Zeyuan Allen-Zhu Ankit Garg Yuanzhi Li Rafael Oliveira Avi Wigderson

Geodesically Convex Optimization & Applications to Operator Scaling and Invariant Theory



# Contents

- 2nd order methods for Matrix Scaling
- Geodesic Convexity
- Operator Scaling Setup & Algorithm
- Application: Orbit Closure Intersection

# **Recap - Non-Negative Matrices & Scaling**

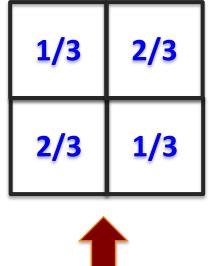
- $X \in M_n(\mathbb{R}_{\geq 0})$  is **doubly stochastic (DS)** if row/column sums of X are equal to 1.
- Y is **scaling** of X if  $\exists$  positive  $\alpha_1, \dots, \alpha_n, \beta_1, \dots, \beta_n$ s.t.  $y_{ij} = \alpha_i x_{ij} \beta_j$ .
- X has DS scaling if  $\exists$  scaling Y of X s.t. all row/column sums of Y equal 1.

$$ds(A) = \sum_{i} (r_{i} - 1)^{2} + \sum_{j} (c_{j} - 1)^{2}$$

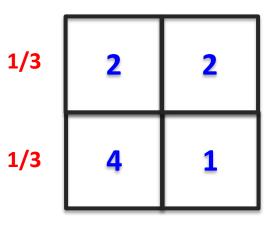
A has approx. DS scaling if  $\forall \epsilon > 0$  there is scaling  $B_{\epsilon}$  of A s.t.  $ds(B_{\epsilon}) < \epsilon$ .

- 1. When does *X* have approx. DS scaling?
- 2. Can we find it efficiently?

### Has **convex** formulation!







 $\mathbf{A} \in M_n(\mathbb{R}_{\geq 0})$  input matrix.

$$f(x) = \sum_{1 \le i \le n} log\left(\sum_{j} A_{ij} e^{x_j}\right) - \sum_{j} x_j$$

Side Note: f(x) is logarithm of [GY'98] capacity for matrix scaling

A has DS scaling iff

$$inf\{f(x): x > 0\} > -\infty$$

How can we solve (really fast) optimization problem above?

- $\nabla^2 f(x)$  not bounded spectral norm bad for 1<sup>st</sup> order methods
- f(x) not self-concordant cannot apply std 2<sup>nd</sup> order methods
- But f(x) "self-robust" still hope for some 2<sup>nd</sup> order methods

# Self Concordance & Self Robustness

Self concordance:  $f : \mathbb{R} \to \mathbb{R}$  is self concordant if  $|f'''(x)| \le 2(f''(x))^{3/2}$ 

 $f : \mathbb{R}^n \to \mathbb{R}$  self concordant if self concordant along each line. "well-approximated" by quadratic function around every pt.

Unfortunately, log of capacity **NOT** self-concordant.

Self robustness [CMTV'18, ALOW'18]:  $f : \mathbb{R} \to \mathbb{R}$  is self robust if  $|f'''(x)| \le 2 \cdot f''(x)$ 

 $f: \mathbb{R}^n \to \mathbb{R}$  self robust if self robust along each line.

"well approximated" by quadratic on <u>small nbhd</u> around each pt.

Log of capacity is self-robust!

**Question:** Can we efficiently optimize self-robust functions?

**Answer:** Yes! Perform "box-constrained Newton Method"

Essentially: optimize "quadratic approx" of fncn on small nbhd

Self robustness [CMTV'18, ALOW'18]:  $f : \mathbb{R} \to \mathbb{R}$  is self robust if  $|f'''(x)| \le 2 \cdot f''(x)$ 

 $f : \mathbb{R}^n \to \mathbb{R}$  self robust if self robust along each line. "well approximated" by quadratic on small nbhd around each pt.

More formally:  $f : \mathbb{R}^n \to \mathbb{R}$  self robust,  $x, \delta \in \mathbb{R}^n$  s.t.  $||\delta||_{\infty} \leq 1$ 

$$f(x + \delta) \leq f(x) + \langle \nabla f(x), \delta \rangle + \delta^T \nabla^2 f(x) \delta$$

$$f(x) + \langle \nabla f(x), \delta \rangle + \frac{1}{6} \delta^T \nabla^2 f(x) \delta \leq f(x + \delta)$$

**Idea:** iteratively solve minimization problem  $\min_{||\delta||_{\infty} \leq 1} \langle \nabla f(x_t), \delta \rangle + \delta^T \nabla^2 f(x_t) \delta$ 

Then update  $x_{t+1} \leftarrow x_t + \delta$ .

$$f(x_{t+1}) - f(x^*) \le (1 - 1/||x_t - x^*||_{\infty})(f(x_t) - f(x^*))$$

# (Kind of) Faster Algorithm & Analysis

### Algorithm [ALOW'17, CMTV'17]

• Start with  $x_0 = 1$ ,  $\ell = O(R \cdot log(1/\epsilon))$ .

• For 
$$t = 0$$
 to  $\ell - 1$   
 $> f^{(t)}(y) = f(x_t + y).$   
 $> q_t$  quadratic-approximation to  $f^{(t)}$ .  
 $> y_t = \operatorname{argmin}_{||y||_{\infty} \le 1} q_t(y).$   
 $> x_{t+1} = x_t + y_t.$   
• Return  $x_{\ell}$ .

## Analysis:

- 1. There is approx. minimizer  $x^* \in B_{\infty}(0, R)$  (add regularizer)
- 2. Each step gets us  $\times (1 1/R)$  closer to OPT
- 3. After  $Rlog(1/\epsilon)$  iterations  $f(x) f(x^*) \le \epsilon$
- 4. This x gives us  $\epsilon$ -approximate scaling

 $\mathbf{A} \in M_n(\mathbb{R}_{\geq 0})$  input matrix.

$$f(x) = \sum_{1 \le i \le n} log\left(\sum_{j} A_{ij} e^{x_j}\right) - \sum_{j} x_j$$

Let

$$(A_x)_{ik} = \frac{A_{ik}e^{x_k}}{\sum_j A_{ij}e^{x_j}}$$

**Claim:**  $||\nabla f(z)||_2^2 = ds(A_z)$ 

If z s.t.  $f(z) \leq inf_{x>0}f(x) + \epsilon$  and  $||\nabla f(z)||_2^2 \leq \epsilon$  thus  $ds(A_z) \leq \epsilon$ 

Thus  $\epsilon$ -close to DS.

A completely positive operator is any map  $T: M_n(\mathbb{C}) \to M_n(\mathbb{C})$ given by  $(A_1, \dots, A_m)$  s.t.

$$T(X) = \sum_{1 \le i \le m} A_i X A_i^{\dagger}$$

Such maps take psd matrices to psd matrices.

Dual of  $\mathbf{T}(\mathbf{X})$  is map  $\mathbf{T}^*: \mathbf{M}_n(\mathbb{C}) \to \mathbf{M}_n(\mathbb{C})$  given by:

$$T^*(X) = \sum_{1 \le i \le m} A_i^{\dagger} X A_i$$

- Analog of scaling?
- Doubly stochastic?

## **Operator Scaling**

A quantum operator  $T: M_n(\mathbb{C}) \to M_n(\mathbb{C})$  is **doubly** stochastic (DS) if  $T(I) = T^*(I) = I$ .

Scaling of T(X) consists of  $L, R \in GL_n(\mathbb{C})$  s.t.

$$(A_1, \ldots, A_m) \rightarrow (LA_1R, \ldots, LA_mR)$$

Distance to doubly-stochastic:

$$ds(T) \stackrel{\text{\tiny def}}{=} \|T(I) - I\|_F^2 + \|T^*(I) - I\|_F^2$$

T(X) has approx. DS scaling if  $\forall \epsilon > 0$ ,  $\exists$  scaling  $L_{\epsilon}$ ,  $R_{\epsilon}$  s.t. operator  $T_{\epsilon}(X)$  given by  $(L_{\epsilon}A_{1}R_{\epsilon}, ..., L_{\epsilon}A_{m}R_{\epsilon})$  has  $ds(T_{\epsilon}) \leq \epsilon$ .

- 1. When does  $(A_1, ..., A_m)$  have approx. DS scaling?
- 2. Can we find it efficiently?

**NO** convex formulation!

### **Previous work**

**Problem:** operator  $\mathbf{T} = (A_1, \dots, A_m)$ ,  $\epsilon > 0$ , can T be  $\epsilon$ -scaled to double stochastic? If yes, find scaling.

### Algorithm G [Gurvits' 04, GGOW'15]:

Repeat  $k = poly(n, 1/\epsilon)$  times:

- 1. Left normalize T(X), i.e.,  $(A_1, \dots, A_m) \leftarrow (LA_1, \dots, LA_m)$ s.t. T(I) = I.
- 2. Right normalize  $\mathbf{T}(\mathbf{X})$ , i.e.,  $(A_1, \dots, A_m) \leftarrow (A_1R, \dots, A_mR)$ s.t.  $T^*(I) = I$ .

If at any point  $ds(T) \le \epsilon$ , output the current scaling. Else output **no scaling**.

### Potential Function (Capacity) [Gur'04]:

$$cap(T) = inf\left\{\frac{det(T(X))}{det(X)}: X > 0\right\}.$$

For  $\epsilon < 1/n^2$ , can scale **T** to  $\epsilon$ -close to DS iff cap(T) > 0.

## **Previous work – Analysis**

### **Algorithm G:**

Repeat *k* times:

1. Left normalize:  $(A_1, \dots, A_m) \leftarrow (RA_1, \dots, RA_m)$  s.t. T(I) = I.

2. Right normalize:  $(A_1, ..., A_m) \leftarrow (A_1C, ..., A_mC)$  s.t.  $T^*(I) = I$ . If at any point T(X) is close to DS, output current scaling. Else output **no scaling**.

Potential Function (Capacity) [Gur'04]:

$$cap(T) = inf\left\{\frac{det(T(X))}{det(X)}: X > 0\right\}.$$

Analysis [Gur'04, GGOW'15]:

- 1.  $cap(T) > 0 \Rightarrow cap(T) > e^{-poly(n)}$  (GGOW'15)
- 2.  $ds(T) \Rightarrow cap(T)$  grows by (1 + 1/n) after normalization
- 3.  $cap(T) \le 1$  for normalized operators.

# Potential Function (Capacity) [Gur'04]: $cap(T) = inf \left\{ \frac{det(T(X))}{det(X)} : X > 0 \right\}.$

For  $\epsilon < 1/n^2$ , can scale **T** to  $\epsilon$ -close to DS iff cap(T) > 0.

How can we decide if cap(T) > 0? Can we approx. capacity?

**[GGOW'15]:** natural scaling algorithm decides whether cap(T) > 0in deterministic poly(n) time. Moreover, it finds  $exp(\epsilon)$ -approx. to capacity in time  $poly(n, 1/\epsilon)$ .

Can we get convergence in  $\log\left(\frac{1}{\epsilon}\right)$ ? Need a different algorithm!

Capacity: optimization problem over *Positive Definite* matrices Is capacity a special function in this manifold? Generalizes Euclidean convexity to Riemannian manifolds.

- $\mathbb{R}^n$  becomes a smooth manifold (locally looks like  $\mathbb{R}^n$ )
- Straight lines become geodesics ("shortest paths")

**Example (our setup):** complex positive definite matrices  $S_+$  with geodesic from A to B given by:

$$\gamma_{A,B}: [0,1] \to \mathcal{S}_+ \qquad \gamma_{A,B}(t) = A^{1/2} (A^{-1/2} B A^{-1/2})^t A^{1/2}$$

### **Convexity**:

- $\mathbf{K} \subseteq S_+$  g-convex if  $\forall A, B \in K$  geodesic from A to B in K
- Function  $f : K \to \mathbb{R}$  is g-convex if univariate function  $f(\gamma_{A,B}(t))$  is convex in t for any  $A, B \in K$

Geodesically convex functions over  $S_+$ :

- $\log(\det(T(X)))$
- log(det(X)) (geodesically linear)

Thus log of capacity  $\stackrel{\text{def}}{=} \log(\det(T(X))) - \log(\det(X))$  g-convex!

For  $log(1/\epsilon)$  convergence, need new opt. tools for g-convex fncs.

Known approaches for g-convex functions:

• **[Folklore]** g-self-concordant functions converge in time  $poly(n \cdot log(1/\epsilon))$ .

No analog of ellipsoid or interior point method known for this setting.

Self concordance:  $f : \mathbb{R} \to \mathbb{R}$  is self concordant if  $|f'''(x)| \le 2(f''(x))^{3/2}$ 

 $f : \mathbb{R}^n \to \mathbb{R}$  self concordant if self concordant along each line.

 $h: \mathcal{S}_+ \to \mathbb{R}$  g-self concordant if self concordant along each geodesic.

Unfortunately, log of capacity **NOT** self-concordant.

Self robustness:  $f : \mathbb{R} \to \mathbb{R}$  is self robust if  $|f'''(x)| \le 2 \cdot f''(x)$ 

 $f : \mathbb{R}^n \to \mathbb{R}$  self robust if self robust along each line.

 $h: \mathcal{S}_+ \to \mathbb{R}$  g-self robust if self robust along each geodesic.

Log of capacity is self-robust!

**Question:** Can we efficiently optimize g-self robust functions?

# This work – g-convex opt for self-robust fcns

**Problem:** given  $f : S_+ \to \mathbb{R}$  g-self robust,  $\epsilon > 0$ , and bound on initial distance R to OPT (diameter) find  $X_{\epsilon} \in S_+$  such that

 $f(X_{\epsilon}) \leq \inf_{Y \in \mathcal{S}_{+}} f(Y) + \epsilon$ 

### Theorem [AGLOW'18]:

There exists a deterministic  $poly(n, R, log(1/\epsilon))$ , algorithm for the problem above.

- Second order method, generalizing recent work of [ALOW'17, CMTV'17] for matrix scaling to g-convex setting (Box constrained Newton method)
- Generalizes to other manifolds and metrics

#### **Remark:**

• For operator scaling,  $X_{\epsilon}$  also gives us scaling  $\epsilon$ -close to DS

## This paper – g-convex opt for self-robust fcns

**Problem:** given  $f : S_+ \to \mathbb{R}$  g-self robust,  $\epsilon > 0$ , and bound on initial distance R to OPT (diameter) find  $X_{\epsilon} \in S_+$  such that

$$f(X_{\epsilon}) \leq \inf_{Y \in S_+} f(Y) + \epsilon$$

#### Algorithm

• Start with  $X_0 = I$ ,  $\ell = O(R \cdot log(1/\epsilon))$ .

• For 
$$t = 0$$
 to  $\ell - 1$   
 $\geq f^{(t)}(D) = f(X_t^{1/2} \exp(D)X_t^{1/2}).$   
 $\geq Q_t$  quadratic-approximation to  $f^{(t)}$ .  
 $\geq D_t = \operatorname{argmin}_{||D||_F \leq 1} Q_t(D).$  (Euclidean convex opt.)  
 $\geq X_{t+1} = X_t^{1/2} \exp(D_t) X_t^{1/2}.$   
• Return  $X_\ell$ .

- Why would we need this instead of regular scaling?
- What is the bound for **R** in operator scaling?
  - **[AGLOW'18]** polynomial bound for **R**

### **Invariant Theory:**

 $G = \mathbb{SL}_n(\mathbb{C})^2$ , vector space  $V = M_n(\mathbb{C})^m$  action by L-R mult:  $(A_1, \dots, A_m) \to (LA_1R, \dots, LA_mR)$ 

**Orbit Closure:** given  $v = (A_1, ..., A_m) \in V$ , orbit closure is  $\overline{\mathcal{O}_v} = \overline{\{(LA_1R, ..., LA_mR) \mid (L, R) \in G\}}$ 

**Orbit Closure Intersection Problem:** given two quantum operators  $u = (A_1, ..., A_m), v = (B_1, ..., B_m)$ , is  $\overline{\mathcal{O}_u} \cap \overline{\mathcal{O}_v} \neq \emptyset$ ?

If v = 0 problem becomes the *null-cone problem*. [GGOW'16]: connections to non-commutative PIT, non-commutative algebra, combinatorics, functional analysis...

How can we solve the orbit intersection problem for L-R action?

[Mum'65]: alg. structure of orbit closures

• 
$$\overline{\mathcal{O}_{(A_1,\dots,A_m)}} \cap \overline{\mathcal{O}_{(B_1,\dots,B_m)}} = \emptyset$$
 iff invariant polynomial s.t.  
 $p((A_1,\dots,A_m)) \neq p((B_1,\dots,B_m))$ 

Randomized algorithm:

Given  $(A_1, \dots, A_m)$  and  $(B_1, \dots, B_m)$ , does  $\overline{\mathcal{O}_{(A_1, \dots, A_m)}} \cap \overline{\mathcal{O}_{(B_1, \dots, B_m)}} \neq \emptyset$ ?

- 1. [IQS'17, DM'17]: Invariants of degree  $n^6$  suffice
- 2. Take random invariant polynomial and evaluate it on  $(A_1, \dots, A_m)$  and  $(B_1, \dots, B_m)$

## KN'79 – Duality Theory

### [KN'79]:

- Elts of min norm in  $\overline{\mathcal{O}_{(A_1,\dots,A_m)}}$ , are DS operators
  - $\epsilon$ -close to DS implies  $\epsilon$ -close to min. norm
- $(B_1, ..., B_m)$  and  $(C_1, ..., C_m)$  elts of min norm in  $\overline{\mathcal{O}_{(A_1, ..., A_m)}}$ then there exist  $U, V \in SU(n)$  s.t.  $C_i = UB_i V$

**[AGLOW'18]:** solving orbit closure intersection problem. Given  $(A_1, \ldots, A_m)$  and  $(B_1, \ldots, B_m)$ , does  $\overline{\mathcal{O}_{(A_1, \ldots, A_m)}} \cap \overline{\mathcal{O}_{(B_1, \ldots, B_m)}} \neq \emptyset$ 

- 1. Our g-convex opt finds  $\epsilon$ -approx to element of min norm (DS)
- 2. With elements of min norm, test if they are SU(n)-equivalent
  - we give efficient algorithm for testing equivalence

Why do we need  $\log(1/\epsilon)$  convergence?

- Orbit closures can be exponentially close and not intersect
  - Need to have  $\epsilon = \exp(-poly(n))$  approximation
  - Not the case for null-cone problem
- **SU**(**n**)-equivalence algorithm also approximate (and lossy)

Independently, **[DM'18]** solved orbit closure intersection for LR-action in algebraic way.

- Solution also works for fields of positive characteristic
  - Our solution works only over  ${\mathbb C}$

Prior to **[AGLOW'18, DM'18]** only *randomized* polynomial time algorithm known for orbit closure intersection (PIT instance).

# **Open questions**

- Efficient algorithms for more classes of g-convex functions?
- Efficient algorithms for null-cone and orbit closure intersection for more general actions?
  - Recent developments for tensor scaling, though still  $poly(1/\epsilon)$
  - Upcoming work gets  $poly(nR \cdot log(1/\epsilon))$ , but still have bad bounds on R
- More applications of g-convexity?
  - Recent work [VY'18] on Brascamp-Lieb showing it is g-convex

# Thank you!

