

# Finding Hall blockers by matrix scaling

based on

*--- convergence analysis on “divergent” Sinkhorn iteration ---*

Hiroshi Hirai

Graduate School of Mathematics

Nagoya University

hirai.hiroshi@math.nagoya-u.ac.jp

Joint work with

Koyo Hayashi, Keiya Sakabe

arXiv:2204.07425 v2

Discrete Optimization and Machine Learning

GRIPS, August 8, 2023

# Matrix Scaling (Sinkhorn 1964)

$A: n \times n$  nonnegative matrix

*Can we scale  $A$  to doubly stochastic matrix  $RAC$   
by multiplying positive diagonal matrices  $R, C$  ?*

**Goal:** Find  $R, C$  s.t.  $RAC \mathbf{1} \approx \mathbf{1}, (RAC)^\top \mathbf{1} \approx \mathbf{1}$

## **Applications:**

- Markov chain
  - Preprocessing for solving linear equation
  - Optimal transport, entropic regularization (Cuturi 2013)
- and more

# Sinkhorn Algorithm

- Row normalization:  $A \leftarrow RA$  s.t.  $(RA)\mathbf{1} = \mathbf{1}$ ;  $R = \text{diag}(\dots (\sum_k A_{ik})^{-1} \dots)$
- Col normalization:  $A \leftarrow AC$  s.t.  $(AC)^\top \mathbf{1} = \mathbf{1}$ ;  $C = \text{diag}(\dots (\sum_k A_{kj})^{-1} \dots)$
- Repeat it

$$\begin{array}{c} 6 \\ 3 \\ 4 \\ 3 \end{array} \begin{pmatrix} 4 & 10 & 1 & 1 \\ & 6 & & \\ 1 & & 1 & 1 \\ & 4 & & \\ 3 & & & \end{pmatrix} \xrightarrow{\text{Row normalize}} \begin{array}{c} 1 \\ 1 \\ 1 \\ 1 \end{array} \begin{pmatrix} 4/3 & 2 & 1/3 & 1/3 \\ & 1 & & \\ 1/3 & & 1/3 & 1/3 \\ & 1 & & \\ 1 & & & \end{pmatrix}$$

$$\begin{array}{c} \text{Col normalize} \\ \rightarrow \end{array} \begin{array}{c} 1/2 \\ 9/4 \\ 1/2 \\ 3/4 \end{array} \begin{pmatrix} 1 & 1 & 1 & 1 \\ & 1/2 & & \\ 1/4 & & 1 & 1 \\ & 1/2 & & \\ 3/4 & & & \end{pmatrix} \xrightarrow{\text{Row normalize}} \dots$$

# Characterization of Scalability

Sinkhorn, Knopp 1967, Rothblum, Schneider 1989

Thm [SK1967, RS1989]: The following are equivalent:

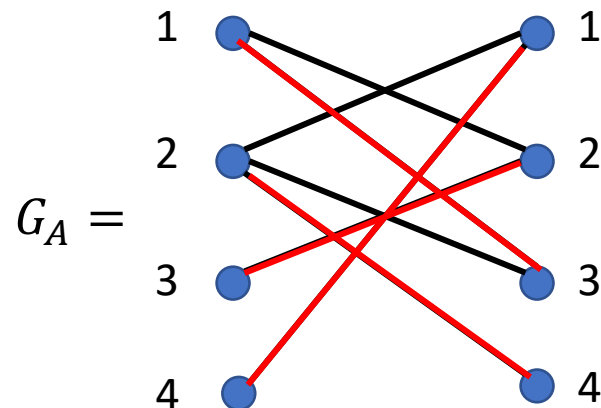
- $A$  is (approximately) doubly stochastic scalable:

$$\forall \epsilon > 0, \exists R, C \text{ s.t. } \|RAC\mathbf{1} - \mathbf{1}\| < \epsilon, \|(RAC)^\top \mathbf{1} - \mathbf{1}\| < \epsilon$$

- Sinkhorn algorithm **converges**;  $A \rightarrow$  doubly stochastic
- $\exists$  **perfect matching** in bipartite graph  $G_A$

$$V(G_A) := [n] \sqcup [n], E(G_A) := \{ij \mid A_{ij} \neq 0\}$$

$$A = \begin{matrix} & \begin{matrix} 1 & 2 & 3 & 4 \end{matrix} \\ \begin{matrix} 1 \\ 2 \\ 3 \\ 4 \end{matrix} & \begin{pmatrix} & & & \\ 1 & 6 & 3 & 4 \\ & 1 & 1 & 1 \\ & & 4 & \\ & 3 & & \end{pmatrix} \end{matrix}$$



# Testing Perfect Matching by Sinkhorn

Linial, Samorodnitsky, Wigderson 2000

$G = (U \sqcup V, E)$ : bipartite graph,  $|U| = |V| = n$

$$A(G)_{ij} := \begin{cases} 1 & \text{if } ij \in E \\ 0 & \text{otherwise} \end{cases}$$

Initialization:  $A \leftarrow A(G)$

1. Row-normalize:  $A \leftarrow \text{diag}(1/\sum_k A_{ik})A$
2. Col-normalize:  $A \leftarrow A \text{diag}(1/\sum_k A_{kj})$
3. If  $\|A\mathbf{1} - \mathbf{1}\|_2 < 1/\sqrt{n}$ , then stop.
4. Go to 1.

$$\begin{array}{c} 1 \ 1 \ 1 \ \dots \ 1 \\ p_1 \\ p_2 \\ \vdots \\ p_n \end{array} \begin{array}{|c|} \hline A \\ \hline \end{array}$$

Thm [LSW2000]

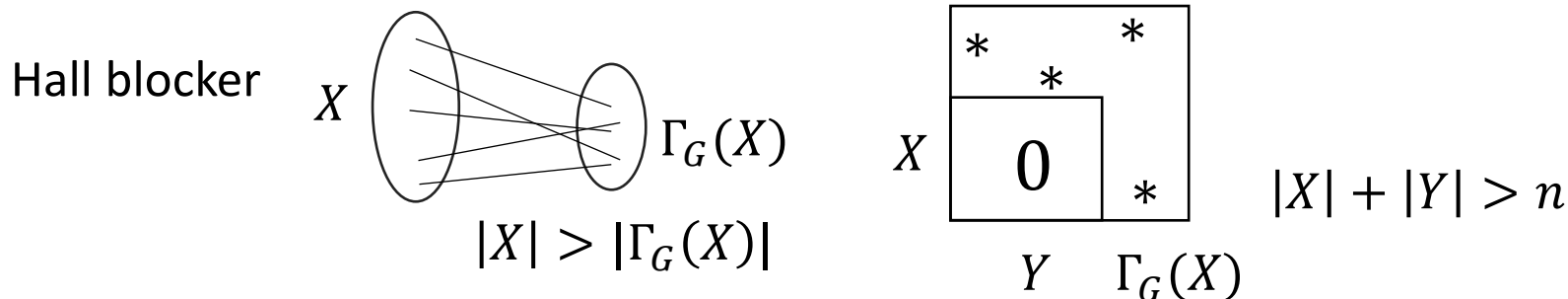
$\exists$  perfect matching in  $G \rightarrow$  terminates within  $O(n^2 \log n)$  iterations

$\nexists$  perfect matching in  $G \rightarrow$  not terminate

- LSW algorithm is slower than augmenting path  
But interesting in its conceptual difference & extraordinary simplicity.
- links to recent development on *operator scaling* & *noncommutative PIT*  
Garg, Gurvits, Oliveira, Wigderson FOCS2016
- LSW algorithm outputs **neither** perfect matching **nor** *certificate of nonexistence*



Hall's marriage theorem :  $\exists$  perfect matching  $\Leftrightarrow \nexists$  Hall blocker



*Can we identify Hall blockers from (divergent !) Sinkhorn iteration ?*

Significant in operator scaling generalization: Franks, Soma, Goemans SODA2023  
But not well-understood even in matrix scaling setting

# Result 1

$O(n^2 \log n)$  Sinkhorn iterations identify a Hall blocker  
from **scaling matrices  $R, C$**

0.  $A \leftarrow A(G), R = C = I$
1. Row-normalize:  $R \leftarrow \text{diag}\left(\left(\sum_j A_{1j}C_{jj}\right)^{-1}, \left(\sum_j A_{2j}C_{jj}\right)^{-1}, \dots, \left(\sum_j A_{nj}C_{jj}\right)^{-1}\right)$
2. Col-normalize:  $C \leftarrow \text{diag}\left(\left(\sum_i R_{ii}A_{i1}\right)^{-1}, \left(\sum_i R_{ii}A_{i2}\right)^{-1}, \dots, \left(\sum_i R_{ii}A_{in}\right)^{-1}\right)$
3. Sort as  $R_{11} \geq R_{22} \geq \dots \geq R_{nn}, C_{11} \leq C_{22} \leq \dots \leq C_{nn}$ ,  
Choose  $k$  s.t.  $(R_{kk}C_{kk})^n \geq R_{11}R_{22} \dots R_{nn}C_{11}C_{22} \dots C_{nn}$ ,  
Output  $X := \{1, 2, \dots, k\}$
4. Go to 1.

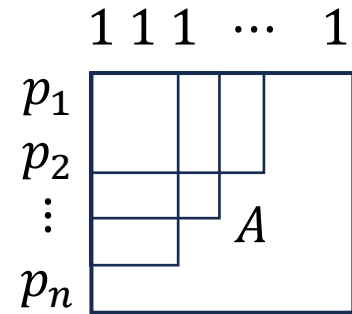
Rounding produce by  
Franks, Soma Goemans SODA2023

Thm [This work] Suppose that  $G$  has no perfect matching  
 $X$  is a Hall blocker within  $O(n^2 \log n)$  iterations

## Result 2

$O(n^6 \log n)$  Sinkhorn iterations identify all parametric Hall blockers  
from **row-marginals**  $\mathbf{p} := A\mathbf{1}$

0.  $A \leftarrow A(G)$
1. Row-normalize:  $A \leftarrow \text{diag}(1/\sum_k A_{ik})A$
2. Col-normalize:  $A \leftarrow A \text{diag}(1/\sum_k A_{kj})$
3. Sort  $\mathbf{p} := A\mathbf{1}$  as  $p_1 \geq p_2 \geq \dots \geq p_n$   
Output  $\mathcal{X} := \{ \{1, 2, \dots, k\} \mid k = 1, 2, \dots, n \}$ .
4. Go to 1.



Thm [This work] Suppose that  $G$  has no perfect matching.

$\mathcal{X}$  contains all *parametric Hall blockers* within  $O(n^6 \log n)$  iterations

maximizing  $|X| - \alpha |\Gamma_G(X)|$  ( $\alpha \in \mathbb{R}_+$ )



## Proof idea

Sinkhorn iteration = Alternating minimization

**Result 1 :  $O(n^2 \log n)$  iterations identify “a” Hall blocker**

~ geometric programming interpretation

$$\inf. \log \frac{(x^\top A y)^n}{\prod_i x_i \prod_j y_j} \quad \text{s.t.} \quad x > 0, y > 0.$$

Fix  $y$  optimize  $x$ . Fix  $x$  optimize  $y$  ...

**Result 2 :  $O(n^6 \log n)$  iterations identify “all parametric” Hall blockers**

~ KL-divergence minimization interpretation

$$\inf. \sum_{ij} M_{ij} \log \frac{M_{ij}}{N_{ij}} \quad \text{s.t.} \quad M \mathbf{1} = \mathbf{1}, N^\top \mathbf{1} = \mathbf{1},$$

$$M, N \geq 0, \quad \text{supp } M, \text{supp } N \subseteq \text{supp } A$$

Fix  $N$  optimize  $M$ . Fix  $M$  optimize  $N$  ...

# Matrix Scaling as Geometric Programming

Thm (SK1967, RS1989) The following are equivalent:

1.  $A$  is approximately scalable.
2. Sinkhorn converges.
3. Geometric programming is **bounded below**:

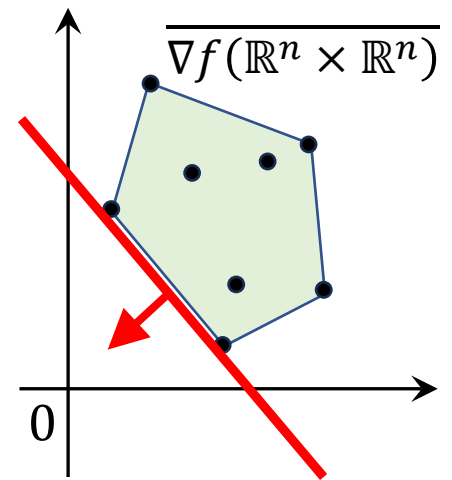
$$\inf_{s,t \in \mathbb{R}^n} f(s,t) := \log \sum_{i,j} A_{ij} e^{s_i + t_j} - \frac{1}{n} \mathbf{1}^\top s - \frac{1}{n} \mathbf{1}^\top t > -\infty$$

Scaling matrix  $R, C \iff \text{diag}(e^{s_i}), \text{diag}(e^{t_j})$

Sinkhorn algorithm  $\iff$  alternating minimization

$$\begin{aligned} \inf_{s,t \in \mathbb{R}^n} f(s,t) > -\infty &\iff 0 \in \overline{\nabla f(\mathbb{R}^n \times \mathbb{R}^n)} \\ &= -\infty \iff 0 \notin \overline{\nabla f(\mathbb{R}^n \times \mathbb{R}^n)} \end{aligned}$$

||



Hall blocker  $\approx$  separating facet of  $\text{conv} \left\{ (e_i, e_j) - \frac{1}{n} (\mathbf{1}, \mathbf{1}) \mid i, j : A_{ij} > 0 \right\}$

# Unbounded Certificate in Geometric Programming

Lem/Obs [This work ?]

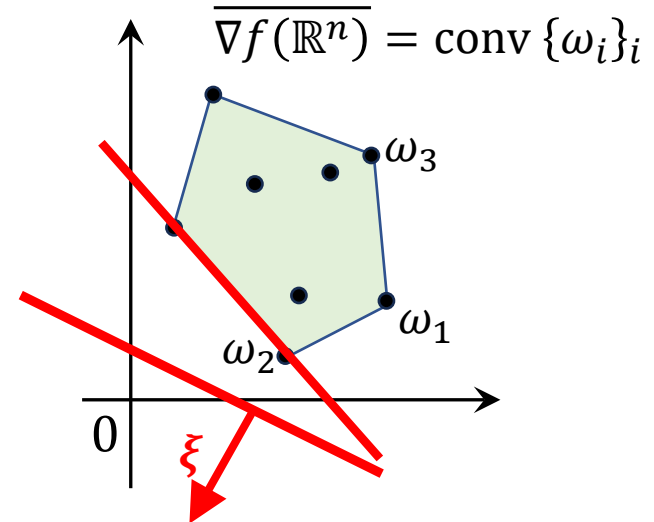
$$g(\xi) := \log \left( \sum_{i=1}^m q_i e^{\omega_i^\top \xi} \right), q_i > 0$$

If  $g(\xi) \leq \log \min\{q_1, \dots, q_m\}$ ,

then  $\{0\}$  and  $\overline{\nabla f(\mathbb{R}^n)}$  are separated

by the hyperplane of normal vector  $\xi$

$$\rightarrow \inf g = -\infty$$



$$\because \omega_k^\top \xi \geq 0 (\exists k) \rightarrow g(\xi) = \log \left( \sum_{i=1}^m q_i e^{\omega_i^\top \xi} \right) > \log q_k$$

We analyze

- If unbounded, “many” iterations make  $\xi = (s, t)$  a separating hyperplane
- Round the hyperplane to a separating facet ( $\approx$  Hall blocker)

Rounding procedure: Franks, Soma, Goemans SODA2023

# Analysis (can be skipped)

Decrement of one Sinkhorn iteration:  $s, t \rightarrow s', t \rightarrow s', t'$ ;  $A \xrightarrow{\text{Row}} A' \xrightarrow{\text{Col}} A''$

Lem [folklore ?; Altschuler, Weed, Rigolle NIPS2017]

$$f(s, t) - f(s', t') = D_{KL}(\mathbf{1}|A\mathbf{1})/n + D_{KL}(\mathbf{1}|A'^T\mathbf{1})/n$$

$$\text{ Pinsker's ineq } \geq \frac{1}{2n^2} \{ \|A\mathbf{1} - \mathbf{1}\|_1^2 + \|A'^T\mathbf{1} - \mathbf{1}\|_1^2 \}$$

Lem [Gurvits, Leake STOC2021 ]:

$A$ : nonnegative,  $\text{supp } A = E(G)$

$$\|A\mathbf{1} - \mathbf{1}\|_1 + \|A^T\mathbf{1} - \mathbf{1}\|_1 \geq 2 \max_X |X| - |\Gamma_G(X)|$$

$\geq 2$  If no perfect matching

- $A := A_G$ : 0-1 matrix for  $G$  without perfect matching
- $f(0,0) = O(\log n)$ , decrement of one iteration =  $\Omega(1/n^2)$
- $f(s, t) \leq \log 1 = 0$  after  $O(n^2 \log n)$  iterations
- $s, t$ : separating hyperplane  $\rightarrow$  Hall blocker

## Proof idea

Sinkhorn iteration = Alternating minimization

**Result 1 :  $O(n^2 \log n)$  iterations identify “a” Hall blocker**

~ geometric programming interpretation

$$\inf. \log \frac{x^\top A y}{\prod_i x_i \prod_j y_j} \quad \text{s.t.} \quad x > 0, y > 0.$$

Fix  $y$  optimize  $x$ . Fix  $x$  optimize  $y$  ...

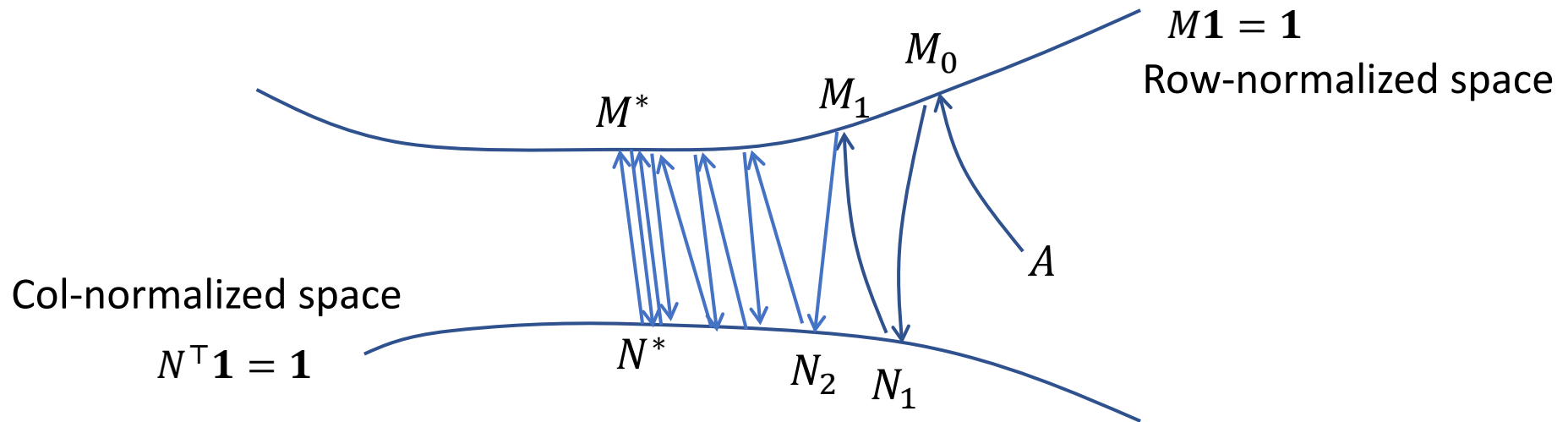
**Result 2 :  $O(n^6 \log n)$  iterations identify “all parametric” Hall blockers**

~ KL-divergence minimization interpretation

$$\inf. \sum_{ij} M_{ij} \log \frac{M_{ij}}{N_{ij}} \quad \text{s.t.} \quad M \mathbf{1} = \mathbf{1}, N^\top \mathbf{1} = \mathbf{1},$$
$$M, N \geq 0, \quad \text{supp } M, \text{supp } N \subseteq \text{supp } A$$

Fix  $N$  optimize  $M$ . Fix  $M$  optimize  $N$  ...

# Information geometry of Sinkhorn iteration



$$\inf. D_{KL}(N|M) \quad \text{s.t. } M\mathbf{1} = \mathbf{1}, N^T\mathbf{1} = \mathbf{1},$$

$$M, N \geq 0, \quad \text{supp } M, \text{supp } N \subseteq \text{supp } A$$

Thm (Csiszar, Tusnady 1984, Gietl, Reffel 2013)

$\{(M_k, N_k)\}_{k=1,2,\dots}$  converges to a minimum KL-divergence pair  $(M^*, N^*)$

- $M^* = N^*$  if Sinkhorn converges
- $M^* \leftrightarrow N^*$  if Sinkhorn does not converge → DEMO  
oscillating

# Rate of Convergence

We point out that proof argument of CT1984 implies:

Lem [CT1984; This work]

$$D_{KL}(N_k | M_k) - D_{KL}(N^* | M^*) \leq \frac{D(N^* | M_0)}{k}$$

||

$$\begin{aligned}
 p^k &:= N_k \mathbf{1} & D_{KL}(p^k | \mathbf{1}) - D_{KL}(p^* | \mathbf{1}) &\geq D_{KL}(p^k | p^*) \geq \frac{1}{2n} \|p^k - p^*\|_1^2 \\
 p^* &:= N^* \mathbf{1} & &\text{Pythagorean thm} \\
 &= \operatorname{argmin}_{N^\top \mathbf{1} = \mathbf{1}} D_{KL}(N \mathbf{1} | \mathbf{1}) & &\text{Pinsker}
 \end{aligned}$$

Lem [This work]

$$\|p^k - p^*\|_1 \leq \sqrt{\frac{2n D_{KL}(N^* | M_0)}{k}}$$

We give an *explicit* formula of the marginal limit  $p^* = N^* \mathbf{1}$   
in terms of **DM-decomposition & parametric Hall blockers**

# The Sinkhorn Limit $N^* \leftrightarrow M^*$

Aas 2014: The limits  $N^* \leftrightarrow M^*$  are **block diagonalized** so that each block is oscillated as

$$\begin{array}{c}
 1 \ 1 \ 1 \ \dots \ 1 \\
 \alpha \ \boxed{\phantom{0000}} \\
 \alpha \\
 \vdots \\
 \alpha
 \end{array}
 \Leftrightarrow
 \begin{array}{c}
 \frac{1}{\alpha} \ \frac{1}{\alpha} \ \frac{1}{\alpha} \ \dots \ \frac{1}{\alpha} \\
 1 \ \boxed{\phantom{0000}} \\
 1 \\
 \vdots \\
 1
 \end{array}$$

$\alpha := \text{col number} / \text{row number}$

This work:

This block diagonalization = extended **Dulmage-Mendelsohn decomposition**

Canonical form of matrices  
under row/column permutation

→ DEMO



# The DM-decomposition

Dulmage-Mendelson 1958

DM-decomposition =

Permute rows & columns along  
maximum-size Hall blockers

Size of zero block

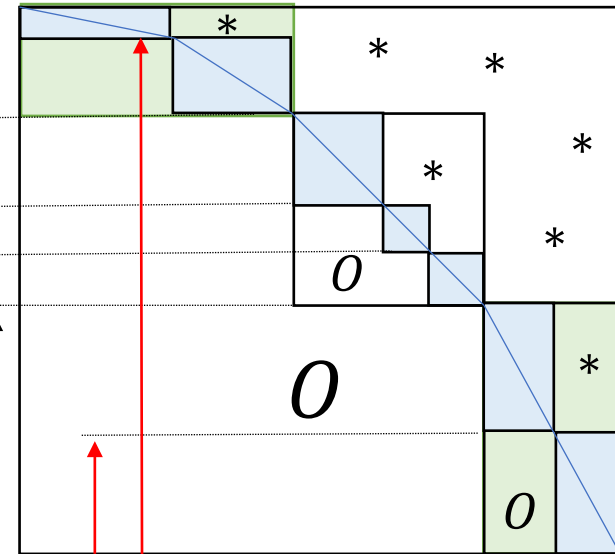
$$:= \# \text{ row} + \# \text{ col}$$

$$\approx |X| - |\Gamma(X)| + \text{const}$$

$A$

$\rightarrow$

Maximum  
Hall blocker



Remaining  
part

Tomizawa 1977 (unpublished)

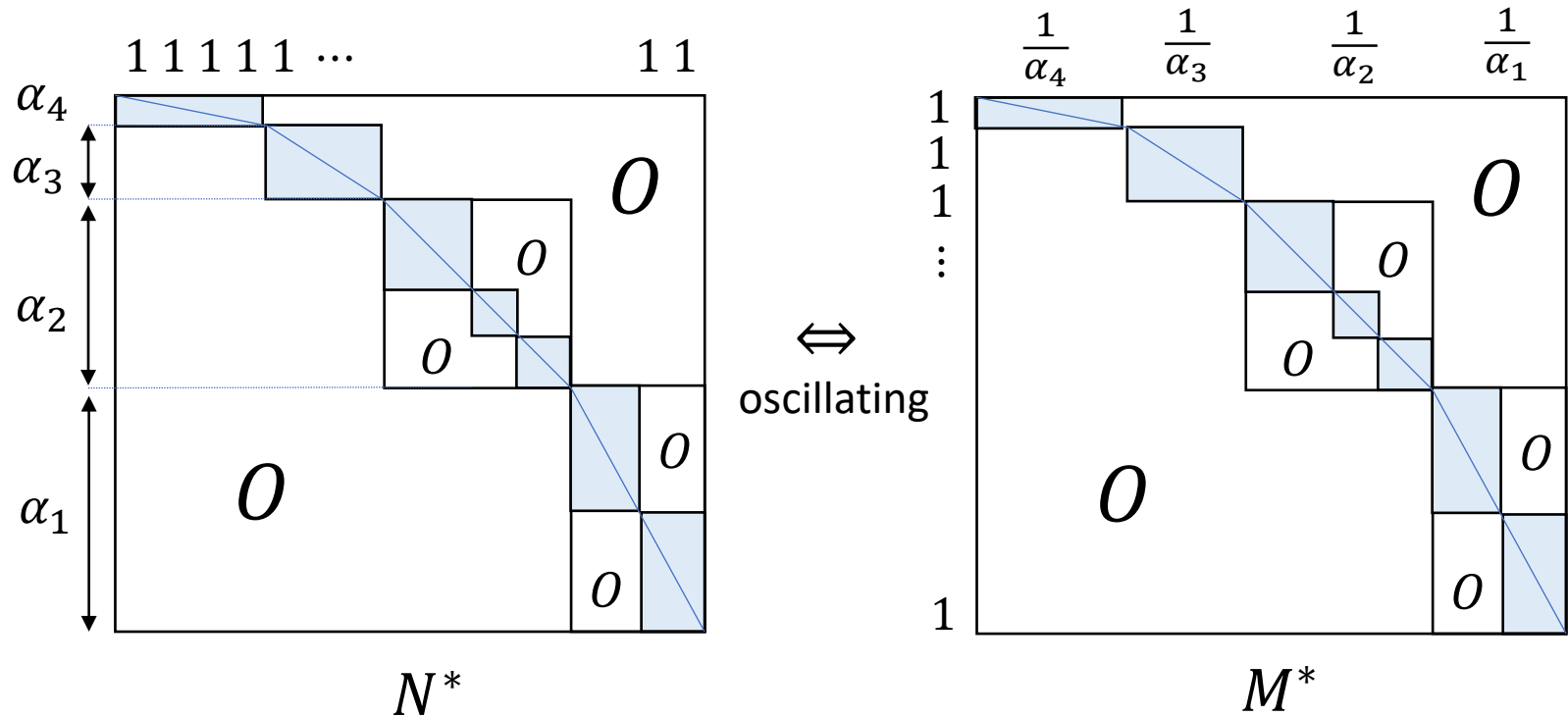
Extended DM decomposition = decompose “remaining parts”

via parametric stable sets in bipartite graph  $G(A)$

Weighted size of zero block  $:= \# \text{ row} + \alpha \times \# \text{ col}$

$$\approx |X| - \alpha |\Gamma(X)| + \text{const}$$

# The Sinkhorn Limit $N^* \leftrightarrow M^*$



- From  $\alpha_1 > \alpha_2 > \alpha_3 > \dots$ , the row-sum vector  $p^* = N^* \mathbf{1}$  identifies the structure of DM-decomposition  $\rightarrow$  parametric Hall blockers
- So does  $p = A \mathbf{1}$  if it is close to  $p^* = N^* \mathbf{1}$  after  $O(n^6 \log n)$  iterations  
 $\leftarrow$  convergence rate  $O(1/\sqrt{k})$

# Summary

- Parametric Hall blockers from divergent Sinkhorn iteration
- Information geometry & DM-decomposition
- Can we improve  $O(n^6 \log n)$  ?

Original motivation: Operator scaling (Gurvits 2004, Garg et al. 2020)

Completely positive operator  $X \mapsto \sum_k A_k X A_k^\dagger$  ( $A_1, A_2, \dots, A_m \in \mathbb{C}^{n \times n}$ )

doubly-stochastic  $\Leftrightarrow$   
def  $\sum_k A_k A_k^\dagger = I, \sum_k A_k^\dagger A_k = I$

Find *scaling*  $g, h \in GL_n(\mathbb{C})$  s.t.  $X \mapsto \sum_k g A_k h X h^\dagger A_k^\dagger g^\dagger$  is doubly-stochastic

- Operator Sinkhorn algorithm (Gurvits algorithm)
- Shrunk subspace = a certificate of nonscalability (analogue of Hall blocker)

## Operator Sinkhorn algorithm (Gurvits algorithm, flip-flop)

1. Row-normalize:  $A_k \leftarrow gA_k \quad : \quad g(\sum_k A_k A_k^\dagger)g^\dagger = I$
2. Col-normalize:  $A_k \leftarrow A_k h^\dagger \quad : \quad h(\sum_k A_k^\dagger A_k)h^\dagger = I$
3. Go to 1.

Thm (Gurvitz 2004): The following conditions are equivalent:

- $X \mapsto \sum_k A_k X A_k^\dagger$  is approximately scalable
- Operator Sinkhorn converges
- No vector subspace (*shrunk subspace*)  $V \subseteq \mathbb{C}^n$ :  $\dim V > \dim \sum_k A_k V$

matrix scaling case  
 $|X| > |\Gamma_G(X)|$

*Can we find shrunk subspaces by operator Sinkhorn iterations ?*

Franks, Soma, Goemans SODA2023: **YES** by *modified* Sinkhorn iterations.

Still complicated & huge (polynomial) iterations

*Thank you for your attention !*