $N^{\text{new}} = (A \cdot (PT^{\text{new}}P^\top)^{\circ 2} \cdot A^\top)^{-1}$ , which takes  $O(\mathcal{T}_{\text{mat}}(U, U, U))$  time. Otherwise we call UPDATEHESSIANINV on Line 19, which takes  $O(\mathcal{T}_{\text{mat}}(U, U, Lr_i))$  time by Lemma 11. In total the if-clause has running time  $O(\mathcal{T}_{\text{mat}}(U, U, \min\{Lr_i, U\}))$ .
3. Computing the gradient  $g = -\eta^{\text{new}} \cdot b + A \cdot \text{diag}(P(P^\top \text{diag}(s^{\text{new}})P)^{-1}P^\top)$  on Line 21 takes  $O(\mathcal{T}_{\text{mat}}(U, U, L))$  time since  $m \leq U$ .

Thus the total time per iteration is  $O(\mathcal{T}_{\text{mat}}(U, U, \min\{Lr_i, U\}))$ . ◀

### 7.2 Amortized time

In this section we bound the amortized running time of Algorithm 1.

Let  $\omega$  be the matrix multiplication exponent, let  $\alpha$  be the dual matrix multiplication exponent. The current best values are  $\omega \approx 2.373$  and  $\alpha \approx 0.314$  [21, 10, 1]. Note that the current best values of  $\omega$  and  $\alpha$  satisfies that  $\alpha \geq 5 - 2\omega$ . We use the following modified lemma from [15]:



## 79:16 A Faster Interior-Point Method for Sum-Of-Squares Optimization

► **Lemma 18** (Helpful lemma for amortization, modified version of Lemma 10.13 of [15]). *Let  $t$  denote the total number of iterations. Let  $r_i \in [L]$  be the rank for the  $i$ -th iteration for  $i \in [t]$ . Assume  $r_i$  satisfies the following condition: for any vector  $g \in \mathbb{R}_+^L$  which is non-increasing, we have  $\sum_{i=1}^t r_i \cdot g_{r_i} \leq O(t \cdot \|g\|_2)$ .*

*If the cost in the  $i$ -th iteration is  $O(\mathcal{T}_{\text{mat}}(U, U, \min\{Lr_i, U\}))$ , when  $\alpha \geq 5 - 2\omega$ , the amortized cost per iteration is  $U^{2+o(1)} + U^{\omega-1/2+o(1)} \cdot L^{1/2}$ .*

We include the proof of Lemma 18 for completeness. The main difference between this proof and that of [15] is that we cut off at  $U/L$  instead of  $L$ . Our proof makes use of the following two facts about  $\omega$  and  $\alpha$  (Lemma A.4 and Lemma A.5 of [9]).

► **Fact 19** (Relation of  $\omega$  and  $\alpha$ ).  $\omega \leq 3 - \alpha$ .

► **Fact 20** (Upper bound of  $\mathcal{T}_{\text{mat}}(n, n, r)$ ). *For any  $r \leq n$ , we have that  $\mathcal{T}_{\text{mat}}(n, n, r) \leq n^{2+o(1)} + r^{\frac{\omega-2}{1-\alpha}} \cdot n^{2-\frac{\alpha(\omega-2)}{1-\alpha}+o(1)}$ .*

**Proof of Lemma 18.** For  $r_i$  that satisfies  $r_i \leq U/L$ , we have

$$\begin{aligned} \mathcal{T}_{\text{mat}}(U, U, Lr_i) &\leq U^{2+o(1)} + (Lr_i)^{\frac{\omega-2}{1-\alpha}} \cdot U^{2-\frac{\alpha(\omega-2)}{1-\alpha}+o(1)} \\ &= U^{2+o(1)} + U^{2-\frac{\alpha(\omega-2)}{1-\alpha}+o(1)} \cdot L^{\frac{\omega-2}{1-\alpha}} \cdot r_i^{\frac{\omega-2}{1-\alpha}}, \end{aligned} \quad (14)$$

where the first step follows from Fact 20.

Define a sequence  $g \in \mathbb{R}_+^L$  such that for  $r \in [L]$ ,

$$g_r = \begin{cases} r^{\frac{\omega-2}{1-\alpha}-1} & \text{if } r \leq U/L, \\ (U/L)^{\frac{\omega-2}{1-\alpha}} \cdot r^{-1} & \text{if } r > U/L. \end{cases}$$

Note that  $g$  is non-increasing because  $\frac{\omega-2}{1-\alpha} \leq 1$  (Fact 19). Then using the condition in the lemma statement, we have

$$\begin{aligned} \sum_{i=1}^t \min\{r_i^{\frac{\omega-2}{1-\alpha}}, (U/L)^{\frac{\omega-2}{1-\alpha}}\} &= \sum_{i=1}^t r_i \cdot g_{r_i} \\ &\leq t \cdot \|g\|_2 \\ &\leq t \cdot \left( \int_{x=1}^{U/L} x^{\frac{2(\omega-2)}{1-\alpha}-2} dx + (U/L)^{\frac{2(\omega-2)}{1-\alpha}} \cdot \int_{x=U/L}^L x^{-2} dx \right)^{1/2} \quad (15) \\ &\leq t \cdot \left( c \cdot (U/L)^{\frac{2(\omega-2)}{1-\alpha}-1} + (U/L)^{2(\frac{\omega-2}{1-\alpha})} \cdot (U/L)^{-1} \right)^{1/2} \\ &= t \cdot O((U/L)^{\frac{(\omega-2)}{1-\alpha}-1/2}), \end{aligned}$$

where the first step follows from the definition of  $g \in \mathbb{R}_+^L$ , the second step follows from the assumption  $\sum_{i=1}^t r_i \cdot g_{r_i} \leq t \cdot \|g\|_2$  in the lemma statement, the third step follows from upper bounding the  $\ell_2$  norm  $\|g\|_2^2 = \sum_{r=1}^L g_r^2$ , and the fourth step follows  $\frac{2(\omega-2)}{1-\alpha} \geq 1$  when  $\alpha \geq 5 - 2\omega$ , so the integral  $\int_{x=1}^{U/L} x^{\frac{2(\omega-2)}{1-\alpha}-2} dx = c \cdot x^{\frac{2(\omega-2)}{1-\alpha}-1} \Big|_1^{U/L} = O((U/L)^{\frac{2(\omega-2)}{1-\alpha}-1})$  where  $c := 1/(\frac{2(\omega-2)}{1-\alpha} - 1)$ .

Thus we have

$$\begin{aligned}
& \sum_{t=1}^t \mathcal{T}_{\text{mat}}(U, U, \min\{Lr_i, U\}) \\
& \leq \sum_{t=1}^t \left( U^{2+o(1)} + U^{2-\frac{\alpha(\omega-2)}{1-\alpha}+o(1)} \cdot L^{\frac{\omega-2}{1-\alpha}} \cdot \min\{r_i^{\frac{\omega-2}{1-\alpha}}, (U/L)^{\frac{\omega-2}{1-\alpha}}\} \right) \\
& = t \cdot U^{2+o(1)} + U^{2-\frac{\alpha(\omega-2)}{1-\alpha}+o(1)} \cdot L^{\frac{\omega-2}{1-\alpha}} \cdot \sum_{t=1}^t \min\{r_i^{\frac{\omega-2}{1-\alpha}}, (U/L)^{\frac{\omega-2}{1-\alpha}}\} \\
& \leq t \cdot U^{2+o(1)} + U^{2-\frac{\alpha(\omega-2)}{1-\alpha}+o(1)} \cdot L^{\frac{\omega-2}{1-\alpha}} \cdot t \cdot (U/L)^{\frac{(\omega-2)}{1-\alpha}-1/2} \\
& = t \cdot (U^{2+o(1)} + U^{\omega-1/2+o(1)} \cdot L^{1/2}),
\end{aligned}$$

where the first step follows from Eq. (14) and  $\mathcal{T}_{\text{mat}}(U, U, U) = U^\omega = U^{2-\frac{\alpha(\omega-2)}{1-\alpha}} \cdot L^{\frac{\omega-2}{1-\alpha}} \cdot (U/L)^{\frac{\omega-2}{1-\alpha}}$ , the second step follows from moving summation inside, the third step follows from Eq. (15), and the last step follows from adding the terms together. ◀

Now we are ready to prove our main theorem for the amortized time of Algorithm 1.

► **Theorem 21** (Time of Algorithm 1). *When  $\alpha \geq 5 - 2\omega$ , the running time of Algorithm 1 is  $(U^2 \cdot L^{1/2} + U^{\omega-1/2} \cdot L) \cdot (\log(1/\delta) + U^{o(1)})$ .*

**Proof.** Using Lemma 17 the initialization time is  $O(U^\omega) \leq O(U^{\omega-1/2} \cdot L)$  since  $U \leq L^2$ .

Using Lemma 14 we know that the ranks  $r_i$  indeed satisfy the requirement of Lemma 18, and since the worst case time per iteration is  $O(\mathcal{T}_{\text{mat}}(U, U, \min\{Lr_i, U\}))$  (Lemma 17), using Lemma 18 the time per iteration is  $U^{2+o(1)} + U^{\omega-1/2+o(1)} \cdot L^{1/2}$ . Since there are in total  $t = 40\epsilon_N^{-1} \sqrt{L} \log(L/\delta)$  iterations, we get the total running time as claimed. ◀

### 7.3 Comparison with previous results

In this section we compare the running time of [27], [16, 15], and our result. We assume that  $m = \Theta(U)$  when making the comparisons.

Ignoring  $\log(1/\delta)$  and  $U^{o(1)}$  factors, and since  $L \leq U \leq L^2$ , the running times are

$$\begin{aligned}
[27] \text{ (SOS)} & : L^{0.5} \cdot U^\omega, \\
[16, 15] \text{ (SDP)}^5 & : L^{0.5} \cdot \min\{UL^2 + U^\omega, L^4 + L^{2\omega-0.5}\}, \\
\text{Ours (SOS)} & : L^{0.5} \cdot (U^2 + U^{\omega-0.5} \cdot L^{0.5}).
\end{aligned}$$

#### Current $\omega$ and $\alpha$

Plugging in the current best values  $\omega \approx 2.373$  and  $\alpha \approx 0.314$ , we have

$$\begin{aligned}
[27] \text{ (SOS)} & : L^{0.5} \cdot U^{2.373}, \\
[16, 15] \text{ (SDP)} & : L^{0.5} \cdot \min\{UL^2 + U^{2.373}, L^{4.246}\} \\
& = L^{0.5} \cdot \begin{cases} UL^2 & \text{when } U \in (L, L^{1.457}], \\ U^{2.373} & \text{when } U \in (L^{1.457}, L^{1.789}], \\ L^{4.246} & \text{when } U \in (L^{1.789}, L^2), \end{cases} \\
\text{Ours (SOS)} & : L^{0.5} \cdot (U^2 + U^{1.873} L^{0.5}).
\end{aligned}$$

<sup>5</sup> When solving SOS, [16] has running time  $O(L^{0.5} \cdot (UL^2 + U^\omega + L^\omega)) \leq O(L^{0.5} \cdot (UL^2 + U^\omega))$ , and [15] has running time  $O(L^{0.5} \cdot (U^2 + L^4) + U^\omega + L^{2\omega}) \leq O(L^{4.5} + L^{2\omega})$  since  $L \leq U \leq L^2$ .

Note that our running time is always better than the previous results, and for several values of  $L$  and  $U$  we improve by a polynomial factor. See Figure 1 for an illustration.

## 8 Initialization

There exist standard techniques to transform a convex program to a form that has an easily obtainable strictly feasible point, see e.g. [39]. We follow the initialization procedure presented by [9] and [16] and adapt to SOS optimization. Similar initialization lemma exists for WSOS optimization. The proof of this lemma can be found in our full version.

Let the matrix  $P \in \mathbb{R}^{U \times L}$  and the operator  $\Lambda : \mathbb{R}^U \rightarrow \mathbb{R}^{L \times L}$  that  $\Lambda(s) = P^\top \text{diag}(s)P$  be defined as in the interpolant basis paragraph of Section 3.

► **Lemma 22** (Initialization). *Given an instance of (SOS) that fulfills Slater's condition, and let  $R$  be an upper bound on the  $\ell_1$ -norm of the primal feasible solutions, i.e. all primal feasible  $x$  of (SOS) fulfill  $\|x\|_1 \leq R$ , and let  $\delta \in (0, 1)$ . We define  $\bar{A} \in \mathbb{R}^{(m+1) \times (U+2)}$ ,  $\bar{b} \in \mathbb{R}^{m+1}$ , and  $\bar{c} \in \mathbb{R}^{U+2}$  as*

$$\bar{A} = \begin{bmatrix} A & 0 & \frac{1}{R}b - A\bar{g}^0 \\ \mathbf{1}_U^\top & 1 & 0 \end{bmatrix}, \quad \bar{b} = \begin{bmatrix} \frac{1}{R}b \\ 1 + \langle \mathbf{1}_U, \bar{g}^0 \rangle \end{bmatrix}, \quad \bar{c} = \begin{bmatrix} \frac{\delta}{\|c\|_\infty}c \\ 0 \\ 1 \end{bmatrix},$$

and let

$$\bar{x}^0 = \begin{bmatrix} \bar{g}^0 \\ 1 \\ 1 \end{bmatrix} \in \mathbb{R}^{U+2}, \quad \bar{y}^0 = \begin{bmatrix} 0_m \\ -1 \end{bmatrix} \in \mathbb{R}^{m+1}, \quad \text{and } \bar{s}^0 = \begin{bmatrix} \mathbf{1}_U + \frac{\delta}{\|c\|_\infty}c \\ 1 \\ 1 \end{bmatrix} \in \mathbb{R}^{U+2},$$

where  $\bar{g}^0 = g_{\Sigma^*}(\bar{s}_{[U]}^0) \in \mathbb{R}^U$  for the gradient function  $g_{\Sigma^*}(s) := \text{diag}(P(P^\top \text{diag}(s)P)^{-1}P^\top)$  that maps from  $\mathbb{R}^U$  to  $\mathbb{R}^U$ . This defines the auxiliary primal-dual system

$$\begin{aligned} \min \langle \bar{c}, \bar{x} \rangle & \qquad \qquad \qquad \max \langle \bar{y}, \bar{b} \rangle \\ \bar{A}\bar{x} = \bar{b} & \qquad \qquad \qquad \bar{A}^\top \bar{y} + \bar{s} = \bar{c} \\ \bar{x} \in \Sigma_{n,2d} \times \mathbb{R}_{\geq 0}^2, & \qquad \qquad \bar{s} \in \Sigma_{n,2d}^* \times \mathbb{R}_{\geq 0}^2. \end{aligned} \tag{Aux-SOS}$$

Then  $(\bar{x}^0, \bar{y}^0, \bar{s}^0)$  are feasible to the auxiliary system (Aux-SOS).

Further, under the canonical barrier (we use  $\bar{a}_i$  to denote the  $i$ -th column of  $\bar{A}$ ):

$$\bar{F}_\eta(\bar{y}) = -\eta \langle \bar{y}, \bar{b} \rangle - \log \det \left( \Lambda((\bar{c} - \bar{A}^\top \bar{y})_{[U]}) \right) - \log(\bar{c}_{U+1} - \langle \bar{a}_{U+1}, \bar{y} \rangle) - \log(\bar{c}_{U+2} - \langle \bar{a}_{U+2}, \bar{y} \rangle),$$

we have that  $\|\bar{g}_{\eta^0}(\bar{y}^0)\|_{\bar{H}(\bar{y}^0)^{-1}} = 0$  for  $\eta^0 = 1$ .

Further, for any solution  $(\bar{x}, \bar{y}, \bar{s})$  to (Aux-SOS) with duality gap  $\leq \delta^2$ , its restriction  $\hat{x} := \bar{x}_{[U]}$  fulfills

$$\begin{aligned} \langle c, \hat{x} \rangle & \leq \min_{Ax=b, x \in \Sigma_{n,2d}} \langle c, x \rangle + \delta \cdot R \|c\|_\infty, \\ \|A\hat{x} - b\|_1 & \leq 8\delta L \cdot (LR\|A\|_\infty + \|b\|_1), \\ \hat{x} & \in \Sigma_{n,2d}. \end{aligned}$$

---

## References

- 1 Josh Alman and Virginia Vassilevska Williams. A refined laser method and faster matrix multiplication. In *Proceedings of the 2021 ACM-SIAM Symposium on Discrete Algorithms (SODA)*, pages 522–539. SIAM, 2021.

- 2 Christine Bachoc and Frank Vallentin. New upper bounds for kissing numbers from semidefinite programming. Technical report, Journal of the American Mathematical Society, 2006.
- 3 Brandon Ballinger, Grigoriy Blekherman, Henry Cohn, Noah Giansiracusa, Elizabeth Kelly, and Achill Schürmann. Experimental study of energy-minimizing point configurations on spheres. *Experimental Mathematics*, 18(3):257–283, 2009. doi:10.1080/10586458.2009.10129052.
- 4 Boaz Barak, Samuel B. Hopkins, Jonathan A. Kelner, Pravesh K. Kothari, Ankur Moitra, and Aaron Potechin. A nearly tight sum-of-squares lower bound for the planted clique problem. *SIAM J. Comput.*, 48(2):687–735, 2019. doi:10.1137/17M1138236.
- 5 Boaz Barak, Prasad Raghavendra, and David Steurer. Rounding semidefinite programming hierarchies via global correlation. In *2011 IEEE 52nd annual symposium on foundations of computer science (FOCS)*, pages 472–481. IEEE, 2011.
- 6 Markus Bläser. Fast matrix multiplication. *Theory of Computing*, pages 1–60, 2013.
- 7 Grigoriy Blekherman, Pablo A Parrilo, and Rekha R Thomas. *Semidefinite optimization and convex algebraic geometry*. SIAM, 2012.
- 8 L. Bos, S. De Marchi, A. Sommariva, and M. Vianello. Computing multivariate fekeete and leja points by numerical linear algebra. *SIAM Journal on Numerical Analysis*, 48(5):1984–1999, 2010. doi:10.1137/090779024.
- 9 Michael B Cohen, Yin Tat Lee, and Zhao Song. Solving linear programs in the current matrix multiplication time. In *Proceedings of the 51st Annual ACM Symposium on Theory of Computing (STOC)*, 2019.
- 10 François Le Gall and Florent Urrutia. Improved rectangular matrix multiplication using powers of the coppersmith-winograd tensor. In *Proceedings of the 2018 ACM-SIAM Symposium on Discrete Algorithms (SODA)*, pages 1029–1046. SIAM, 2018.
- 11 Bissan Ghaddar, Jakub Marecek, and M. Mevissen. Optimal power flow as a polynomial optimization problem. *IEEE Transactions on Power Systems*, 31:539–546, 2016.
- 12 Roxana Heß, Didier Henrion, Jean-Bernard Lasserre, and Tien Son Pham. Semidefinite approximations of the polynomial abscissa. *SIAM J. Control. Optim.*, 54(3):1633–1656, 2016. doi:10.1137/15M1033198.
- 13 Samuel B. Hopkins, Pravesh K. Kothari, Aaron Potechin, Prasad Raghavendra, Tselil Schramm, and David Steurer. The power of sum-of-squares for detecting hidden structures. In *58th IEEE Annual Symposium on Foundations of Computer Science, (FOCS)*, pages 720–731. IEEE Computer Society, 2017. doi:10.1109/FOCS.2017.72.
- 14 Samuel B. Hopkins and Jerry Li. Mixture models, robustness, and sum of squares proofs. In *Proceedings of the 50th Annual ACM SIGACT Symposium on Theory of Computing, (STOC)*, pages 1021–1034. ACM, 2018. doi:10.1145/3188745.3188748.
- 15 Baihe Huang, Shunhua Jiang, Zhao Song, Runzhou Tao, and Ruizhe Zhang. Solving sdp faster: A robust ipm framework and efficient implementation, 2021. arXiv:2101.08208.
- 16 Haotian Jiang, Tarun Kathuria, Yin Tat Lee, Swati Padmanabhan, and Zhao Song. A faster interior point method for semidefinite programming. In *2020 IEEE 61st annual symposium on foundations of computer science (FOCS)*, pages 910–918. IEEE, 2020.
- 17 Shunhua Jiang, Yunze Man, Zhao Song, Zheng Yu, and Danyang Zhuo. Fast graph neural tangent kernel via kronecker sketching. *arXiv preprint AAAI’22*, 2021. arXiv:2112.02446.
- 18 Narendra Karmarkar. A new polynomial-time algorithm for linear programming. In *Proceedings of the 16th annual ACM symposium on Theory of computing (STOC)*, pages 302–311, 1984.
- 19 Jean Bernard Lasserre. *An Introduction to Polynomial and Semi-Algebraic Optimization*. Cambridge Texts in Applied Mathematics. Cambridge University Press, 2015. doi:10.1017/CB09781107447226.
- 20 M. Laurent. *Sums of squares, moment matrices and optimization over polynomials*, pages 155–270. Number 149 in The IMA Volumes in Mathematics and its Applications Series. Springer Verlag, Germany, 2009.
- 21 François Le Gall. Powers of tensors and fast matrix multiplication. In *Proceedings of the 39th international symposium on symbolic and algebraic computation*, pages 296–303, 2014.

- 22 Yin Tat Lee and Aaron Sidford. Path finding methods for linear programming: Solving linear programs in  $\tilde{O}(\sqrt{\text{rank}})$  iterations and faster algorithms for maximum flow. In *2014 IEEE 55th Annual Symposium on Foundations of Computer Science (FOCS)*, pages 424–433. IEEE, 2014.
- 23 Yin Tat Lee, Zhao Song, and Qiuyu Zhang. Solving empirical risk minimization in the current matrix multiplication time. In *Conference on Learning Theory (COLT)*, pages 2140–2157. PMLR, 2019.
- 24 Yurii Nesterov. Squared functional systems and optimization problems. In *High performance optimization*, pages 405–440. Springer, 2000.
- 25 Yurii Nesterov and Arkadi Nemirovski. Interior-point polynomial algorithms in convex programming. In *Siam Studies in Applied Mathematics*, 1987. doi:10.1137/1.9781611970791.
- 26 Dávid Papp. Optimal designs for rational function regression. *Journal of the American Statistical Association*, 107(497):400–411, 2012. doi:10.1080/01621459.2012.656035.
- 27 Dávid Papp and Sercan Yildiz. Sum-of-squares optimization without semidefinite programming. *SIAM Journal on Optimization*, 29(1):822–851, 2019.
- 28 Pablo Parrilo. *Sum of squares : theory and applications : AMS short course, sum of squares : theory and applications, January 14-15, 2019, Baltimore, Maryland*. American Mathematical Society, Providence, Rhode Island, 2020.
- 29 Mihai Putinar and Florian-Horia Vasilescu. Positive polynomials on semi-algebraic sets. *Comptes Rendus de l'Académie des Sciences - Series I - Mathematics*, 328(7):585–589, 1999. doi:10.1016/S0764-4442(99)80251-1.
- 30 James Renegar. *A Mathematical View of Interior-Point Methods in Convex Optimization*. Society for Industrial and Applied Mathematics, January 2001. doi:10.1137/1.9780898718812.
- 31 Tae Roh, Bogdan Dumitrescu, and Lieven Vandenbergh. Multidimensional FIR filter design via trigonometric sum-of-squares optimization. *J. Sel. Topics Signal Processing*, 1(4):641–650, 2007. doi:10.1109/JSTSP.2007.910261.
- 32 Alvis Sommariva and Marco Vianello. Computing approximate fekete points by qr factorizations of vandermonde matrices. *Computers & Mathematics with Applications*, 57(8):1324–1336, 2009.
- 33 Zhao Song, Shuo Yang, and Ruizhe Zhang. Does preprocessing help training over-parameterized neural networks? *Advances in Neural Information Processing Systems*, 34, 2021.
- 34 Zhao Song, Lichen Zhang, and Ruizhe Zhang. Training multi-layer over-parametrized neural network in subquadratic time. *arXiv preprint*, 2021. arXiv:2112.07628.
- 35 Gilbert Strang. Karmarkar’s algorithm and its place in applied mathematics. *The Mathematical Intelligencer*, 9(2):4–10, 1987.
- 36 Ning Tan. *On the Power of Lasserre SDP Hierarchy*. PhD thesis, EECS Department, University of California, Berkeley, December 2015. URL: <http://www2.eecs.berkeley.edu/Pubs/TechRpts/2015/EECS-2015-236.html>.
- 37 Pravin M Vaidya. Speeding-up linear programming using fast matrix multiplication. In *30th Annual Symposium on Foundations of Computer Science (FOCS)*, pages 332–337. IEEE, 1989.
- 38 Jan van den Brand, Binghui Peng, Zhao Song, and Omri Weinstein. Training (overparametrized) neural networks in near-linear time. In *12th Innovations in Theoretical Computer Science Conference (ITCS 2021)*, volume 185, pages 63:1–63:15, 2021. doi:10.4230/LIPIcs.ITCS.2021.63.
- 39 Yinyu Ye, Michael J Todd, and Shinji Mizuno. An  $O(\sqrt{n}L)$ -iteration homogeneous and self-dual linear programming algorithm. *Mathematics of operations research*, 19(1):53–67, 1994.