

Rounding Sum of Squares Relaxations

Boaz Barak – Microsoft Research

Joint work with Jonathan Kelner (MIT) and David Steurer (Cornell)

SSS



ICERM workshop on semidefinite programming and graph algorithms

February 10-14, 2014

This talk is about

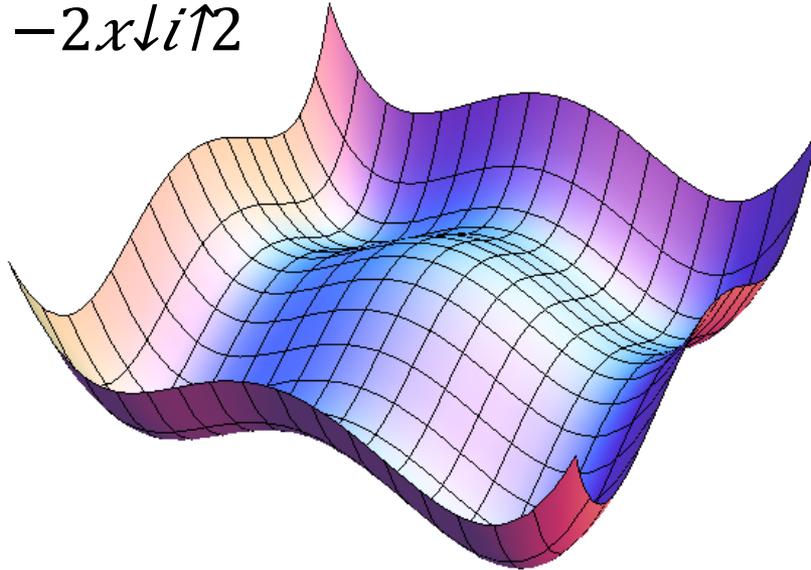
- Semi-definite programming , SOS/Positivstellensatz method
- Proof complexity
- The Unique Games Conjecture
- Graph partitioning, small set expansion
- Machine Learning
- *Cryptography.. (in spirit).*

Sum-of-Squares (SOS) Algorithm

[Shor'87, Parillo '00, Nesterov '00, Lasserre '01]

Motivation: Sometimes a polynomial P can have exponentially many local minima...

E.g. $P(\mathbf{x}) = n + \sum_{i=1}^n (x_i^4 - 2x_i^2)$



... but there is still a short proof that $P \geq 0$

E.g. $P(\mathbf{x}) = \sum_{i=1}^n (x_i^2 - 1)^2$

... and this proof can be found efficiently via semidefinite programming (SDP)

SOS Algorithm:

For low degree $P, P \downarrow 1, \dots, P \downarrow k$ we consider the program \mathcal{P} :

$$\max_{x \in \mathbb{R}^n} P(x) \quad \text{s.t.}$$

$$P \downarrow 1(x) = \dots = P \downarrow k(x) = 0$$

SOS Proof that $\mathcal{P} < \nu$: Polynomials $Q \downarrow 1, \dots, Q \downarrow k$ and SOS polys S, S' s.t.

$$(\nu - P)S = \sum P \downarrow i Q \downarrow i + S' + 1$$

Positivstellensatz: All true bounds have SOS proof. [Artin '27, Krivine '64, Stengle '74]

Degree of proof: max degree of $Q \downarrow 1, \dots, Q \downarrow k, S, S'$ [Gregoriev-Vorobjov'99]

Theorem: [Shor '87, Parillo '00, Nesterov '00, Lasserre '01]

1) A proof of degree d can be found in $n \uparrow O(d)$ time.

2) Can find in $n \uparrow O(d)$ time the min ν s.t. \exists degree d proof that $\mathcal{P} < \nu$

Program \mathcal{P} :

SOS Algorithm:

For low degree $P, P \downarrow 1, \dots, P \downarrow k$ we consider the program \mathcal{P} :

$$\max_{x \in \mathbb{R}^n} P(x) \quad \text{s.t.}$$

Can optimize in $n \uparrow O(d)$ time over programs with degree d proofs.

SOS Proof that $P < \nu$: Polynomials $Q \downarrow 1, \dots, Q \downarrow k$ and SOS polys S, S' s.t.

$$(\nu - P)S = \sum_{i=1}^k Q \downarrow i + S' + 1$$

Positivstellensatz: All true bounds have SOS proof. [Artin '27, Krivine '64, Stengle '74]

Degree of proof: max degree of $Q \downarrow 1, \dots, Q \downarrow k, S, S'$ [Gregoriev-Vorobjov'99]

Theorem: [Shor '87, Parillo '00, Nesterov '00, Lasserre '01]

1) A proof of degree d can be found in $n \uparrow O(d)$ time.

2) Can find in $n \uparrow O(d)$ time the min ν s.t. \exists degree d proof that $P < \nu$

Program \mathcal{P} :

$$\max_{x \in \mathbb{R}^n} P(x) \quad s.t.$$
$$P \downarrow 1(x) = \dots = P \downarrow k(x) = 0$$

SOS Proof that $\mathcal{P} \Rightarrow P$ $S = \sum P \downarrow i Q \downarrow i + S^T + 1$

Can optimize in $n \uparrow O(d)$ time over programs with degree d proofs.

Can't hope for $d \ll n$ always: Captures SAT, CLIQUE, 3COL, MAX-CUT, etc...

But maybe $d \ll n$ often? Essentially only one (robust) lower bound showing $d \geq \Omega(n)$
[Grigoriev '01]

This talk: General method to analyze the SOS algorithm. [B-Kelner-Steurer'13]

- Applications:**
- Optimizing polynomials w/ non-negative coefficients over sphere.
 - Algorithms for **quantum separability problem** [Brandao-Harrow'13]
 - **Sparse coding:** learning dictionaries beyond the \sqrt{n} barrier.
 - Finding **sparse vectors in subspaces.**
 - Approach to refute the **Unique Games Conjecture.**

Program \mathcal{P} :

$$\max_{x \in \mathbb{R}^n} P(x) \quad s.t.$$
$$P \downarrow 1(x) = \dots = P \downarrow k(x) = 0$$

SOS Proof that $\mathcal{P} \Rightarrow P$ $S = \sum P \downarrow i Q \downarrow i + S^T + 1$

Can optimize in $n \uparrow O(d)$ time over programs with degree d proofs.

Can't hope for $d \ll n$ always: Captures SAT, CLIQUE, 3COL, MAX-CUT, etc...

But maybe $d \ll n$ often? Essentially only one (robust) lower bound showing $d \geq \Omega(n)$
[Grigoriev '01]

This talk: General method to analyze the SOS algorithm. [B-Kelner-Steurer'13]

- Applications:**
- Optimizing polynomials w/ non-negative coefficients over sphere.
 - Algorithms for **quantum separability problem** [Brandao-Harrow'13]
 - **Sparse coding:** learning dictionaries beyond the \sqrt{n} barrier.
 - Finding **sparse vectors in subspaces.**
 - Approach to refute the **Unique Games Conjecture.**

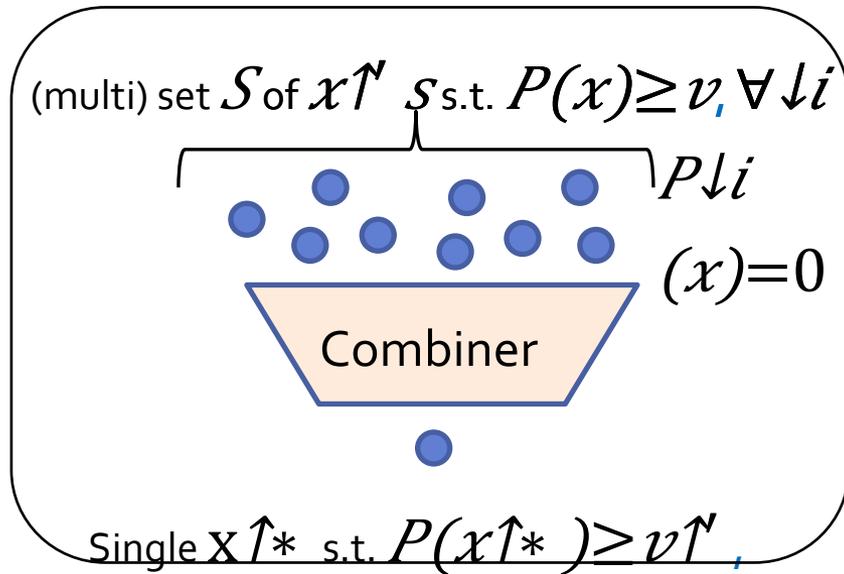
Program \mathcal{P} :

$$\max_{x \in \mathbb{R}^n} P(x) \quad \text{s.t.}$$

$$P \downarrow 1(x) = \dots = P \downarrow k(x) = 0$$

Finding x is hard. We consider easier problem:

Given many x 's maximizing \mathcal{P} , find a single x^* with value close to maximum.



"Finding a needle in a needle-stack"

Non-trivial combiner:

Only depends on low degree marginals of \mathcal{S}

$$\{ \mathbb{E} \downarrow x \sim \mathcal{S} \ x \downarrow i \downarrow 1 \ \dots \ x \downarrow i \downarrow k \} \downarrow$$

$$i \downarrow 1, \dots, i \downarrow k$$

[B-Kelner-Steurer'13]: Transform "simple" non-trivial combiners to algorithm for original problem.

Idea in a nutshell: Simple combiners will output a solution even when fed "fake marginals".

Pseudoexpectations (aka "Fake Marginals")

- Def:** [Lasserre '01] *Degree d pseudoexpectation* is operator mapping any degree $\leq d$ poly P into a number $\mathbb{E} P(X)$ satisfying:
- **Normalization:** $\mathbb{E} 1 = 1$
 - **Linearity:** $\mathbb{E} [aP(X) + bQ(X)] = a\mathbb{E} P(X) + b\mathbb{E} Q(X) \quad \forall P, Q$ of $\deg \leq d$
 - **Positivity:** $\mathbb{E} P^2(X) \geq 0 \quad \forall P$ of $\deg \leq d/2$

Fundamental Fact: \exists deg d SOS proof for $P > 0 \Leftrightarrow$

$\mathbb{E} P(X) > 0$ for any deg $O(d)$ pseudoexpectation operator

Take home message:

- Pseudoexpectation "looks like" real expectation to low degree polynomials.
- Can efficiently find pseudoexpectation matching any polynomial constraints.
- Proofs about real random vars can often be "lifted" to pseudoexpectation.



Idea in a nutshell: Simple combiners will output a solution even when fed "fake marginals".

Pseudoexpectations (aka "Fake Marginals")

Problem: Given low degree $P_1, \dots, P_k: \mathbb{R}^n \rightarrow \mathbb{R}$, maximize $P(x)$ s.t. $\forall i$

Def: Degree d pseudoexpectation is operator mapping any degree $\leq d$ poly P into a number $\mathbb{E}(P(X))$ satisfying:

- Normalization: $\mathbb{E} 1 = 1$
- Linearity: $\mathbb{E} [aP(X) + bQ(X)] = a\mathbb{E} P(X) + b\mathbb{E} Q(X) \quad \forall P, Q$ of $\deg \leq d$
- Positivity: $\mathbb{E} P^2(X) \geq 0 \quad \forall P$ of $\deg \leq d/2$

Fundamental Fact: \exists deg d SOS proof for $P > 0 \iff \mathbb{E} P(X) > 0$ for any deg $O(d)$ pseudoexpectation operator

Take home message:

- Pseudoexpectation "looks like" real expectation to low degree polynomials.
- Can efficiently find pseudoexpectation matching any polynomial constraints.
- Proofs about real random vars can often be "lifted" to pseudoexpectation.



Idea in a nutshell: Simple combiners will output a solution even when fed "fake marginals".

Combining \Rightarrow Rounding

Problem: Given low degree $P, P_{\downarrow 1}, \dots, P_{\downarrow k} : \mathbb{R}^n \rightarrow \mathbb{R}$ maximize $P(x)$ s.t. $\forall i$

[B-Kelner-Steurer'13]: Transform "simple" non-trivial combiners to algorithm for original problem.

Non-trivial combiner: Alg \mathcal{C} with

INPUT: $\{ \mathbb{E} X_{\downarrow i \downarrow 1} \dots X_{\downarrow i \downarrow k} \}_{i \downarrow 1 \dots i \downarrow k \in [n]}$, X r.v. over \mathbb{R}^n s.t. \mathbb{E}

OUTPUT: $x^* \in \mathbb{R}^n$ s.t. $\mathbb{E} P(x(X)) \geq v \neq 0, \forall i$

$P_{\downarrow i}(x^*) = 0$

Crucial Observation: If proof that x^* is good solution is in SOS framework, then it holds even if \mathcal{C} is fed with a **pseudoexpectation**.

Corollary: In this case, we can find x^* efficiently:

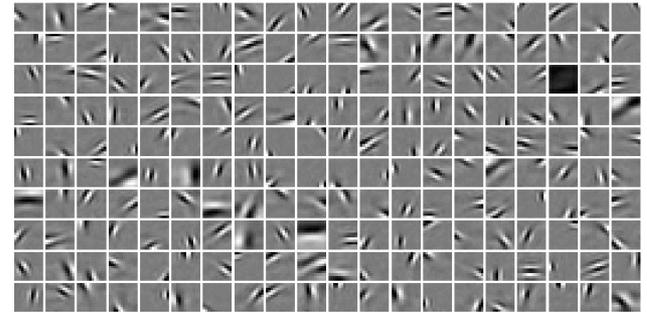
- Use SOS PSD to find pseudoexpectation matching input conditions.
- Use \mathcal{C} to **round** the PSD solution into an actual solution x^*

Example Application: Dictionary Learning / Sparse Coding

Let $a_1, \dots, a_m \in \mathbb{R}^n$ set of vectors.

Goal: Given examples of form $y = \sum W_{\downarrow i} a_i$, where $\mu = \Pr[W_{\downarrow i} \neq 0] \ll 1$ recover

Find the "right" representation of observed data $a_1 \dots a_m$ [Olhausen-Field '96]



LOTS of work: important primitive in Machine Learning, Vision, Neuroscience...

Previous best (rigorous) results: $\mu \ll 1/\sqrt{n}$

[Spielman-Wang-Wright '12, Arora-Moitra-Ge '13, Agrawal-Anandkumar-Jain-Netrapalli-Tandon '13]

We show: $\mu \ll 1$ is sufficient* (even in non-independent, overcomplete case)

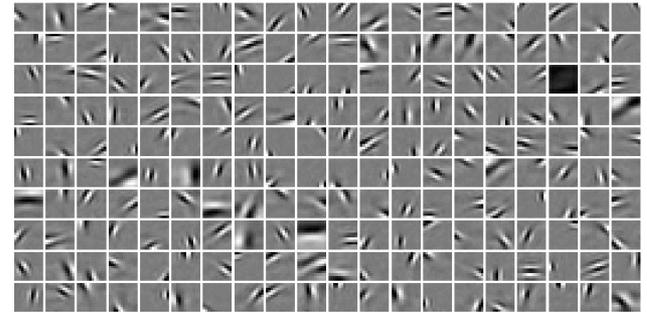
*In succinctly time, we show $\mu \ll n^{-1/2} = \epsilon$ is sufficient in poly

Example Application: Dictionary Learning / Sparse Coding

Let $a_1, \dots, a_m \in \mathbb{R}^n$ set of vectors.

Goal: Given examples of form $y = \sum W_{\downarrow i} a_i$, where $\mu = \Pr[W_{\downarrow i} \neq 0] \ll 1$ recover

Find the "right" representation of observed data $a_1 \dots a_m$ [Olhausen-Field '96]



LOTS of work: important primitive in Machine Learning, Vision, Neuroscience,...

Previous best (rigorous) results: $\mu \ll 1/\sqrt{n}$

[Spielman-Wang-Wright '12, Arora-Moitra-Ge '13, Agrawal-Anandkumar-Jain-Netrapalli-Tandon '13]

We show: $\mu \ll 1$ is sufficient* (even in non-independent, overcomplete case)

*In succinctly time, we show $\mu \ll n^{-\epsilon}$ is sufficient in poly

Let $a_1, \dots, a_m \in \mathbb{R}^n$ set of vectors.

Goal: Given examples of form $y = \sum W_i a_i$, where $\mu = \Pr[W_i \neq 0] \ll 1$ recover a_1, \dots, a_m

Achieve in 3 steps:

(1) Find a program \mathcal{P} s.t. every x maximizing \mathcal{P} is close to one of a_i 's

(2) Give combining alg \mathcal{C} taking moments of dist X over maximizers into a vector x^* close to one of a_i 's.

(3) Show that arguments in (1) and (2) fall under the

Result generalizes to overcomplete, **non independent** case.

Step 1. For simplicity, assume $m=n$, a_i 's orthonormal basis, W_i i.i.d. random vars

over $\{0, \pm 1\}$ s.t. $\Pr[W_i = +1] = \Pr[W_i = -1] = \mu/2$

Consider the polynomial $P(x) = \mathbb{E} \langle y, x \rangle^4 = \mathbb{E} (\sum W_i \langle a_i, x \rangle)^4$ (can approximate P arbitrarily well from examples)

Opening parenthesis we get $P(x) \leq \mu \sum \langle a_i, x \rangle^4 + 2\mu^2 (\sum \langle a_i, x \rangle^2)^2 = \mu \sum \langle a_i, x \rangle^4 + 2\mu^2 \|x\|^4$

Corollary: x unit, $P(x) \geq \mu - o(\mu) \Rightarrow \max \langle a_i, x \rangle^2 \geq 1 - o(1)$

Establishes (1) !

Let $a_1, \dots, a_m \in \mathbb{R}^n$ set of vectors.

Goal: Given examples of form $y = \sum W_i a_i$, where $\mu = \Pr[W_i \neq 0] \ll 1$ recover a_1, \dots, a_m

Achieve in 3 steps: $a_m \downarrow$

(1) Find a program \mathcal{P} s.t. every x maximizing \mathcal{P} is close to one of a_i 's

(2) Give combining alg \mathcal{C} taking moments of dist X over maximizers into a vector x^* close to one of a_i 's.

(3) Show that arguments in (1) and (2) fall under the

Result generalizes to overcomplete, **non independent** case.

Step 1. For simplicity, assume $m=n$, a_i 's orthonormal basis, W_i i.i.d. random vars

over $\{0, \pm 1\}$ s.t. $\Pr[W_i = +1] = \Pr[W_i = -1] = \mu/2$

Consider the polynomial $P(x) = \mathbb{E} \langle y, x \rangle^4 = \mathbb{E} (\sum W_i \langle a_i, x \rangle)^4$ (can approximate P arbitrarily well from examples)

Opening parenthesis we get $P(x) \leq \mu \sum \langle a_i, x \rangle^4 + 2\mu^2 (\sum \langle a_i, x \rangle^2)^2 = \mu \sum \langle a_i, x \rangle^4 + 2\mu^2 \|x\|^4$

Corollary: x unit, $P(x) \geq \mu - o(\mu) \Rightarrow \max \langle a_i, x \rangle^2 \geq 1 - o(1)$

Establishes (1) !

Let $a_1, \dots, a_m \in \mathbb{R}^n$ set of vectors.

Goal: Given examples of form $y = \sum W_i a_i$, where $\mu = \Pr[W_i \neq 0] \ll 1$ recover a_i

Achieve in 3 steps:

a_m

(1) Find a program \mathcal{P} s.t. every x maximizing \mathcal{P} is close to one of a_i 's

(2) Give combining alg \mathcal{C} taking mon into a vector x close to one of a_i 's.

Slightly tedious but straightforward computations.

(3) Show that arguments in (1) and (2) fall under the SOS framework.

Step 2. Let X be dist over unit vectors s.t. every $x \in \text{Supp}(X)$ satisfies $\langle a_i, x \rangle \geq 1 - o(1)$ for some i

Pick set $U = \{u_1 \dots u_\ell\}$ of random (std gaussian) vectors.

Let M be matrix s.t. $M_{i,j} = \mathbb{E} f_{i,j}(U)$ for $f_{i,j}(x) = \prod_{u \in U} \langle u, x \rangle$. Note that $\mathbb{E} f_{i,j}(x) = \prod_{u \in U} \mathbb{E} \langle u, x \rangle = 0$

Our combining algorithm outputs the top e-vec of M

Happens w $\exp(-O(\ell))$ prob

Suppose that $\Pr[X \text{ close to } a_1] \geq 1/n$ and for every $t \in [\ell]$, $\langle u_t, a_1 \rangle \geq 1/n$

Then if $\ell \gg \log n$ then (up to scaling) $M \approx (a_1 a_1^T) \otimes 2$ and we'll

Establishes (2) !

A personal overview of the Unique Games Conjecture

Unique Games Conjecture: UG/SSE problem is NP-hard. [Khot'02,Raghavendra-Steurer'08]

reasons to believe

~~"Standard crypto heuristic":
Tried to solve it and couldn't.~~

~~Very clean picture of complexity landscape:
simple algorithms are optimal
[Khot'02...Raghavendra'08....]~~

~~Simple poly algorithms can't refute it
[Khot-Vishnoi'04]~~

~~Simple subexp' algorithms can't refute it
[B-Gopalan-Håstad-Meka-Raghavendra-Steurer'12]~~

SOS proof system

reasons to suspect

~~Random instances are easy via simple
algorithm
[Arora-Khot-Kolla-Steurer-Tulsiani-Vishnoi'05]~~

~~Quasipoly algo on KV instance
[Kolla '10]~~

~~Subexponential algorithm
[Arora-B-Steurer '10]~~

SOS solves all candidate hard instances
[B-Brandao-Harrow-Kelner-Steurer-Zhou '12]

SOS useful for sparse vector problem
Candidate algorithm for search problem
[B-Kelner-Steurer '13]

Conclusions

- Sum of Squares is a powerful algorithmic framework that can yield strong results for the **right** problems.

(contrast with previous results on SDP/LP hierarchies, showing lower bounds when using either wrong hierarchy or wrong problem.)

- “Combiner” view allows to focus on the features of the problem rather than details of relaxation.
- SOS seems particularly useful for problems with some geometric structure, includes several problems related to unique games and machine learning.
- Still have only rudimentary understanding when SOS works or not.
- Other proof complexity \leftrightarrow approximation algorithms connections?

Other Results

Sparse vector problem:

Recover μ -sparse vector in d -dimensional subspace given arbitrary basis.

(motivation: machine learning, optimization, [Demagnet-Hand 13])

worst-case variant is algorithmic bottleneck in UG/SSE alg [Arora-B-Steurer'10])

Random case: Recovery for any $\mu \ll 1$

(Improving on $\mu \ll 1/\sqrt{d}$ [Demagnet-Hand '13])

Worst case: Recovery* for $\mu \ll 1/d^{1/3}$

[Brandao-Harrow'12]: Using our techniques, find separable quantum state maximizing a "local operations classical communication" (*LOCC*) measurement.