

Matrix Scaling, Entropic Optimal Transport, and Sinkhorn Algorithm: Theory, Applications, and New Perspectives

Hanbaek Lyu

Department of Mathematics
University of Wisconsin – Madison

UW-Madison AI Seminar

Apr. 20, 2026

Collaborators and Students



Jakwang Kim
(U. HK, Shenzhen)



Sumit Mukherjee
(Columbia)



William Powell
(5th year Ph.D.)



Danny Duan
(5th year Ph.D.)



Shuqi Bi
(2nd year Ph.D.)



Rahul Choudhary
(5th year CS Ph.D.)

Introduction

Matrix Scaling
Sinkhorn Algorithm
Entropic OT
Wasserstein Distances

▶ Matrix Scaling Problem:

Given a nonnegative matrix \mathbf{A} and target row and column sums (\mathbf{r}, \mathbf{c}) , find diagonal matrices \mathbf{D}_1 and \mathbf{D}_2 s.t. $\mathbf{D}_1\mathbf{A}\mathbf{D}_2$ satisfies margin (\mathbf{r}, \mathbf{c})

▶ Matrix Scaling Problem:

Given a nonnegative matrix \mathbf{A} and target row and column sums (\mathbf{r}, \mathbf{c}) , find diagonal matrices \mathbf{D}_1 and \mathbf{D}_2 s.t. $\mathbf{D}_1\mathbf{A}\mathbf{D}_2$ satisfies margin (\mathbf{r}, \mathbf{c})

- Sinkhorn's matrix scaling algorithm (1964)
 - Normalize the rows of \mathbf{A} to match the target row sum \mathbf{r} ; Obtain \mathbf{A}'
 - Normalize the columns of \mathbf{A}' to match the target column sum \mathbf{c} ; Obtain \mathbf{A}''
 - Repeat

► Matrix Scaling Problem:

Given a nonnegative matrix \mathbf{A} and target row and column sums (\mathbf{r} , \mathbf{c}), find diagonal matrices \mathbf{D}_1 and \mathbf{D}_2 s.t. $\mathbf{D}_1\mathbf{A}\mathbf{D}_2$ satisfies margin (\mathbf{r} , \mathbf{c})

- Sinkhorn's matrix scaling algorithm (1964)
 - Normalize the rows of \mathbf{A} to match the target row sum \mathbf{r} ; Obtain \mathbf{A}'
 - Normalize the columns of \mathbf{A}' to match the target column sum \mathbf{c} ; Obtain \mathbf{A}''
 - Repeat
- Sinkhorn's algorithm is known to solve the following relative entropy minimization problem:

$$\arg \min_{\mathbf{Z} \in \mathcal{F}(\mathbf{r}, \mathbf{c})} [D_{KL}(\mathbf{Z} \| \mathbf{\Lambda}) = \sum_{ij} \mathbf{Z}_{ij} \log(\mathbf{Z}_{ij} / \mathbf{\Lambda}_{ij})]$$

└──┘
Transportation polytope
with margin (\mathbf{r} , \mathbf{c})

► Matrix Scaling Problem:

Given a nonnegative matrix \mathbf{A} and target row and column sums (\mathbf{r}, \mathbf{c}) , find diagonal matrices \mathbf{D}_1 and \mathbf{D}_2 s.t. $\mathbf{D}_1\mathbf{A}\mathbf{D}_2$ satisfies margin (\mathbf{r}, \mathbf{c})

- Sinkhorn's matrix scaling algorithm (1964)
 - Normalize the rows of \mathbf{A} to match the target row sum \mathbf{r} ; Obtain \mathbf{A}'
 - Normalize the columns of \mathbf{A}' to match the target column sum \mathbf{c} ; Obtain \mathbf{A}''
 - Repeat
- Sinkhorn's algorithm is known to solve the following relative entropy minimization problem:

$$\arg \min_{\mathbf{Z} \in \mathcal{T}(\mathbf{r}, \mathbf{c})} [D_{KL}(\mathbf{Z} \parallel \mathbf{\Lambda}) = \sum_{ij} \mathbf{Z}_{ij} \log(\mathbf{Z}_{ij} / \mathbf{\Lambda}_{ij})]$$

└──┘
Transportation polytope
with margin (\mathbf{r}, \mathbf{c})

- Why? It is in fact the **alternating maximization** on the "dual"!

► Kantorovich Duality and Sinkhorn's algorithm

Primal

$$\min_{\mathbf{Z} \in \mathcal{T}(\mathbf{r}, \mathbf{c})} [D_{KL}(\mathbf{Z} \parallel \mathbf{W})]$$

Rescaled Matrix

► Kantorovich Duality and Sinkhorn's algorithm

Primal

Dual

$$\min_{\mathbf{Z} \in \mathcal{T}(\mathbf{r}, \mathbf{c})} [D_{KL}(\mathbf{Z} \parallel \mathbf{W})] = \max_{\boldsymbol{\alpha}, \boldsymbol{\beta}} [\langle \mathbf{r}, \boldsymbol{\alpha} \rangle + \langle \boldsymbol{\beta}, \mathbf{c} \rangle - \langle \exp(\boldsymbol{\alpha} \oplus \boldsymbol{\beta}), \mathbf{W} \rangle]$$

Rescaled Matrix

$\boldsymbol{\alpha}$ = Lagrange multipliers for the row-sum constraint

$\boldsymbol{\beta}$ = Lagrange multipliers for the column-sum constraint

(Schrödinger potentials)

► Kantorovich Duality and Sinkhorn's algorithm

Primal

Dual

$$\min_{\mathbf{Z} \in \mathcal{T}(\mathbf{r}, \mathbf{c})} [D_{KL}(\mathbf{Z} \parallel \mathbf{W})] = \max_{\boldsymbol{\alpha}, \boldsymbol{\beta}} [\langle \mathbf{r}, \boldsymbol{\alpha} \rangle + \langle \boldsymbol{\beta}, \mathbf{c} \rangle - \langle \exp(\boldsymbol{\alpha} \oplus \boldsymbol{\beta}), \mathbf{W} \rangle]$$

Rescaled Matrix

$\boldsymbol{\alpha}$ = Lagrange multipliers for the row-sum constraint

$\boldsymbol{\beta}$ = Lagrange multipliers for the column-sum constraint

(Schrödinger potentials)

$$\boldsymbol{\pi}^* = \text{Diag}(\exp(\boldsymbol{\alpha}^*)) \mathbf{W} \text{Diag}(\exp(\boldsymbol{\beta}^*))$$

► Kantorovich Duality and Sinkhorn's algorithm

Primal

Dual

$$\min_{\mathbf{Z} \in \mathcal{T}(\mathbf{r}, \mathbf{c})} [D_{KL}(\mathbf{Z} \parallel \mathbf{W})] = \max_{\boldsymbol{\alpha}, \boldsymbol{\beta}} [\langle \mathbf{r}, \boldsymbol{\alpha} \rangle + \langle \boldsymbol{\beta}, \mathbf{c} \rangle - \langle \exp(\boldsymbol{\alpha} \oplus \boldsymbol{\beta}), \mathbf{W} \rangle]$$

Rescaled Matrix

$\boldsymbol{\alpha}$ = Lagrange multipliers for the row-sum constraint

$\boldsymbol{\beta}$ = Lagrange multipliers for the column-sum constraint

(Schrödinger potentials)

$$\boldsymbol{\pi}^* = \text{Diag}(\exp(\boldsymbol{\alpha}^*)) \mathbf{W} \text{Diag}(\exp(\boldsymbol{\beta}^*))$$

Matrix scaling algorithm
(Iterative Proportional Fitting)

Fit the row sums;
Fit the column sums; etc.

► Kantorovich Duality and Sinkhorn's algorithm

Primal

Dual

$$\min_{\mathbf{Z} \in \mathcal{T}(\mathbf{r}, \mathbf{c})} [D_{KL}(\mathbf{Z} \parallel \mathbf{W})] = \max_{\boldsymbol{\alpha}, \boldsymbol{\beta}} [\langle \mathbf{r}, \boldsymbol{\alpha} \rangle + \langle \boldsymbol{\beta}, \mathbf{c} \rangle - \langle \exp(\boldsymbol{\alpha} \oplus \boldsymbol{\beta}), \mathbf{W} \rangle]$$

Rescaled Matrix

$\boldsymbol{\alpha}$ = Lagrange multipliers for the row-sum constraint

$\boldsymbol{\beta}$ = Lagrange multipliers for the column-sum constraint

(Schrödinger potentials)

$$\pi^* = \text{Diag}(\exp(\boldsymbol{\alpha}^*)) \mathbf{W} \text{Diag}(\exp(\boldsymbol{\beta}^*))$$

Alternating maximization

Matrix scaling algorithm
(Iterative Proportional Fitting)

Fit the row sums;
Fit the column sums; etc.

Sinkhorn Algorithm

$$\begin{cases} \forall 1 \leq i \leq n, \boldsymbol{\beta}_k(j) \leftarrow \log \frac{\mathbf{c}(j)}{\sum_{i=1}^m \mathbf{W}_{ij} \exp(\boldsymbol{\alpha}_{k-1}(i))}, \\ \forall 1 \leq i \leq m, \boldsymbol{\alpha}_k(i) \leftarrow \log \frac{\mathbf{r}(i)}{\sum_{j=1}^n \mathbf{W}_{ij} \exp(\boldsymbol{\beta}_k(j))}. \end{cases}$$

- ▶ Quick math for the convergence of Sinkhorn's Algorithm (Dual interpretation of Franklin and Lorenz '89)

Sinkhorn Algorithm

$$\begin{cases} \forall 1 \leq i \leq n, \beta_k(j) \leftarrow \log \frac{c(j)}{\sum_{i=1}^m w_{ij} \exp(\alpha_{k-1}(i))}, \\ \forall 1 \leq i \leq m, \alpha_k(i) \leftarrow \log \frac{r(i)}{\sum_{j=1}^n w_{ij} \exp(\beta_k(j))}. \end{cases}$$

Alternating maximization

$$= \begin{cases} \beta_k \leftarrow \arg \max_{\beta} f(\alpha_{k-1}, \beta) \\ \alpha_k \leftarrow \arg \max_{\alpha} f(\alpha, \beta_k) \end{cases}$$

- ▶ Quick math for the convergence of Sinkhorn's Algorithm (Dual interpretation of Franklin and Lorenz '89)

Sinkhorn Algorithm

$$\begin{cases} \forall 1 \leq i \leq n, \beta_k(j) \leftarrow \log \frac{c(j)}{\sum_{i=1}^m w_{ij} \exp(\alpha_{k-1}(i))}, \\ \forall 1 \leq i \leq m, \alpha_k(i) \leftarrow \log \frac{r(i)}{\sum_{j=1}^n w_{ij} \exp(\beta_k(j))}. \end{cases} =$$

Alternating maximization

$$\begin{cases} \beta_k \leftarrow \arg \max_{\beta} f(\alpha_{k-1}, \beta) \\ \alpha_k \leftarrow \arg \max_{\alpha} f(\alpha, \beta_k) \end{cases}$$

Step 1: ℓ_{∞} - non-expansion of the potentials

$$\|\alpha_k - \alpha^*\|_{\infty} = \|\xi(\beta_{k-1}) - \xi(\beta^*)\|_{\infty}$$

($\xi: \beta \mapsto \alpha$ update map)

- ▶ Quick math for the convergence of Sinkhorn's Algorithm (Dual interpretation of Franklin and Lorenz '89)

Sinkhorn Algorithm

Alternating maximization

$$\begin{cases} \forall 1 \leq i \leq n, \beta_k(j) \leftarrow \log \frac{c(j)}{\sum_{i=1}^m w_{ij} \exp(\alpha_{k-1}(i))}, \\ \forall 1 \leq i \leq m, \alpha_k(i) \leftarrow \log \frac{r(i)}{\sum_{j=1}^n w_{ij} \exp(\beta_k(j))}. \end{cases} = \begin{cases} \beta_k \leftarrow \arg \max_{\beta} f(\alpha_{k-1}, \beta) \\ \alpha_k \leftarrow \arg \max_{\alpha} f(\alpha, \beta_k) \end{cases}$$

Step 1: ℓ_{∞} - non-expansion of the potentials

$$\begin{aligned} \|\alpha_k - \alpha^*\|_{\infty} &= \|\xi(\beta_{k-1}) - \xi(\beta^*)\|_{\infty} && (\xi: \beta \mapsto \alpha \text{ update map}) \\ &= \left\| \int_0^1 J_{\xi; \beta}(\beta_{k-1}; s) ds (\beta_{k-1} - \beta^*) \right\|_{\infty} \end{aligned}$$

- ▶ Quick math for the convergence of Sinkhorn's Algorithm (Dual interpretation of Franklin and Lorenz '89)

Sinkhorn Algorithm

Alternating maximization

$$\begin{cases} \forall 1 \leq i \leq n, \boldsymbol{\beta}_k(j) \leftarrow \log \frac{\mathbf{c}(j)}{\sum_{i=1}^m \mathbf{W}_{ij} \exp(\boldsymbol{\alpha}_{k-1}(i))}, \\ \forall 1 \leq i \leq m, \boldsymbol{\alpha}_k(i) \leftarrow \log \frac{\mathbf{r}(i)}{\sum_{j=1}^n \mathbf{W}_{ij} \exp(\boldsymbol{\beta}_k(j))}. \end{cases} = \begin{cases} \boldsymbol{\beta}_k \leftarrow \arg \max_{\boldsymbol{\beta}} f(\boldsymbol{\alpha}_{k-1}, \boldsymbol{\beta}) \\ \boldsymbol{\alpha}_k \leftarrow \arg \max_{\boldsymbol{\alpha}} f(\boldsymbol{\alpha}, \boldsymbol{\beta}_k) \end{cases}$$

Step 1: ℓ_∞ - non-expansion of the potentials

$$\begin{aligned} \|\boldsymbol{\alpha}_k - \boldsymbol{\alpha}^*\|_\infty &= \|\xi(\boldsymbol{\beta}_{k-1}) - \xi(\boldsymbol{\beta}^*)\|_\infty && (\xi: \boldsymbol{\beta} \mapsto \boldsymbol{\alpha} \text{ update map}) \\ &= \left\| \int_0^1 J_{\xi; \boldsymbol{\beta}}(\boldsymbol{\beta}_{k-1}; s) ds (\boldsymbol{\beta}_{k-1} - \boldsymbol{\beta}^*) \right\|_\infty \\ &\leq \left(\int_0^1 \underbrace{\|J_{\xi; \boldsymbol{\beta}}(\boldsymbol{\beta}_{k-1}; s)\|_{\infty \rightarrow \infty}}_{=\text{row-stochastic}} ds \right) \|\boldsymbol{\beta}_{k-1} - \boldsymbol{\beta}^*\|_\infty \end{aligned}$$

- ▶ Quick math for the convergence of Sinkhorn's Algorithm (Dual interpretation of Franklin and Lorenz '89)

Sinkhorn Algorithm

Alternating maximization

$$\begin{cases} \forall 1 \leq i \leq n, \boldsymbol{\beta}_k(j) \leftarrow \log \frac{\mathbf{c}(j)}{\sum_{i=1}^m \mathbf{w}_{ij} \exp(\boldsymbol{\alpha}_{k-1}(i))}, \\ \forall 1 \leq i \leq m, \boldsymbol{\alpha}_k(i) \leftarrow \log \frac{\mathbf{r}(i)}{\sum_{j=1}^n \mathbf{w}_{ij} \exp(\boldsymbol{\beta}_k(j))}. \end{cases} = \begin{cases} \boldsymbol{\beta}_k \leftarrow \arg \max_{\boldsymbol{\beta}} f(\boldsymbol{\alpha}_{k-1}, \boldsymbol{\beta}) \\ \boldsymbol{\alpha}_k \leftarrow \arg \max_{\boldsymbol{\alpha}} f(\boldsymbol{\alpha}, \boldsymbol{\beta}_k) \end{cases}$$

Step 1: ℓ_∞ - non-expansion of the potentials

$$\begin{aligned} \|\boldsymbol{\alpha}_k - \boldsymbol{\alpha}^*\|_\infty &= \|\xi(\boldsymbol{\beta}_{k-1}) - \xi(\boldsymbol{\beta}^*)\|_\infty && (\xi: \boldsymbol{\beta} \mapsto \boldsymbol{\alpha} \text{ update map}) \\ &= \left\| \int_0^1 J_{\xi; \boldsymbol{\beta}}(\boldsymbol{\beta}_{k-1}; s) ds (\boldsymbol{\beta}_{k-1} - \boldsymbol{\beta}^*) \right\|_\infty \\ &\leq \left(\int_0^1 \underbrace{\|J_{\xi; \boldsymbol{\beta}}(\boldsymbol{\beta}_{k-1}; s)\|_{\infty \rightarrow \infty}}_{=\text{row-stochastic}} ds \right) \|\boldsymbol{\beta}_{k-1} - \boldsymbol{\beta}^*\|_\infty \\ &\leq \|\boldsymbol{\beta}_{k-1} - \boldsymbol{\beta}^*\|_\infty \end{aligned}$$

- ▶ Quick math for the convergence of Sinkhorn's Algorithm (Dual interpretation of Franklin and Lorenz '89)

Sinkhorn Algorithm

Alternating maximization

$$\begin{cases} \forall 1 \leq i \leq n, \boldsymbol{\beta}_k(j) \leftarrow \log \frac{\mathbf{c}(j)}{\sum_{i=1}^m \mathbf{w}_{ij} \exp(\boldsymbol{\alpha}_{k-1}(i))}, \\ \forall 1 \leq i \leq m, \boldsymbol{\alpha}_k(i) \leftarrow \log \frac{\mathbf{r}(i)}{\sum_{j=1}^n \mathbf{w}_{ij} \exp(\boldsymbol{\beta}_k(j))}. \end{cases} = \begin{cases} \boldsymbol{\beta}_k \leftarrow \arg \max_{\boldsymbol{\beta}} f(\boldsymbol{\alpha}_{k-1}, \boldsymbol{\beta}) \\ \boldsymbol{\alpha}_k \leftarrow \arg \max_{\boldsymbol{\alpha}} f(\boldsymbol{\alpha}, \boldsymbol{\beta}_k) \end{cases}$$

Step 1: ℓ_∞ - non-expansion of the potentials

$$\begin{aligned} \|\boldsymbol{\alpha}_k - \boldsymbol{\alpha}^*\|_\infty &= \|\xi(\boldsymbol{\beta}_{k-1}) - \xi(\boldsymbol{\beta}^*)\|_\infty && (\xi: \boldsymbol{\beta} \mapsto \boldsymbol{\alpha} \text{ update map}) \\ &= \left\| \int_0^1 J_{\xi; \boldsymbol{\beta}}(\boldsymbol{\beta}_{k-1}; s) ds (\boldsymbol{\beta}_{k-1} - \boldsymbol{\beta}^*) \right\|_\infty \\ &\leq \left(\int_0^1 \underbrace{\|J_{\xi; \boldsymbol{\beta}}(\boldsymbol{\beta}_{k-1}; s)\|_{\infty \rightarrow \infty}}_{=\text{row-stochastic}} ds \right) \|\boldsymbol{\beta}_{k-1} - \boldsymbol{\beta}^*\|_\infty \\ &\leq \|\boldsymbol{\beta}_{k-1} - \boldsymbol{\beta}^*\|_\infty \\ &\leq \dots \leq \|\boldsymbol{\alpha}_0 - \boldsymbol{\alpha}^*\|_\infty \end{aligned}$$

- ▶ Quick math for the convergence of Sinkhorn's Algorithm

Step 2: span-contraction of the potentials

$$\|\alpha_k - \alpha^*\|_{sp} = \|\xi(\beta_{k-1}) - \xi(\beta^*)\|_{sp} \quad (sp(x) := \max x - \min x)$$

► Quick math for the convergence of Sinkhorn's Algorithm

Step 2: span-contraction of the potentials

$$\|\alpha_k - \alpha^*\|_{sp} = \|\xi(\beta_{k-1}) - \xi(\beta^*)\|_{sp} \quad (sp(x) := \max x - \min x)$$

$$= \left\| \int_0^1 J_{\xi; \beta}(\beta_{k-1}; s) ds (\beta_{k-1} - \beta^*) \right\|_{sp}$$

$$\leq \left(\int_0^1 \left\| \underbrace{J_{\xi; \beta}(\beta_{k-1}; s)}_{\substack{=\text{row-stochastic, nonzero} \\ \leq \tanh(\Delta(J_{\xi; \beta}(\beta_{k-1}; s)))/2 < 1}} \right\|_{sp \rightarrow sp} ds \right) \|\beta_{k-1} - \beta^*\|_{sp}$$

(Birkoff's contraction Thm.)

$$\Delta(P) := \max_{i, i', j, j'} \log \frac{P_{ij} P_{i'j'}}{P_{ij'} P_{i'j}}$$

► Quick math for the convergence of Sinkhorn's Algorithm

Step 2: span-contraction of the potentials

$$\begin{aligned}
 \|\alpha_k - \alpha^*\|_{sp} &= \|\xi(\beta_{k-1}) - \xi(\beta^*)\|_{sp} && (sp(x) := \max x - \min x) \\
 &= \left\| \int_0^1 J_{\xi; \beta}(\beta_{k-1}; s) ds (\beta_{k-1} - \beta^*) \right\|_{sp} \\
 &\leq \left(\int_0^1 \underbrace{\left\| \underbrace{J_{\xi; \beta}(\beta_{k-1}; s)}_{\substack{=\text{row-stochastic, nonzero} \\ \leq \tanh(\Delta(J_{\xi; \beta}(\beta_{k-1}; s)))/2 < 1}} \right\|_{sp \rightarrow sp} ds}_{\text{Birkoff's contraction Thm.}} \right) \|\beta_{k-1} - \beta^*\|_{sp} \\
 &\leq \rho(\text{range of potentials}) \|\beta_{k-1} - \beta^*\|_{sp} && \Delta(P) := \max_{i, i', j, j'} \log \frac{P_{ij} P_{i'j'}}{P_{ij'} P_{i'j}}
 \end{aligned}$$

► Quick math for the convergence of Sinkhorn's Algorithm

Step 2: span-contraction of the potentials

$$\begin{aligned}
 \|\alpha_k - \alpha^*\|_{sp} &= \|\xi(\beta_{k-1}) - \xi(\beta^*)\|_{sp} && (sp(x) := \max x - \min x) \\
 &= \left\| \int_0^1 J_{\xi; \beta}(\beta_{k-1}; s) ds (\beta_{k-1} - \beta^*) \right\|_{sp} \\
 &\leq \left(\int_0^1 \underbrace{\left\| \underbrace{J_{\xi; \beta}(\beta_{k-1}; s)}_{\substack{=\text{row-stochastic, nonzero} \\ \leq \tanh(\Delta(J_{\xi; \beta}(\beta_{k-1}; s)))/2 < 1}} \right\|_{sp \rightarrow sp} ds}_{\text{(Birkoff's contraction Thm.)}} \right) \|\beta_{k-1} - \beta^*\|_{sp} \\
 &\leq \rho(\text{range of potentials}) \|\beta_{k-1} - \beta^*\|_{sp} && \Delta(P) := \max_{i, i', j, j'} \log \frac{P_{ij} P_{i'j'}}{P_{ij'} P_{i'j}} \\
 &\leq \dots \leq \rho(\text{range of potentials})^k \|\alpha_0 - \alpha^*\|_{sp}
 \end{aligned}$$

► Monge-Kantorovich Optimal Transport

Given a margin (\mathbf{r}, \mathbf{c}) and cost matrix $\mathbf{C} = (c_{ij})$, c_{ij} = cost of moving a unit mass from location i to j , find the most efficient coupling π^* :

$$\pi^* = \operatorname{argmin}_{\pi \in \Pi(\mathbf{r}, \mathbf{c})} \langle \pi, \mathbf{C} \rangle$$

► Monge-Kantorovich Optimal Transport

Given a margin (\mathbf{r}, \mathbf{c}) and cost matrix $\mathbf{C} = (c_{ij})$, c_{ij} = cost of moving a unit mass from location i to j , find the most efficient coupling π^* :

$$\pi^* = \operatorname{argmin}_{\pi \in \Pi(\mathbf{r}, \mathbf{c})} \langle \pi, \mathbf{C} \rangle$$

► Entropic Optimal Transport

$$\pi^* = \operatorname{argmin}_{\pi \in \Pi(\mathbf{r}, \mathbf{c})} \langle \pi, \mathbf{C} \rangle + \varepsilon \operatorname{KL}(\pi \| \mathbf{r} \otimes \mathbf{c})$$

- Entropic Regularization (Cuturi, NeurIPS '13)

► Monge-Kantorovich Optimal Transport

Given a margin (\mathbf{r}, \mathbf{c}) and cost matrix $\mathbf{C} = (c_{ij})$, c_{ij} = cost of moving a unit mass from location i to j , find the most efficient coupling π^* :

$$\pi^* = \operatorname{argmin}_{\pi \in \Pi(\mathbf{r}, \mathbf{c})} \langle \pi, \mathbf{C} \rangle$$

► Entropic Optimal Transport

- Entropic Regularization (Cuturi, NeurIPS '13)

$$\pi^* = \operatorname{argmin}_{\pi \in \Pi(\mathbf{r}, \mathbf{c})} \langle \pi, \mathbf{C} \rangle + \varepsilon \operatorname{KL}(\pi \| \mathbf{r} \otimes \mathbf{c})$$

$$= \operatorname{argmin}_{\pi \in \Pi(\mathbf{r}, \mathbf{c})} [D_{KL}(\pi \| \mathbf{W}) \text{ with } \mathbf{W} = \exp(-\mathbf{C}/\varepsilon) \odot (\mathbf{r} \otimes \mathbf{c})]$$

(Matrix Scaling Problem, so solve it by Sinkhorn!)

► Wasserstein distances

μ, ν = Probability measures on spaces X and Y

$\kappa: X \times Y \rightarrow [0, \infty]$ cost function

$$W_{\kappa}(\mu, \nu) := \min_{\pi \in \Pi(\mu, \nu)} \left[\langle \pi, \kappa \rangle = \int \kappa(x, y) d\pi(x, y) \right]$$

► Wasserstein distances

μ, ν = Probability measures on spaces X and Y

$\kappa: X \times Y \rightarrow [0, \infty]$ cost function

$$W_{\kappa}(\mu, \nu) := \min_{\pi \in \Pi(\mu, \nu)} \left[\langle \pi, \kappa \rangle = \int \kappa(x, y) d\pi(x, y) \right]$$

Cons: Nonsmooth and high sample complexity ($n^{-1/\text{dim}}$), LP computation

→ Not suitable for ML uses

► Wasserstein distances

μ, ν = Probability measures on spaces X and Y

$\kappa: X \times Y \rightarrow [0, \infty]$ cost function

$$W_{\kappa}(\mu, \nu) := \min_{\pi \in \Pi(\mu, \nu)} \left[\langle \pi, \kappa \rangle = \int \kappa(x, y) d\pi(x, y) \right]$$

Cons: Nonsmooth and high sample complexity ($n^{-1/\text{dim}}$), LP computation

→ Not suitable for ML uses

► Sinkhorn Divergence (Cuturi' 13)

$$S_{\varepsilon}(\mu, \nu) = \text{OT}_{\varepsilon}(\mu, \nu) - \frac{1}{2} \text{OT}_{\varepsilon}(\mu, \mu) - \frac{1}{2} \text{OT}_{\varepsilon}(\nu, \nu).$$

$$\text{OT}_{\kappa}(\mu, \nu) := \min_{\pi \in \Pi(\mu, \nu)} \left[\langle \pi, \kappa \rangle + \varepsilon D_{KL}(\pi, \mu \otimes \nu) \right] \quad (\text{Computed by Empirical Sinkhorn})$$

▶ Wasserstein distances

μ, ν = Probability measures on spaces X and Y

$\kappa: X \times Y \rightarrow [0, \infty]$ cost function

$$W_{\kappa}(\mu, \nu) := \min_{\pi \in \Pi(\mu, \nu)} \left[\langle \pi, \kappa \rangle = \int \kappa(x, y) d\pi(x, y) \right]$$

Cons: Nonsmooth and high sample complexity ($n^{-1/\text{dim}}$), LP computation

→ Not suitable for ML uses

▶ Sinkhorn Divergence (Cuturi' 13)

$$S_{\varepsilon}(\mu, \nu) = \text{OT}_{\varepsilon}(\mu, \nu) - \frac{1}{2} \text{OT}_{\varepsilon}(\mu, \mu) - \frac{1}{2} \text{OT}_{\varepsilon}(\nu, \nu).$$

$$\text{OT}_{\kappa}(\mu, \nu) := \min_{\pi \in \Pi(\mu, \nu)} \left[\langle \pi, \kappa \rangle + \varepsilon D_{KL}(\pi, \mu \otimes \nu) \right] \quad (\text{Computed by Empirical Sinkhorn})$$

Pros: smooth, low sample complexity ($n^{-1/2}$), easily computed, $W_{\kappa}(\mu, \nu) = \lim_{\varepsilon \rightarrow 0} S_{\kappa}(\mu, \nu)$

-
- Cuturi, *Sinkhorn Distances: Lightspeed Computation of OT*, NeurIPS 2013.
 - Genevay, Peyré, Cuturi, *Learning Generative Models with Sinkhorn Divergences*, AISTATS 2018.
 - Feydy, Sejourn, Vialard, Amari, Trouv, Peyr, *Interpolating between OT and MMD using Sinkhorn divergences*, AISTATS 2019.
 - Genevay, Chizat, Bach, Cuturi, Peyr, *Sample Complexity of Sinkhorn divergences*, AISTATS 2019.
 - Peyr & Cuturi, *Computational Optimal Transport*, FnT in ML, 2019 (the standard reference book).

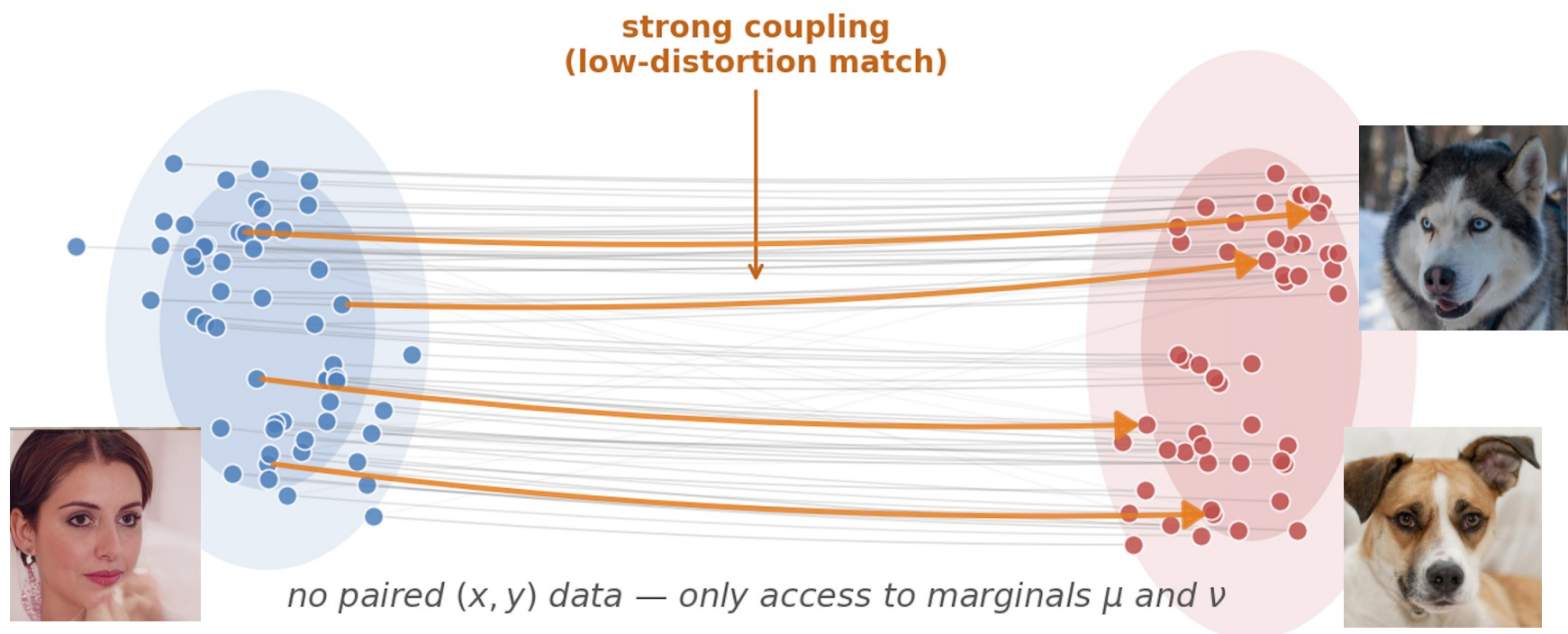
ML Applications on Generative modeling

Image translation
Knowledge distillation

► Generative modeling: Unpaired Image Translation

μ : human face images

ν : dog images



► Generative modeling: Unpaired Image Translation



Results from REU 2025 @ UW-Madison on generative modeling

Dog ↔ Human image translation

► Generative modeling: Unpaired Image Translation

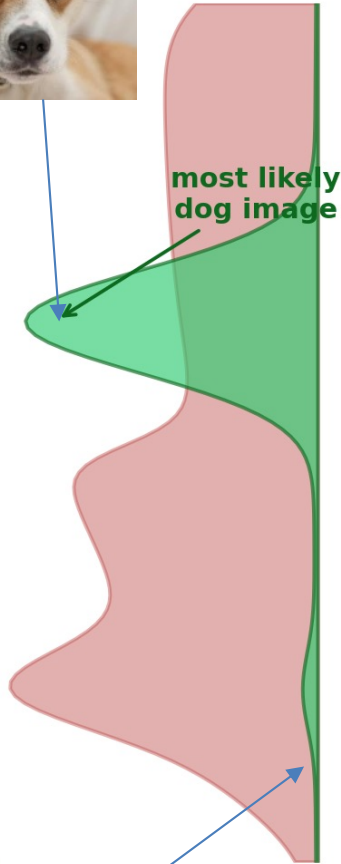


Results from REU 2025 @ UW-Madison on generative modeling
Adult ↔ Children image translation

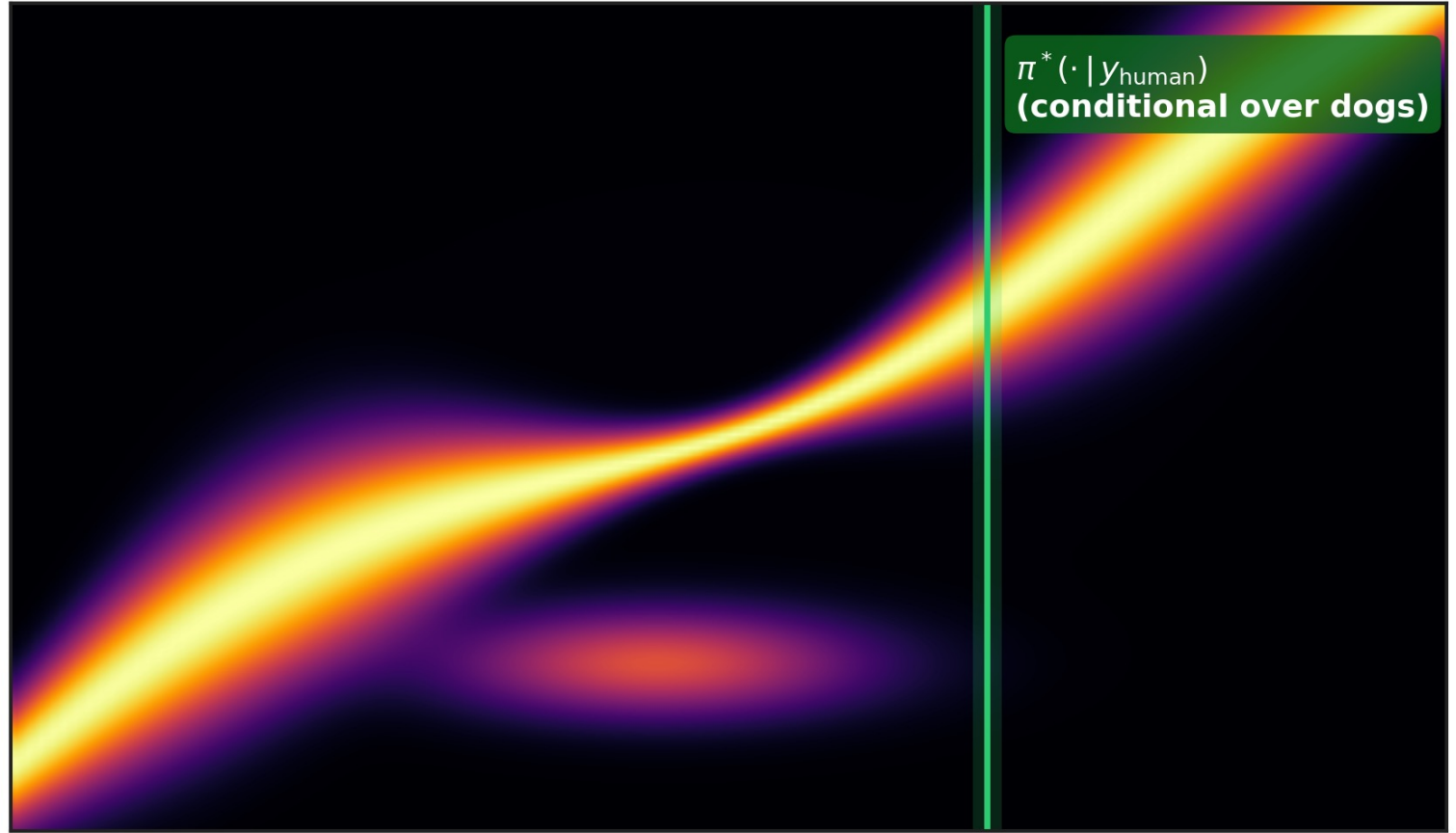
Entropic transport plan $\pi^*(\text{dog}, \text{human})$ — concentrates on low-distortion matches



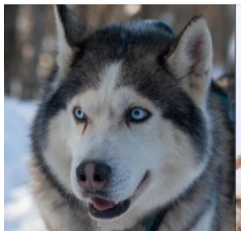
μ : dog image distribution



most likely dog image



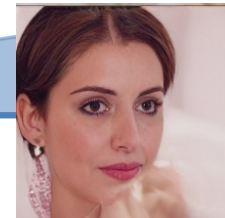
$\pi^*(\cdot | y_{\text{human}})$
(conditional over dogs)



fixed human face y_{human}

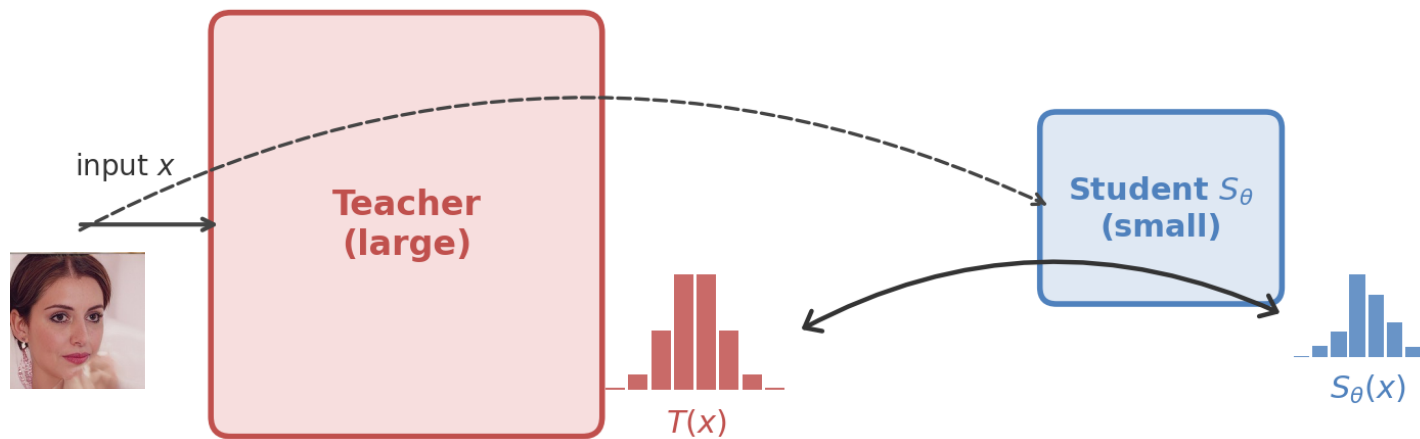


ν : human face image distribution



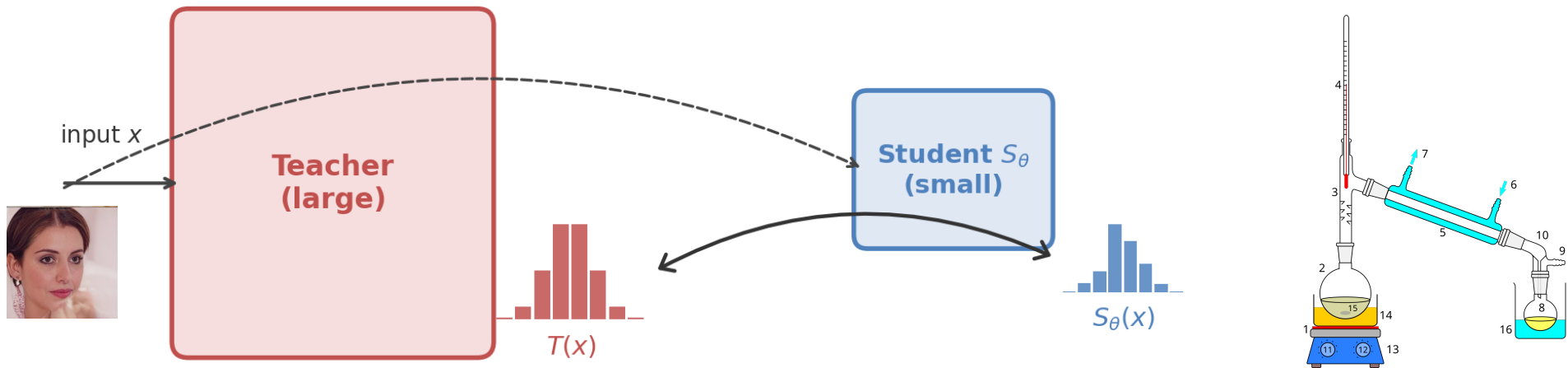
► Generative modeling: Knowledge Distillation

- Compress a large teacher T into a small student S_θ while preserving behavior



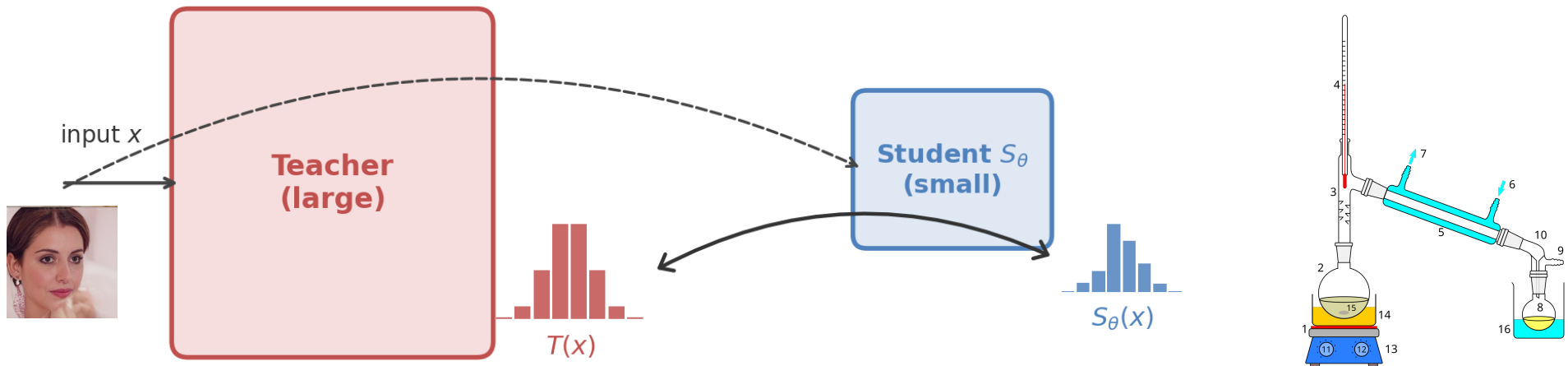
▶ Generative modeling: Knowledge Distillation

- Compress a large teacher T into a small student S_θ while preserving behavior



▶ Generative modeling: Knowledge Distillation

- Compress a large teacher T into a small student S_θ while preserving behavior



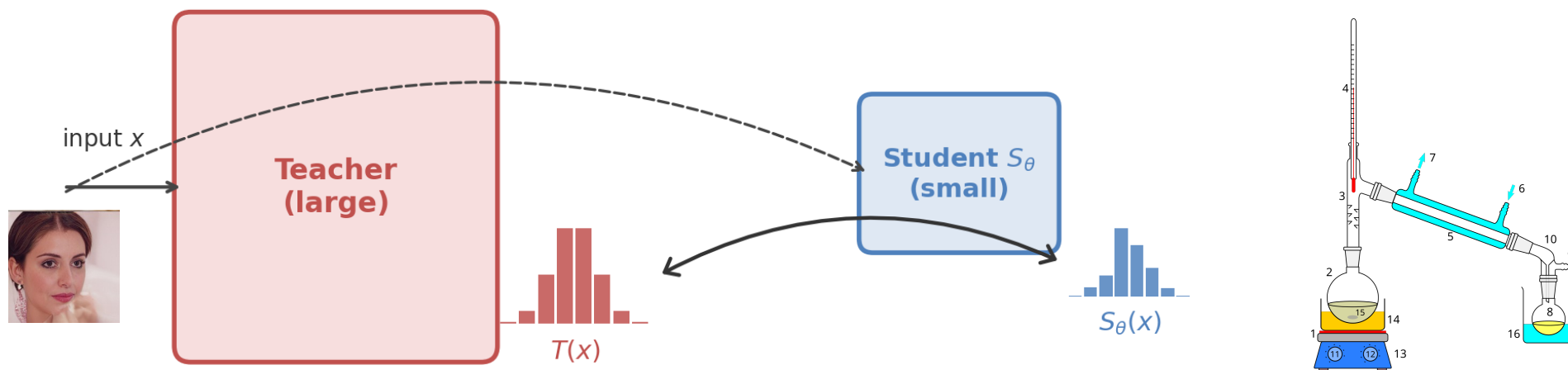
Distillation is being used to make capable small LLM from large models

Refs: Cui–Wu–Chen '22 (Wasserstein KD); Park–Kim–Lu '19 (Relational KD); Bhardwaj et al. '22 (OT for LLM distillation)

- **Llama 3.1 → Llama 3.2** — the 1B and 3B were distilled from the 8B/70B.

▶ Generative modeling: Knowledge Distillation

- Compress a large teacher T into a small student S_θ while preserving behavior



Distillation is being used to make capable small LLM from large models

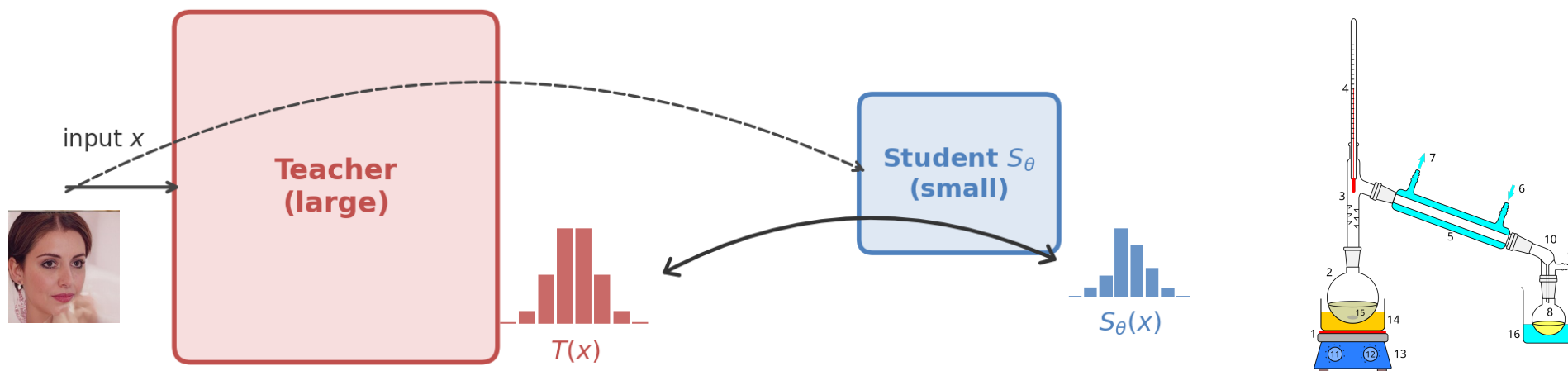
Refs: Cui–Wu–Chen '22 (Wasserstein KD); Park–Kim–Lu '19 (Relational KD); Bhardwaj et al. '22 (OT for LLM distillation)

• **Llama 3.1 → Llama 3.2** — the 1B and 3B were distilled from the 8B/70B.

- Classical (Hinton '15): Find θ by minimizing $D_{KL}(S_\theta(x) \parallel T(x))$

▶ Generative modeling: Knowledge Distillation

- Compress a large teacher T into a small student S_θ while preserving behavior



Distillation is being used to make capable small LLM from large models

Refs: Cui–Wu–Chen '22 (Wasserstein KD); Park–Kim–Lu '19 (Relational KD); Bhardwaj et al. '22 (OT for LLM distillation)

• **Llama 3.1 → Llama 3.2** — the 1B and 3B were distilled from the 8B/70B.

- Classical (Hinton '15): Find θ by minimizing $D_{KL}(S_\theta(x) \parallel T(x))$

EOT-based KD: Sinkhorn divergence handles disjoint supports; KL blows up

$$\mathcal{L}_{\text{KD}}(\theta) = \text{Sinkhorn}_\varepsilon(S_\theta(x), T(x); C)$$

New Perspective I: Random Coordinate Sinkhorn

► Motivating Question: EOT on large graph $G=(V,E)$

$G = (V, E)$ a large graph, μ, ν two probability distributions on the node set V .

Find the most efficient coupling π^* in terms of the shortest-path distance on G

$$\pi^* = \operatorname{argmin}_{\pi \in \Pi(\mu, \nu)} \langle \pi, d_G \rangle + \varepsilon D_{\text{KL}}(\pi \| \mu \otimes \nu)$$

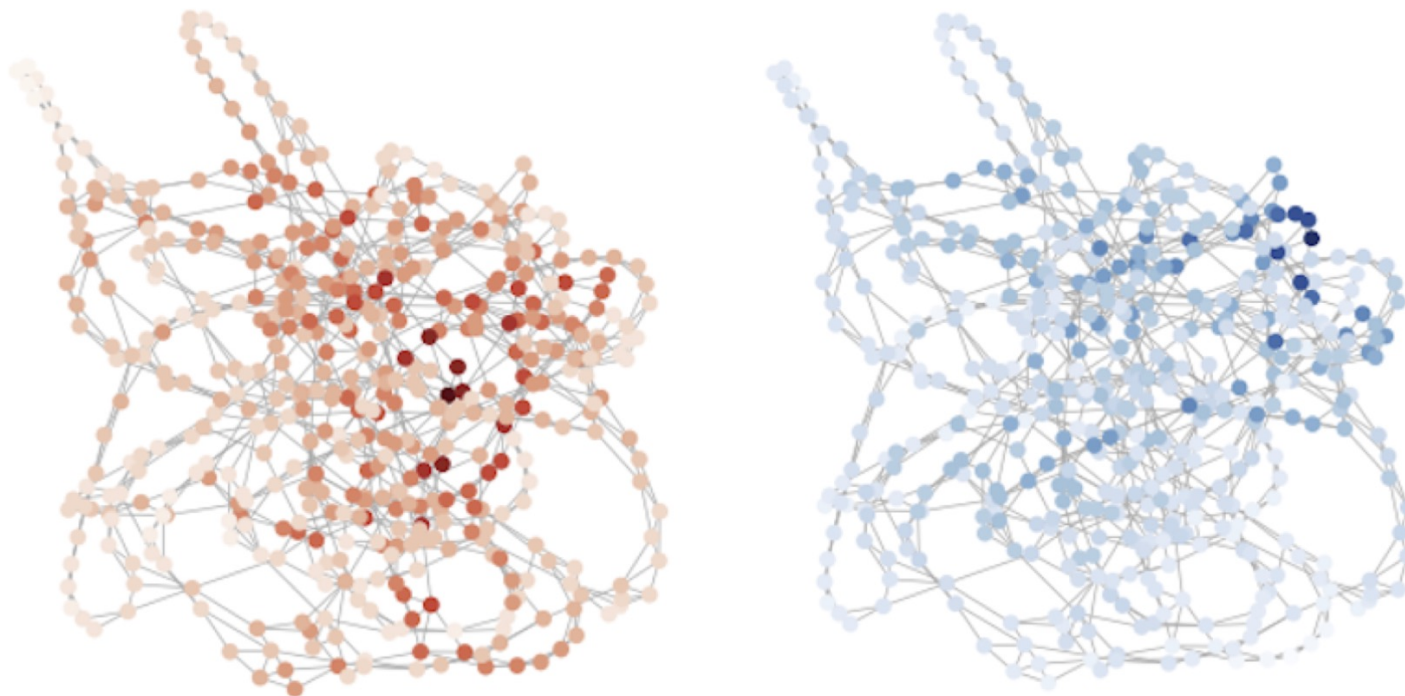


FIGURE 1. A 500-node graph G and two probability distributions μ and ν on the node set shown. Color shade indicates mass at each node.

► Motivating Question: EOT on large graph $G=(V,E)$

$G = (V, E)$ a large graph, μ, ν two probability distributions on the node set V .

Find the most efficient coupling π^* in terms of the shortest-path distance on G

$$\pi^* = \operatorname{argmin}_{\pi \in \Pi(\mu, \nu)} \langle \pi, d_G \rangle + \varepsilon D_{\text{KL}}(\pi \| \mu \otimes \nu)$$

| Shortest-path distance matrix

For large $G = (V, E)$, even evaluating the objective is challenging

Computational Cost: The cost matrix d_G requires an All-Pairs Shortest Path (APSP) computation = $O(n^3)$ or $O(n^2 \log n + n|E|)$

Memory Cost: The matrix d_G is dense, requiring $O(n^2)$ storage, which is prohibitive for graphs with millions of nodes.

- ▶ Open Question: EOT on large graph $G=(V,E)$

$$(\boldsymbol{\alpha}^*, \boldsymbol{\beta}^*) \in \operatorname{argmax}_{\boldsymbol{\alpha}, \boldsymbol{\beta}: V \rightarrow \mathbb{R}} f_G(\boldsymbol{\alpha}, \boldsymbol{\beta}) := \langle \boldsymbol{\alpha}, \boldsymbol{\mu} \rangle + \langle \boldsymbol{\beta}, \boldsymbol{\nu} \rangle - \langle \exp(\boldsymbol{\alpha} \oplus \boldsymbol{\beta}), K_G \odot (\boldsymbol{\mu} \otimes \boldsymbol{\nu}) \rangle.$$

$$K_G(x, y) := \exp(-d_G(x, y)/\varepsilon)$$

- ▶ Open Question: EOT on large graph $G=(V,E)$

$$(\boldsymbol{\alpha}^*, \boldsymbol{\beta}^*) \in \operatorname{argmax}_{\boldsymbol{\alpha}, \boldsymbol{\beta}: V \rightarrow \mathbb{R}} f_G(\boldsymbol{\alpha}, \boldsymbol{\beta}) := \langle \boldsymbol{\alpha}, \boldsymbol{\mu} \rangle + \langle \boldsymbol{\beta}, \boldsymbol{\nu} \rangle - \langle \exp(\boldsymbol{\alpha} \oplus \boldsymbol{\beta}), K_G \odot (\boldsymbol{\mu} \otimes \boldsymbol{\nu}) \rangle.$$

$$K_G(x, y) := \exp(-d_G(x, y)/\varepsilon)$$

Main Idea: Solve an easier stochastic surrogate

$$\begin{aligned} (\hat{\boldsymbol{\alpha}}, \hat{\boldsymbol{\beta}}) \in \operatorname{argmax}_{\boldsymbol{\alpha}, \boldsymbol{\beta}: V \rightarrow \mathbb{R}} (\hat{f}_G(\boldsymbol{\alpha}, \boldsymbol{\beta}) &:= \mathbb{E}_{H \sim \pi} [\hat{f}_H(\boldsymbol{\alpha}, \boldsymbol{\beta})]) \\ &= \langle \boldsymbol{\alpha}, \boldsymbol{\mu} \rangle + \langle \boldsymbol{\beta}, \boldsymbol{\nu} \rangle - \langle \exp(\boldsymbol{\alpha} \oplus \boldsymbol{\beta}), \hat{K} \odot (\boldsymbol{\mu} \otimes \boldsymbol{\nu}) \rangle. \end{aligned}$$

where \hat{f}_H = local dual objective on a random subgraph H :

- ▶ Open Question: EOT on large graph $G=(V,E)$

$$(\boldsymbol{\alpha}^*, \boldsymbol{\beta}^*) \in \operatorname{argmax}_{\boldsymbol{\alpha}, \boldsymbol{\beta}: V \rightarrow \mathbb{R}} f_G(\boldsymbol{\alpha}, \boldsymbol{\beta}) := \langle \boldsymbol{\alpha}, \boldsymbol{\mu} \rangle + \langle \boldsymbol{\beta}, \boldsymbol{\nu} \rangle - \langle \exp(\boldsymbol{\alpha} \oplus \boldsymbol{\beta}), K_G \odot (\boldsymbol{\mu} \otimes \boldsymbol{\nu}) \rangle.$$

$$K_G(x, y) := \exp(-d_G(x, y)/\varepsilon)$$

Main Idea: Solve an easier stochastic surrogate

$$\begin{aligned} (\hat{\boldsymbol{\alpha}}, \hat{\boldsymbol{\beta}}) \in \operatorname{argmax}_{\boldsymbol{\alpha}, \boldsymbol{\beta}: V \rightarrow \mathbb{R}} (\hat{f}_G(\boldsymbol{\alpha}, \boldsymbol{\beta}) &:= \mathbb{E}_{H \sim \pi} [\hat{f}_H(\boldsymbol{\alpha}, \boldsymbol{\beta})]) \\ &= \langle \boldsymbol{\alpha}, \boldsymbol{\mu} \rangle + \langle \boldsymbol{\beta}, \boldsymbol{\nu} \rangle - \langle \exp(\boldsymbol{\alpha} \oplus \boldsymbol{\beta}), \hat{K} \odot (\boldsymbol{\mu} \otimes \boldsymbol{\nu}) \rangle. \end{aligned}$$

where \hat{f}_H = local dual objective on a random subgraph H :

$$\hat{f}_H(\boldsymbol{\alpha}, \boldsymbol{\beta}) = \sum_{x \in V(H)} \boldsymbol{\alpha}(x) \frac{\mu(x)}{p_x} + \sum_{y \in V(H)} \boldsymbol{\beta}(y) \frac{\nu(y)}{p_y} - \sum_{x, y \in V(H)} e^{\boldsymbol{\alpha}(x) + \boldsymbol{\beta}(y)} K_H(x, y) \frac{\mu(x)\nu(y)}{p_{xy}}$$

\hat{K} = Expected surrogate Gibbs kernel

$$\hat{K}(x, y) = \mathbb{E}_H [K_H(x, y) \mid x, y \in V(H)]$$

$$K_H(x, y) = \exp(-d_H(x, y)/\varepsilon) \quad \text{for } x, y \in V(H)$$

Local distance matrix on H

- ▶ Open Question: EOT on large graph $G=(V,E)$

$$(\boldsymbol{\alpha}^*, \boldsymbol{\beta}^*) \in \operatorname{argmax}_{\boldsymbol{\alpha}, \boldsymbol{\beta}: V \rightarrow \mathbb{R}} f_G(\boldsymbol{\alpha}, \boldsymbol{\beta}) := \langle \boldsymbol{\alpha}, \boldsymbol{\mu} \rangle + \langle \boldsymbol{\beta}, \boldsymbol{\nu} \rangle - \langle \exp(\boldsymbol{\alpha} \oplus \boldsymbol{\beta}), \boxed{K_G} \odot (\boldsymbol{\mu} \otimes \boldsymbol{\nu}) \rangle.$$

$$K_G(x, y) := \exp(-d_G(x, y)/\varepsilon)$$

Main Idea: Solve an easier stochastic surrogate

$$\begin{aligned} (\hat{\boldsymbol{\alpha}}, \hat{\boldsymbol{\beta}}) \in \operatorname{argmax}_{\boldsymbol{\alpha}, \boldsymbol{\beta}: V \rightarrow \mathbb{R}} (\hat{f}_G(\boldsymbol{\alpha}, \boldsymbol{\beta}) &:= \mathbb{E}_{H \sim \pi} [\hat{f}_H(\boldsymbol{\alpha}, \boldsymbol{\beta})]) \\ &= \langle \boldsymbol{\alpha}, \boldsymbol{\mu} \rangle + \langle \boldsymbol{\beta}, \boldsymbol{\nu} \rangle - \langle \exp(\boldsymbol{\alpha} \oplus \boldsymbol{\beta}), \boxed{\hat{K}} \odot (\boldsymbol{\mu} \otimes \boldsymbol{\nu}) \rangle. \end{aligned}$$

Q1: How do we control the stochastic relaxation bias $f_G - \hat{f}_G$?

- ▶ Open Question: EOT on large graph $G=(V,E)$

$$(\boldsymbol{\alpha}^*, \boldsymbol{\beta}^*) \in \operatorname{argmax}_{\boldsymbol{\alpha}, \boldsymbol{\beta}: V \rightarrow \mathbb{R}} f_G(\boldsymbol{\alpha}, \boldsymbol{\beta}) := \langle \boldsymbol{\alpha}, \boldsymbol{\mu} \rangle + \langle \boldsymbol{\beta}, \boldsymbol{\nu} \rangle - \langle \exp(\boldsymbol{\alpha} \oplus \boldsymbol{\beta}), K_G \odot (\boldsymbol{\mu} \otimes \boldsymbol{\nu}) \rangle.$$

$$K_G(x, y) := \exp(-d_G(x, y)/\varepsilon)$$

Main Idea: Solve an easier stochastic surrogate

$$\begin{aligned} (\hat{\boldsymbol{\alpha}}, \hat{\boldsymbol{\beta}}) \in \operatorname{argmax}_{\boldsymbol{\alpha}, \boldsymbol{\beta}: V \rightarrow \mathbb{R}} (\hat{f}_G(\boldsymbol{\alpha}, \boldsymbol{\beta}) &:= \mathbb{E}_{H \sim \pi} [\hat{f}_H(\boldsymbol{\alpha}, \boldsymbol{\beta})]) \\ &= \langle \boldsymbol{\alpha}, \boldsymbol{\mu} \rangle + \langle \boldsymbol{\beta}, \boldsymbol{\nu} \rangle - \langle \exp(\boldsymbol{\alpha} \oplus \boldsymbol{\beta}), \hat{K} \odot (\boldsymbol{\mu} \otimes \boldsymbol{\nu}) \rangle. \end{aligned}$$

Q1: How do we control the stochastic relaxation bias $f_G - \hat{f}_G$?

Q2: How do we compute the stochastic surrogate potentials $(\hat{\boldsymbol{\alpha}}, \hat{\boldsymbol{\beta}})$?

- Q1: How do we control the stochastic relaxation bias $f_G - \hat{f}_G$?

$$(\boldsymbol{\alpha}^*, \boldsymbol{\beta}^*) \in \operatorname{argmax}_{\boldsymbol{\alpha}, \boldsymbol{\beta}: V \rightarrow \mathbb{R}} f_G(\boldsymbol{\alpha}, \boldsymbol{\beta}) := \langle \boldsymbol{\alpha}, \boldsymbol{\mu} \rangle + \langle \boldsymbol{\beta}, \boldsymbol{\nu} \rangle - \langle \exp(\boldsymbol{\alpha} \oplus \boldsymbol{\beta}), K_G \odot (\boldsymbol{\mu} \otimes \boldsymbol{\nu}) \rangle$$

$$(\hat{\boldsymbol{\alpha}}, \hat{\boldsymbol{\beta}}) \in \operatorname{argmax}_{\boldsymbol{\alpha}, \boldsymbol{\beta}: V \rightarrow \mathbb{R}} (\hat{f}_G(\boldsymbol{\alpha}, \boldsymbol{\beta}) = \langle \boldsymbol{\alpha}, \boldsymbol{\mu} \rangle + \langle \boldsymbol{\beta}, \boldsymbol{\nu} \rangle - \langle \exp(\boldsymbol{\alpha} \oplus \boldsymbol{\beta}), \hat{K} \odot (\boldsymbol{\mu} \otimes \boldsymbol{\nu}) \rangle)$$

► **Q1: How do we control the stochastic relaxation bias $f_G - \hat{f}_G$?**

$$(\boldsymbol{\alpha}^*, \boldsymbol{\beta}^*) \in \operatorname{argmax}_{\boldsymbol{\alpha}, \boldsymbol{\beta}: V \rightarrow \mathbb{R}} f_G(\boldsymbol{\alpha}, \boldsymbol{\beta}) := \langle \boldsymbol{\alpha}, \boldsymbol{\mu} \rangle + \langle \boldsymbol{\beta}, \boldsymbol{\nu} \rangle - \langle \exp(\boldsymbol{\alpha} \oplus \boldsymbol{\beta}), K_G \odot (\boldsymbol{\mu} \otimes \boldsymbol{\nu}) \rangle$$

$$(\hat{\boldsymbol{\alpha}}, \hat{\boldsymbol{\beta}}) \in \operatorname{argmax}_{\boldsymbol{\alpha}, \boldsymbol{\beta}: V \rightarrow \mathbb{R}} (\hat{f}_G(\boldsymbol{\alpha}, \boldsymbol{\beta}) = \langle \boldsymbol{\alpha}, \boldsymbol{\mu} \rangle + \langle \boldsymbol{\beta}, \boldsymbol{\nu} \rangle - \langle \exp(\boldsymbol{\alpha} \oplus \boldsymbol{\beta}), \hat{K} \odot (\boldsymbol{\mu} \otimes \boldsymbol{\nu}) \rangle)$$

Thm. (Bi, Kim, L. 26+) Let $\pi^* = e^{\boldsymbol{\alpha} \oplus \boldsymbol{\beta}} \odot (\boldsymbol{\mu} \otimes \boldsymbol{\nu})$, $\hat{\pi}^* = e^{\hat{\boldsymbol{\alpha}} \oplus \hat{\boldsymbol{\beta}}} \odot (\boldsymbol{\mu} \otimes \boldsymbol{\nu})$. Then

$$d_{\text{Hell}}(\pi^*, \hat{\pi}^*) \leq 2e^{\kappa_{\max}} d_{\text{Hell}; \boldsymbol{\mu} \otimes \boldsymbol{\nu}}(K_G, \hat{K})$$

where $\kappa_{\max} = \max_{x, y} K_G(x, y) \vee \hat{K}(x, y)$, $d_{\text{Hell}; \boldsymbol{\mu} \otimes \boldsymbol{\nu}}(K, \tilde{K}) = \left(\sum_{i, j} (\sqrt{K_{ij}} - \sqrt{\tilde{K}_{ij}})^2 \mu(i) \nu(j) \right)^{1/2}$

► **Q1: How do we control the stochastic relaxation bias $f_G - \hat{f}_G$?**

$$(\boldsymbol{\alpha}^*, \boldsymbol{\beta}^*) \in \operatorname{argmax}_{\boldsymbol{\alpha}, \boldsymbol{\beta}: V \rightarrow \mathbb{R}} f_G(\boldsymbol{\alpha}, \boldsymbol{\beta}) := \langle \boldsymbol{\alpha}, \boldsymbol{\mu} \rangle + \langle \boldsymbol{\beta}, \boldsymbol{\nu} \rangle - \langle \exp(\boldsymbol{\alpha} \oplus \boldsymbol{\beta}), K_G \odot (\boldsymbol{\mu} \otimes \boldsymbol{\nu}) \rangle$$

$$(\hat{\boldsymbol{\alpha}}, \hat{\boldsymbol{\beta}}) \in \operatorname{argmax}_{\boldsymbol{\alpha}, \boldsymbol{\beta}: V \rightarrow \mathbb{R}} (\hat{f}_G(\boldsymbol{\alpha}, \boldsymbol{\beta}) = \langle \boldsymbol{\alpha}, \boldsymbol{\mu} \rangle + \langle \boldsymbol{\beta}, \boldsymbol{\nu} \rangle - \langle \exp(\boldsymbol{\alpha} \oplus \boldsymbol{\beta}), \hat{K} \odot (\boldsymbol{\mu} \otimes \boldsymbol{\nu}) \rangle)$$

Thm. (Bi, Kim, L. 26+) Let $\pi^* = e^{\boldsymbol{\alpha} \oplus \boldsymbol{\beta}} \odot (\boldsymbol{\mu} \otimes \boldsymbol{\nu})$, $\hat{\pi}^* = e^{\hat{\boldsymbol{\alpha}} \oplus \hat{\boldsymbol{\beta}}} \odot (\boldsymbol{\mu} \otimes \boldsymbol{\nu})$. Then

$$d_{\text{Hell}}(\pi^*, \hat{\pi}^*) \leq 2e^{\kappa_{\max}} d_{\text{Hell}; \boldsymbol{\mu} \otimes \boldsymbol{\nu}}(K_G, \hat{K})$$

where $\kappa_{\max} = \max_{x, y} K_G(x, y) \vee \hat{K}(x, y)$, $d_{\text{Hell}; \boldsymbol{\mu} \otimes \boldsymbol{\nu}}(K, \tilde{K}) = \left(\sum_{ij} (\sqrt{K_{ij}} - \sqrt{\tilde{K}_{ij}})^2 \mu(i) \nu(j) \right)^{1/2}$

Follows directly from a stability bound on matrix scaling w.r.t. input matrix due to William Powell and Danny Duan (26+)

$d_{\text{Hell}; \boldsymbol{\mu} \otimes \boldsymbol{\nu}}(K_G, \hat{K})$ depends on G and the random subgraph H

e.g. $H = k$ -ball neighborhood of a random node

- Q2: How do we compute the stochastic surrogate potentials $(\hat{\alpha}, \hat{\beta})$?

$$\begin{aligned}(\hat{\alpha}, \hat{\beta}) &\in \arg \max_{\alpha, \beta: V \rightarrow \mathbb{R}} (\hat{f}_G(\alpha, \beta) := \mathbb{E}_{H \sim \pi} [\hat{f}_H(\alpha, \beta)]) \\ &= \langle \alpha, \mu \rangle + \langle \beta, \nu \rangle - \langle \exp(\alpha \oplus \beta), \hat{K} \odot (\mu \otimes \nu) \rangle.\end{aligned}$$

Alg. (Full Sinkhorn)

$$\begin{cases} \alpha_k \leftarrow \arg \max_{\alpha} f(\alpha, \beta_{k-1}) \\ \beta_k \leftarrow \arg \max_{\beta} f(\alpha_k, \beta) \end{cases}$$

- **Q2: How do we compute the stochastic surrogate potentials $(\hat{\alpha}, \hat{\beta})$?**

$$\begin{aligned}
 (\hat{\alpha}, \hat{\beta}) &\in \operatorname{argmax}_{\alpha, \beta: V \rightarrow \mathbb{R}} (\hat{f}_G(\alpha, \beta) := \mathbb{E}_{H \sim \pi} [\hat{f}_H(\alpha, \beta)]) \\
 &= \langle \alpha, \mu \rangle + \langle \beta, \nu \rangle - \langle \exp(\alpha \oplus \beta), \hat{K} \odot (\mu \otimes \nu) \rangle.
 \end{aligned}$$

Alg. (Random Coordinate Sinkhorn)

$$\left\{ \begin{array}{l}
 \alpha_k \leftarrow \alpha_{k;N_k} \approx \operatorname{argmax}_{\alpha} \hat{f}_G(\alpha, \beta_{k-1}) \\
 \quad \text{for } i = 1, \dots, N_k: \\
 \quad \text{Sample } H_i \sim \pi \\
 \quad \alpha_{k;i+\frac{1}{2}} \leftarrow \operatorname{argmax}_{\alpha} \hat{f}_{H_i}(\alpha, \beta_{k-1}) \\
 \quad \alpha_{k;i} \leftarrow (1 - \eta) \alpha_{k;i-1} + \eta \alpha_{k;i+\frac{1}{2}} \\
 \beta_k \leftarrow \beta_{k;N_k} \approx \operatorname{argmax}_{\beta} \hat{f}_G(\alpha_k, \beta) \\
 \quad \text{Similarly}
 \end{array} \right.$$

Alg. (Full Sinkhorn)

$$\left\{ \begin{array}{l}
 \alpha_k \leftarrow \operatorname{argmax}_{\alpha} f(\alpha, \beta_{k-1}) \\
 \beta_k \leftarrow \operatorname{argmax}_{\beta} f(\alpha_k, \beta)
 \end{array} \right.$$

- **Q2: How do we compute the stochastic surrogate potentials $(\hat{\alpha}, \hat{\beta})$?**

$$\begin{aligned}
 (\hat{\alpha}, \hat{\beta}) &\in \operatorname{argmax}_{\alpha, \beta: V \rightarrow \mathbb{R}} (\hat{f}_G(\alpha, \beta) := \mathbb{E}_{H \sim \pi} [\hat{f}_H(\alpha, \beta)]) \\
 &= \langle \alpha, \mu \rangle + \langle \beta, \nu \rangle - \langle \exp(\alpha \oplus \beta), \hat{K} \odot (\mu \otimes \nu) \rangle.
 \end{aligned}$$

Alg. (Random Coordinate Sinkhorn)

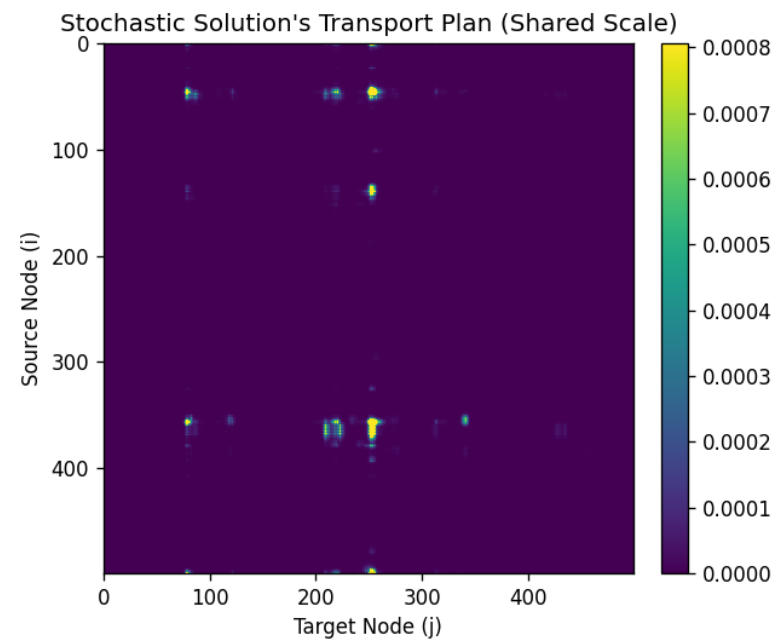
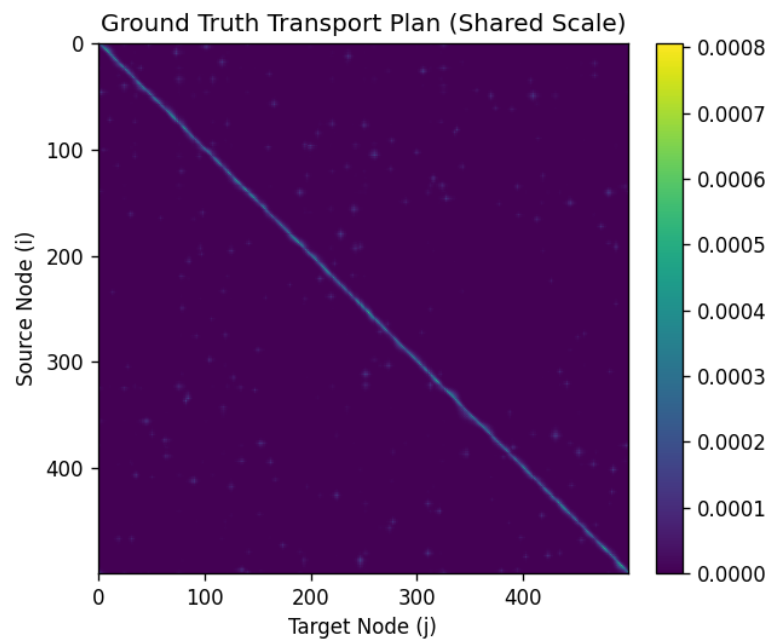
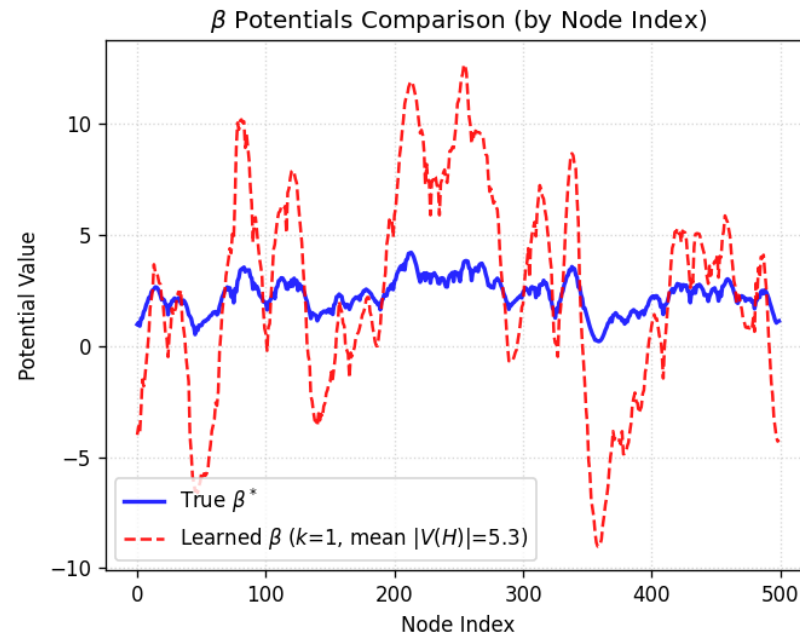
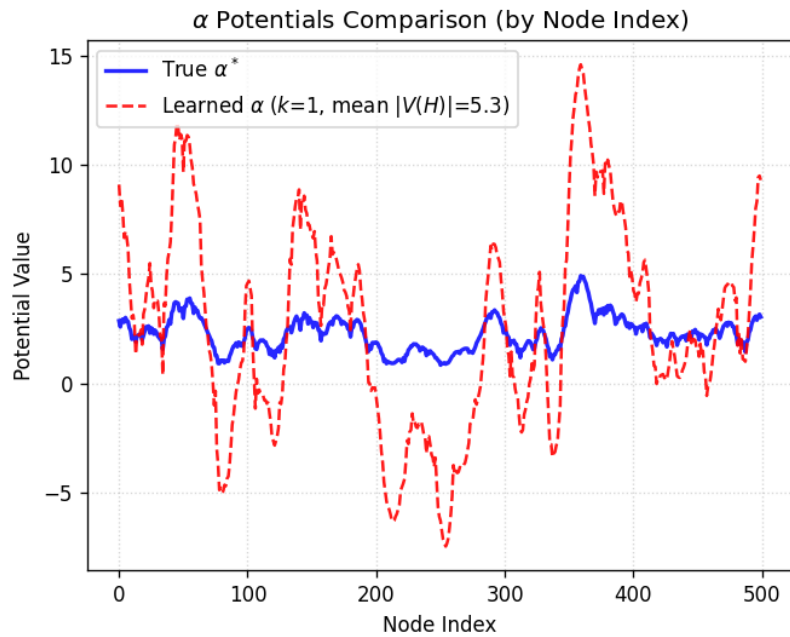
$$\left\{ \begin{array}{l}
 \alpha_k \leftarrow \alpha_{k;N_k} \approx \operatorname{argmax}_{\alpha} \hat{f}_G(\alpha, \beta_{k-1}) \\
 \quad \text{for } i = 1, \dots, N_k: \\
 \quad \text{Sample } H_i \sim \pi \\
 \quad \alpha_{k;i+\frac{1}{2}} \leftarrow \operatorname{argmax}_{\alpha} \hat{f}_{H_i}(\alpha, \beta_{k-1}) \\
 \quad \alpha_{k;i} \leftarrow (1 - \eta) \alpha_{k;i-1} + \eta \alpha_{k;i+\frac{1}{2}} \\
 \beta_k \leftarrow \beta_{k;N_k} \approx \operatorname{argmax}_{\beta} \hat{f}_G(\alpha_k, \beta) \\
 \quad \text{Similarly}
 \end{array} \right.$$

Alg. (Full Sinkhorn)

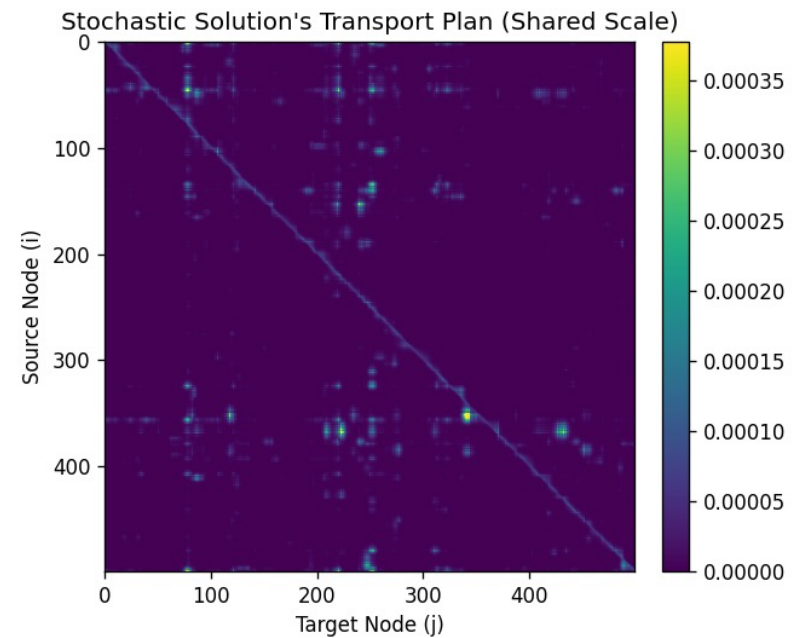
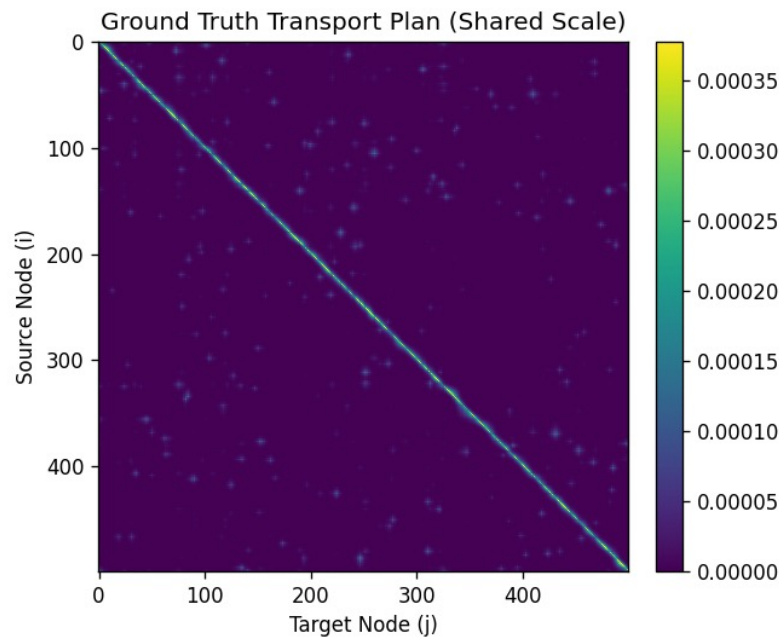
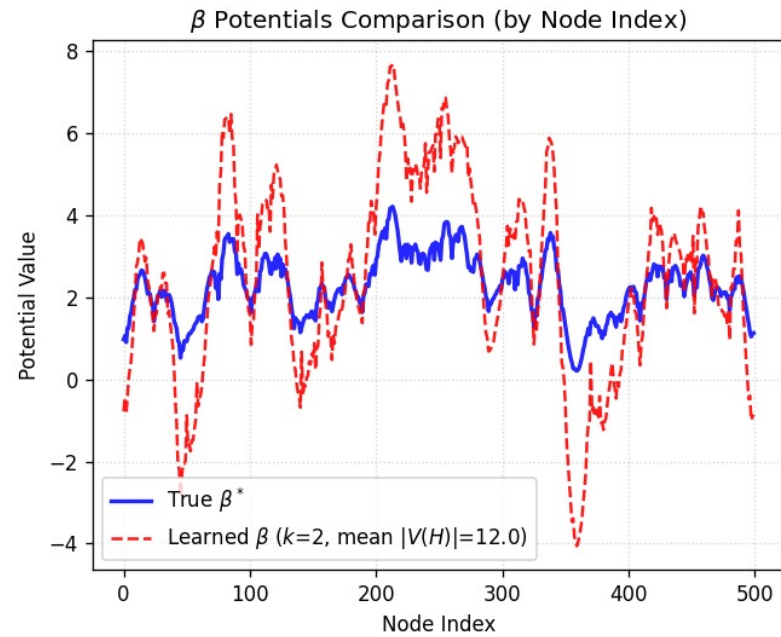
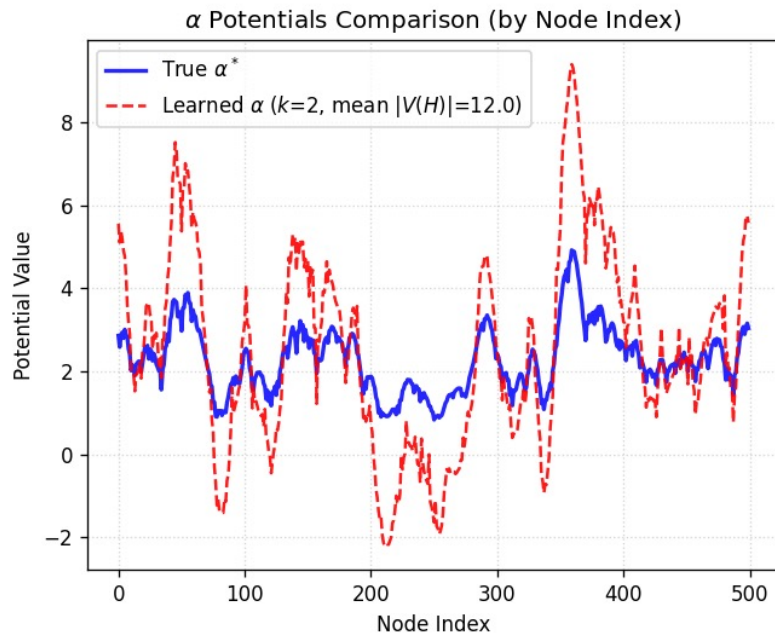
$$\left\{ \begin{array}{l}
 \alpha_k \leftarrow \operatorname{argmax}_{\alpha} f(\alpha, \beta_{k-1}) \\
 \beta_k \leftarrow \operatorname{argmax}_{\beta} f(\alpha_k, \beta)
 \end{array} \right.$$

**A practical choice of random subgraph H
= k-ball at a uniformly random center**

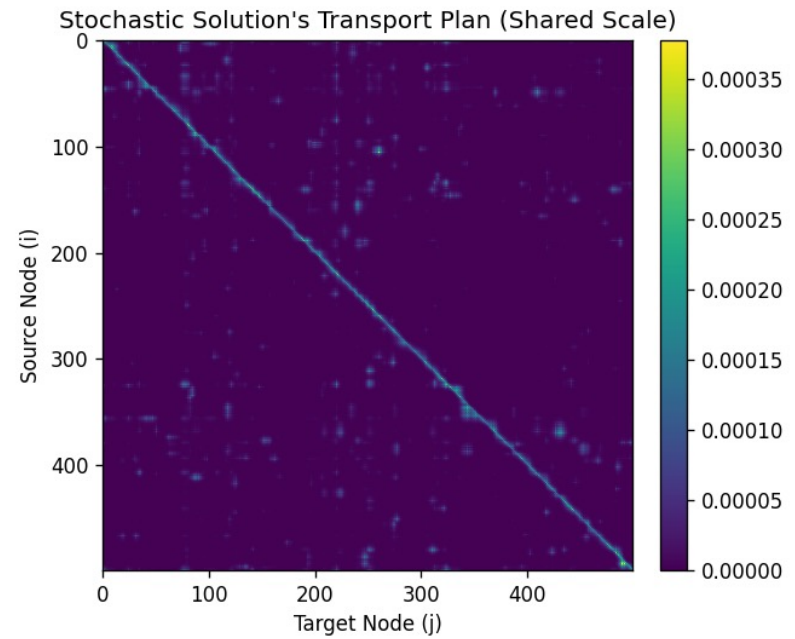
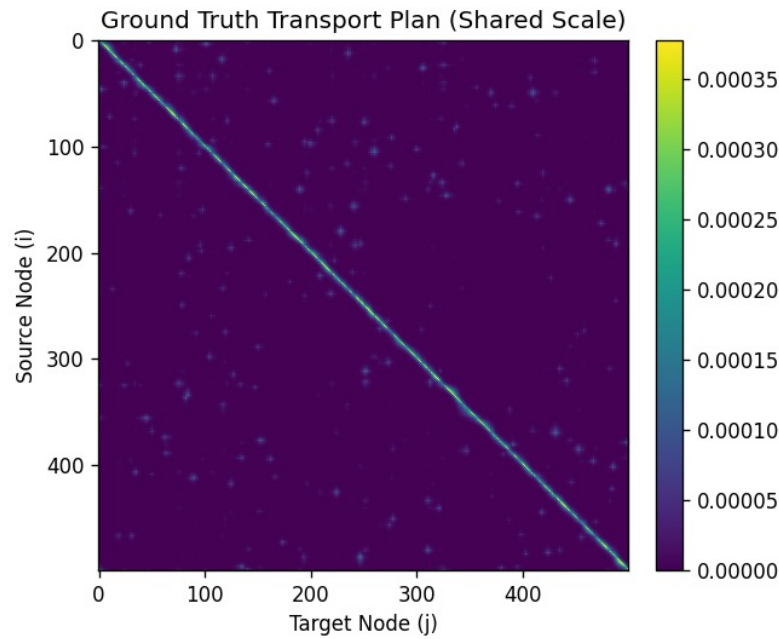
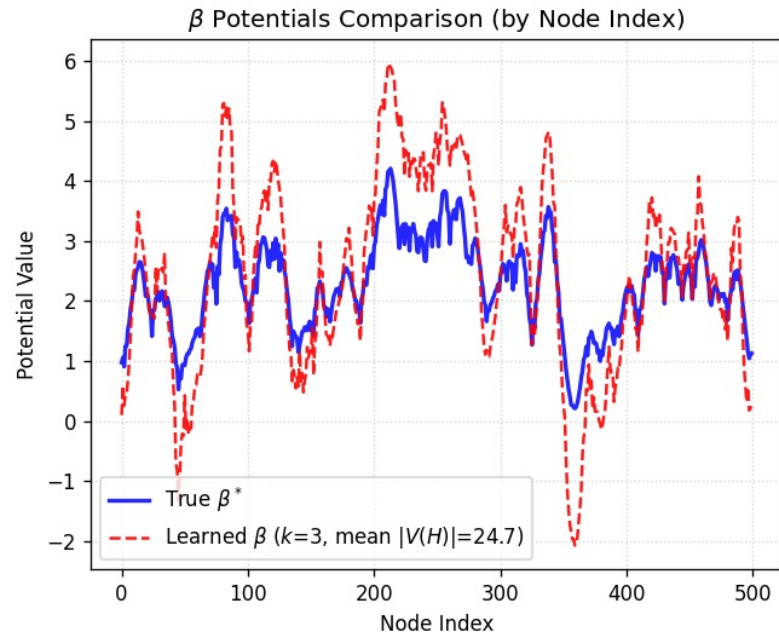
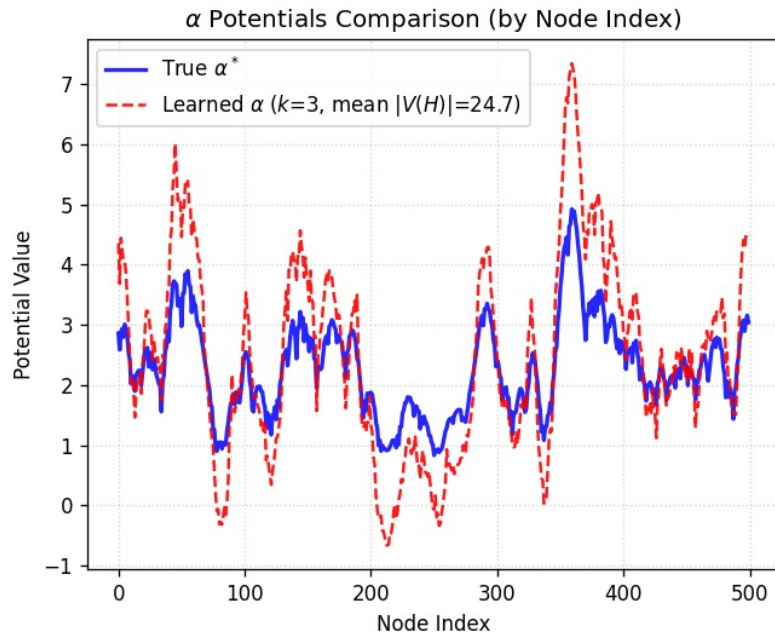
RC-Sinkhorn (k-ball, k=1) vs. Ground Truth $N=500$, $\varepsilon=1.0$, $K_{\text{outer}}=3000$, $N_{\text{inner}}=10$, $\eta=0.25$ [K3000_N10_eta0p25]



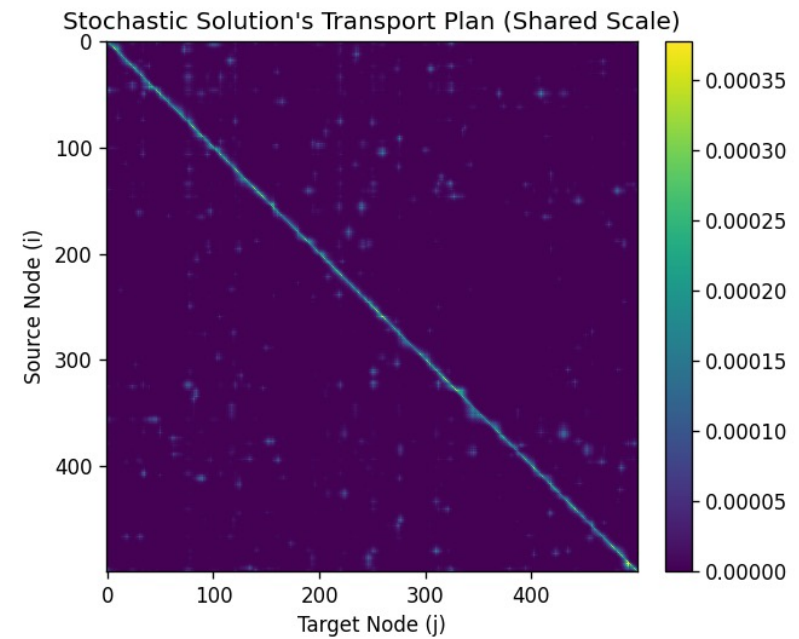
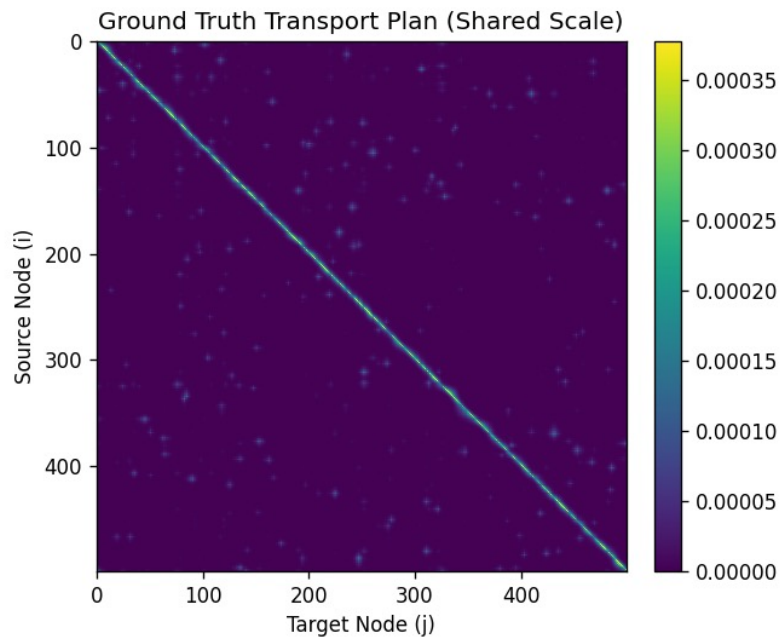
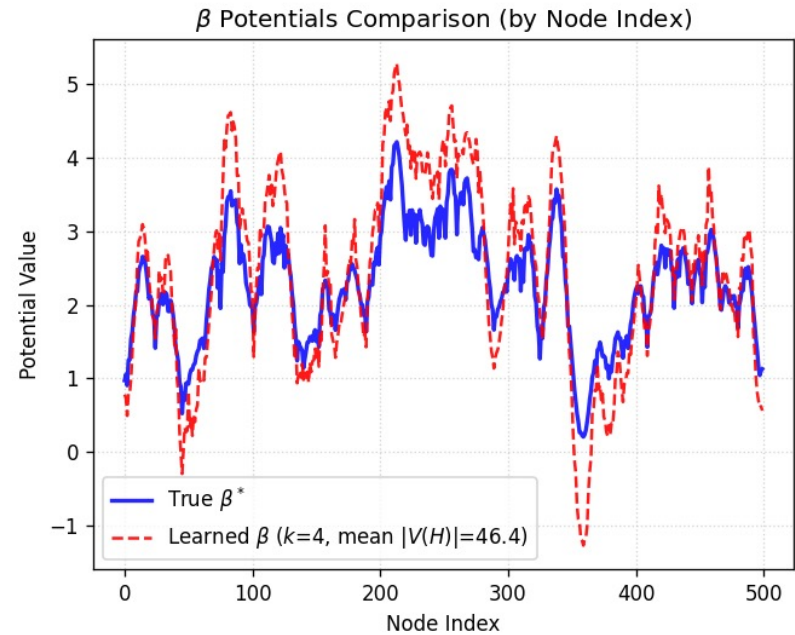
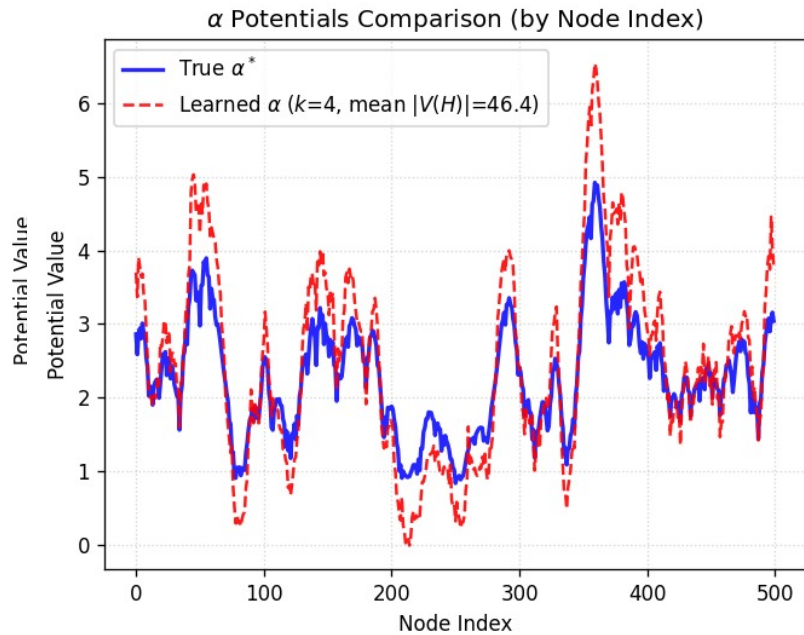
RC-Sinkhorn (k-ball, k=2) vs. Ground Truth $N=500, \epsilon=1.0, K_{\text{outer}}=3000, N_{\text{inner}}=10, \eta=0.25$ [K3000_N10_eta0p25]



RC-Sinkhorn (k-ball, k=3) vs. Ground Truth $N=500$, $\epsilon=1.0$, $K_{\text{outer}}=3000$, $N_{\text{inner}}=10$, $\eta=0.25$ [K3000_N10_eta0p25]



RC-Sinkhorn (k-ball, k=4) vs. Ground Truth $N=500$, $\varepsilon=1.0$, $K_{\text{outer}}=3000$, $N_{\text{inner}}=10$, $\eta=0.25$ [K3000_N10_eta0p25]



► Convergence analysis

Alg. (Random Coordinate Sinkhorn)

$$\left\{ \begin{array}{l} \boldsymbol{\alpha}_k \leftarrow \boldsymbol{\alpha}_{k;N_k} \approx \operatorname{argmax}_{\boldsymbol{\alpha}} \hat{f}_G(\boldsymbol{\alpha}, \boldsymbol{\beta}_{k-1}) \\ \text{for } i = 1, \dots, N_k: \\ \quad \text{Sample } H_i \sim \pi \\ \quad \boldsymbol{\alpha}_{k;i+\frac{1}{2}} \leftarrow \operatorname{argmax}_{\boldsymbol{\alpha}} \hat{f}_{H_i}(\boldsymbol{\alpha}, \boldsymbol{\beta}_{k-1}) \\ \quad \boldsymbol{\alpha}_{k;i} \leftarrow (1 - \eta) \boldsymbol{\alpha}_{k;i-1} + \eta \boldsymbol{\alpha}_{k;i+\frac{1}{2}} \\ \boldsymbol{\beta}_k \leftarrow \boldsymbol{\beta}_{k;N_k} \approx \operatorname{argmax}_{\boldsymbol{\beta}} \hat{f}_G(\boldsymbol{\alpha}_k, \boldsymbol{\beta}) \\ \text{Similarly} \end{array} \right. \cup$$

Alg. (Inexact Sinkhorn)

$$\left\{ \begin{array}{l} \boldsymbol{\alpha}_k \approx \operatorname{argmax}_{\boldsymbol{\alpha}} f(\boldsymbol{\alpha}, \boldsymbol{\beta}_{k-1}) \\ \boldsymbol{\beta}_k \approx \operatorname{argmax}_{\boldsymbol{\beta}} f(\boldsymbol{\alpha}_k, \boldsymbol{\beta}) \end{array} \right.$$

► Convergence analysis

Alg. (Random Coordinate Sinkhorn)

$$\left\{ \begin{array}{l} \boldsymbol{\alpha}_k \leftarrow \boldsymbol{\alpha}_{k;N_k} \approx \operatorname{argmax}_{\boldsymbol{\alpha}} \hat{f}_G(\boldsymbol{\alpha}, \boldsymbol{\beta}_{k-1}) \\ \text{for } i = 1, \dots, N_k: \\ \quad \text{Sample } H_i \sim \pi \\ \quad \boldsymbol{\alpha}_{k;i+\frac{1}{2}} \leftarrow \operatorname{argmax}_{\boldsymbol{\alpha}} \hat{f}_{H_i}(\boldsymbol{\alpha}, \boldsymbol{\beta}_{k-1}) \\ \quad \boldsymbol{\alpha}_{k;i} \leftarrow (1 - \eta) \boldsymbol{\alpha}_{k;i-1} + \eta \boldsymbol{\alpha}_{k;i+\frac{1}{2}} \\ \boldsymbol{\beta}_k \leftarrow \boldsymbol{\beta}_{k;N_k} \approx \operatorname{argmax}_{\boldsymbol{\beta}} \hat{f}_G(\boldsymbol{\alpha}_k, \boldsymbol{\beta}) \\ \text{Similarly} \end{array} \right. \quad \sqcup$$

Alg. (Inexact Sinkhorn)

$$\left\{ \begin{array}{l} \boldsymbol{\alpha}_k \approx \operatorname{argmax}_{\boldsymbol{\alpha}} f(\boldsymbol{\alpha}, \boldsymbol{\beta}_{k-1}) \\ \boldsymbol{\beta}_k \approx \operatorname{argmax}_{\boldsymbol{\beta}} f(\boldsymbol{\alpha}_k, \boldsymbol{\beta}) \end{array} \right.$$

A key issue in full Sinkhorn analysis:

The dual objective $f(\boldsymbol{\alpha}, \boldsymbol{\beta}) = \langle \mathbf{r}, \boldsymbol{\alpha} \rangle + \langle \boldsymbol{\beta}, \mathbf{c} \rangle - \langle \exp(\boldsymbol{\alpha} \oplus \boldsymbol{\beta}), \mathbf{W} \rangle$

is concave but has **unbounded level sets**

Unbounded smoothness and strong concavity parameters

► Convergence analysis

Alg. (Random Coordinate Sinkhorn)

$$\left\{ \begin{array}{l} \boldsymbol{\alpha}_k \leftarrow \boldsymbol{\alpha}_{k;N_k} \approx \operatorname{argmax}_{\boldsymbol{\alpha}} \hat{f}_G(\boldsymbol{\alpha}, \boldsymbol{\beta}_{k-1}) \\ \text{for } i = 1, \dots, N_k: \\ \quad \text{Sample } H_i \sim \pi \\ \quad \boldsymbol{\alpha}_{k;i+\frac{1}{2}} \leftarrow \operatorname{argmax}_{\boldsymbol{\alpha}} \hat{f}_{H_i}(\boldsymbol{\alpha}, \boldsymbol{\beta}_{k-1}) \\ \quad \boldsymbol{\alpha}_{k;i} \leftarrow (1 - \eta) \boldsymbol{\alpha}_{k;i-1} + \eta \boldsymbol{\alpha}_{k;i+\frac{1}{2}} \\ \boldsymbol{\beta}_k \leftarrow \boldsymbol{\beta}_{k;N_k} \approx \operatorname{argmax}_{\boldsymbol{\beta}} \hat{f}_G(\boldsymbol{\alpha}_k, \boldsymbol{\beta}) \\ \text{Similarly} \end{array} \right. \quad \sqcup$$

Alg. (Inexact Sinkhorn)

$$\left\{ \begin{array}{l} \boldsymbol{\alpha}_k \approx \operatorname{argmax}_{\boldsymbol{\alpha}} f(\boldsymbol{\alpha}, \boldsymbol{\beta}_{k-1}) \\ \boldsymbol{\beta}_k \approx \operatorname{argmax}_{\boldsymbol{\beta}} f(\boldsymbol{\alpha}_k, \boldsymbol{\beta}) \end{array} \right.$$

A key issue in full Sinkhorn analysis:

The dual objective $f(\boldsymbol{\alpha}, \boldsymbol{\beta}) = \langle \mathbf{r}, \boldsymbol{\alpha} \rangle + \langle \boldsymbol{\beta}, \mathbf{c} \rangle - \langle \exp(\boldsymbol{\alpha} \oplus \boldsymbol{\beta}), \mathbf{W} \rangle$

is concave but has **unbounded level sets**

Unbounded smoothness and strong concavity parameters

A priori confinement analysis: $\|\boldsymbol{\alpha}_k - \boldsymbol{\alpha}^*\|_{\infty}, \|\boldsymbol{\beta}_k - \boldsymbol{\beta}^*\|_{\infty} \leq \text{non-increasing in } k$

(The Jacobian of the Sinkhorn map is row-stochastic)

► Convergence analysis

Alg. (Inexact Sinkhorn)

$$\begin{cases} \alpha_k \approx \operatorname{argmax}_{\alpha} f(\alpha, \beta_{k-1}) \\ \beta_k \approx \operatorname{argmax}_{\beta} f(\alpha_k, \beta) \end{cases}$$

Thm. (Bi, Kim, L. 26+)

$$\Delta_k = \sup f - f(\alpha_k, \beta_k) = O(k^{-2})$$

provided asymptotically exact updates

► Convergence analysis

Alg. (Inexact Sinkhorn)

$$\begin{cases} \alpha_k \approx \arg \max_{\alpha} f(\alpha, \beta_{k-1}) \\ \beta_k \approx \arg \max_{\beta} f(\alpha_k, \beta) \end{cases}$$

Thm. (Bi, Kim, L. 26+)

$$\Delta_k = \sup f - f(\alpha_k, \beta_k) = O(k^{-2})$$

provided asymptotically exact updates

Can we show boundedness of iterates? $\|\alpha_k - \alpha^*\|_{\infty}, \|\beta_k - \beta^*\|_{\infty} \leq M$?

WTS: $D_k = \max \left\{ \|\alpha_k - \alpha^*\|_{sp}, \|\beta_k - \beta^*\|_{sp} \right\}$ is uniformly bdd

► Convergence analysis

Alg. (Inexact Sinkhorn)

$$\begin{cases} \alpha_k \approx \operatorname{argmax}_{\alpha} f(\alpha, \beta_{k-1}) \\ \beta_k \approx \operatorname{argmax}_{\beta} f(\alpha_k, \beta) \end{cases}$$

Thm. (Bi, Kim, L. 26+)

$$\Delta_k = \sup f - f(\alpha_k, \beta_k) = O(k^{-2})$$

provided asymptotically exact updates

Can we show boundedness of iterates? $\|\alpha_k - \alpha^*\|_{\infty}, \|\beta_k - \beta^*\|_{\infty} \leq M$?

WTS: $D_k = \max \left\{ \|\alpha_k - \alpha^*\|_{sp}, \|\beta_k - \beta^*\|_{sp} \right\}$ is uniformly bdd

We show the following recursion:

$$D_{k+1} \leq \tau(D_k)D_k + c_k$$

where $c_k = \|\alpha_k - \alpha_k^{\text{exact}}(\hat{\beta}_{k-1})\|_{sp} + \|\beta_k - \beta_k^{\text{exact}}(\hat{\alpha}_{k-1})\|_{sp}$

$\tau(D) =$ Explicit sp contraction ratio (< 1 for $D < \infty$)

If c_k is summable (can be controlled), then the recursion gives $\sup_k D_k \leq M$

New Perspective II: Random Matrix Scaling

► Recall matrix scaling..

Primal

Dual

$$\min_{\mathbf{Z} \in \mathcal{T}(\mathbf{r}, \mathbf{c})} [D_{KL}(\mathbf{Z} \parallel \mathbf{W})] = \max_{\boldsymbol{\alpha}, \boldsymbol{\beta}} [\langle \mathbf{r}, \boldsymbol{\alpha} \rangle + \langle \boldsymbol{\beta}, \mathbf{c} \rangle - \langle \exp(\boldsymbol{\alpha} \oplus \boldsymbol{\beta}), \mathbf{W} \rangle]$$

Rescaled Matrix

$\boldsymbol{\alpha}$ = Lagrange multipliers for the row-sum constraint

$\boldsymbol{\beta}$ = Lagrange multipliers for the column-sum constraint

(Schrödinger potentials)

$$\boldsymbol{\pi}^* = \text{Diag}(\exp(\boldsymbol{\alpha}^*)) \mathbf{W} \text{Diag}(\exp(\boldsymbol{\beta}^*))$$

Alternating maximization

Matrix scaling algorithm
(Iterative Proportional Fitting)

Fit the row sums;
Fit the column sums; etc.

Sinkhorn Algorithm

$$\begin{cases} \forall 1 \leq i \leq n, \boldsymbol{\beta}_k(j) \leftarrow \log \frac{\mathbf{c}(j)}{\sum_{i=1}^m \mathbf{W}_{ij} \exp(\boldsymbol{\alpha}_{k-1}(i))}, \\ \forall 1 \leq i \leq m, \boldsymbol{\alpha}_k(i) \leftarrow \log \frac{\mathbf{r}(i)}{\sum_{j=1}^n \mathbf{W}_{ij} \exp(\boldsymbol{\beta}_k(j))}. \end{cases}$$

► Sinkhorn-rescale a Random matrix X ??

Primal

Dual

$$\min_{\mathbf{Z} \in \mathcal{T}(\mathbf{r}, \mathbf{c})} [D_{KL}(\mathbf{Z} \parallel \mathbf{X})] = \max_{\boldsymbol{\alpha}, \boldsymbol{\beta}} [\langle \mathbf{r}, \boldsymbol{\alpha} \rangle + \langle \boldsymbol{\beta}, \mathbf{c} \rangle - \langle \exp(\boldsymbol{\alpha} \oplus \boldsymbol{\beta}), \mathbf{X} \rangle]$$

Rescaled Matrix

$\boldsymbol{\alpha}$ = Lagrange multipliers for the row-sum constraint

$\boldsymbol{\beta}$ = Lagrange multipliers for the column-sum constraint

(Schrödinger potentials)

$$\mathbf{X}^{r,c} := \pi^* = \text{Diag}(\exp(\boldsymbol{\alpha}^*)) \mathbf{X} \text{Diag}(\exp(\boldsymbol{\beta}^*))$$

Alternating maximization

Matrix scaling algorithm
(Iterative Proportional Fitting)

Fit the row sums;
Fit the column sums; etc.

Sinkhorn Algorithm

$$\begin{cases} \forall 1 \leq i \leq n, \boldsymbol{\beta}_k(j) \leftarrow \log \frac{\mathbf{c}(j)}{\sum_{i=1}^m \mathbf{X}_{ij} \exp(\boldsymbol{\alpha}_{k-1}(i))}, \\ \forall 1 \leq i \leq m, \boldsymbol{\alpha}_k(i) \leftarrow \log \frac{\mathbf{r}(i)}{\sum_{j=1}^n \mathbf{X}_{ij} \exp(\boldsymbol{\beta}_k(j))}. \end{cases}$$

► Sinkhorn-rescale a Random matrix X ??

Primal

Dual

$$\min_{\mathbf{Z} \in \mathcal{T}(\mathbf{r}, \mathbf{c})} [D_{KL}(\mathbf{Z} \parallel \mathbf{X})] = \max_{\boldsymbol{\alpha}, \boldsymbol{\beta}} [\langle \mathbf{r}, \boldsymbol{\alpha} \rangle + \langle \boldsymbol{\beta}, \mathbf{c} \rangle - \langle \exp(\boldsymbol{\alpha} \oplus \boldsymbol{\beta}), \mathbf{X} \rangle]$$

Rescaled Matrix

$\boldsymbol{\alpha}$ = Lagrange multipliers for the row-sum constraint

$\boldsymbol{\beta}$ = Lagrange multipliers for the column-sum constraint

(Schrödinger potentials)

$$\mathbf{X}^{\mathbf{r}, \mathbf{c}} = \text{Diag}(\exp(\boldsymbol{\alpha}^*)) \mathbf{X} \text{Diag}(\exp(\boldsymbol{\beta}^*))$$

Depends on X !

Alternating maximization

Matrix scaling algorithm
(Iterative Proportional Fitting)

Fit the row sums;
Fit the column sums; etc.

Sinkhorn Algorithm

$$\begin{cases} \forall 1 \leq i \leq n, \boldsymbol{\beta}_k(j) \leftarrow \log \frac{\mathbf{c}(j)}{\sum_{i=1}^m \mathbf{X}_{ij} \exp(\boldsymbol{\alpha}_{k-1}(i))}, \\ \forall 1 \leq i \leq m, \boldsymbol{\alpha}_k(i) \leftarrow \log \frac{\mathbf{r}(i)}{\sum_{j=1}^n \mathbf{X}_{ij} \exp(\boldsymbol{\beta}_k(j))}. \end{cases}$$

► Sinkhorn-rescale a Random matrix \mathbf{X} ??

$$\mathbf{X}^{r,c} = \text{Diag}(\exp(\boldsymbol{\alpha}^*)) \mathbf{X} \text{Diag}(\exp(\boldsymbol{\beta}^*))$$

Depends on \mathbf{X} !

Thm. (Duan, Powell, L. 26+)

$$\mathbf{X}^{r,c} \approx (\mathbb{E}\mathbf{X})^{r,c} \quad (\text{High-probability concentration in L1, L2 norm})$$

$$\mathbf{X}^{r,c} \rightarrow \pi^* = \underset{\pi \in \Pi(\mathbf{r}^\infty, \mathbf{c}^\infty)}{\text{argmin}} D_{KL}(\pi \parallel \mathbb{E}\mathbf{X}^\infty) \quad (\text{weak convergence w/ explicit rates})$$

Empirical Singular Value Dist. $(\frac{mn}{N\sqrt{n}} (\mathbf{X}^{r,c} - (\mathbb{E}\mathbf{X})^{r,c})) \rightarrow$ limiting distribution

CLT for potentials $(\boldsymbol{\alpha}, \boldsymbol{\beta})$ for $\text{mean}(\mathbf{X}_1, \dots, \mathbf{X}_M)$ and margin (\mathbf{r}, \mathbf{c})

► Sinkhorn-rescale a Random matrix \mathbf{X} ??

$$\mathbf{X}^{r,c} = \text{Diag}(\exp(\boldsymbol{\alpha}^*)) \mathbf{X} \text{Diag}(\exp(\boldsymbol{\beta}^*)) \approx \hat{\mathbf{X}}^{r,c} = \text{Diag}(\exp(\bar{\boldsymbol{\alpha}})) \mathbf{X} \text{Diag}(\exp(\bar{\boldsymbol{\beta}}))$$

Depends on \mathbf{X} !
Depends on $\mathbb{E}\mathbf{X}$!

Thm. (Duan, Powell, L. 26+)

$$\mathbf{X}^{r,c} \approx (\mathbb{E}\mathbf{X})^{r,c} \quad (\text{High-probability concentration in L1, L2 norm})$$

$$\mathbf{X}^{r,c} \rightarrow \pi^* = \underset{\pi \in \Pi(\mathbf{r}^\infty, \mathbf{c}^\infty)}{\text{argmin}} D_{KL}(\pi \parallel \mathbb{E}\mathbf{X}^\infty) \quad (\text{weak convergence w/ explicit rates})$$

Empirical Singular Value Dist. $(\frac{mn}{N\sqrt{n}} (\mathbf{X}^{r,c} - (\mathbb{E}\mathbf{X})^{r,c})) \rightarrow$ limiting distribution

CLT for potentials $(\boldsymbol{\alpha}, \boldsymbol{\beta})$ for $\text{mean}(\mathbf{X}_1, \dots, \mathbf{X}_M)$ and margin (\mathbf{r}, \mathbf{c})

► Key ingredients: Stability Theory for Matrix Scaling

$$\Lambda^{\mathbf{r}, \mathbf{c}} = \operatorname{argmin}_{\mathbf{Z} \in \mathcal{T}(\mathbf{r}, \mathbf{c})} D_{KL}(\mathbf{Z} \parallel \Lambda = \operatorname{Diag}(\exp(\boldsymbol{\alpha})) \Lambda \operatorname{Diag}(\exp(\boldsymbol{\beta})))$$

How much does the rescaled matrix / potentials change when we perturb the input Λ or margin (\mathbf{r}, \mathbf{c}) ?

► Key ingredients: Stability Theory for Matrix Scaling

$$\Lambda^{\mathbf{r}, \mathbf{c}} = \operatorname{argmin}_{\mathbf{Z} \in \mathcal{F}(\mathbf{r}, \mathbf{c})} D_{KL}(\mathbf{Z} \parallel \Lambda) = \operatorname{Diag}(\exp(\boldsymbol{\alpha})) \Lambda \operatorname{Diag}(\exp(\boldsymbol{\beta}))$$

How much does the rescaled matrix / potentials change when we perturb the input Λ or margin (\mathbf{r}, \mathbf{c}) ?

Thm. (Duan, Powell, L. 26+)

$$d_{\text{Hellinger}}(\Lambda^{r,c}, \Sigma^{r,c}) \leq C d_{\text{Hellinger}}(\Lambda, \Sigma)$$

$$d_{\text{Hellinger}}(\Lambda^{r,c}, \Lambda^{r',c'}) \leq C' \|(\mathbf{r}, \mathbf{c}) - (\mathbf{r}', \mathbf{c}')\|_1$$

$$\|(\boldsymbol{\alpha}_\Lambda^{r,c}, \boldsymbol{\beta}_\Lambda^{r,c}) - (\boldsymbol{\alpha}_\Lambda^{r',c'}, \boldsymbol{\beta}_\Lambda^{r',c'})\|_\infty \leq C'' (\text{relative margin error})$$

New Perspective III: Poisson Contingency Tables

► $X_{|r,c} = X \sim \text{Poisson}(\Lambda)$ conditioned on having margin (\mathbf{r}, \mathbf{c}) ??

If $\Lambda = \text{rank } 1$, then $X_{|r,c} \sim \text{Multivariate Hypergeometric Dist.}$

For general Λ , no explicit form, no approximation scheme, etc.

► $\mathbf{X}_{|r,c} = \mathbf{X} \sim \text{Poisson}(\mathbf{\Lambda})$ conditioned on having margin (\mathbf{r}, \mathbf{c}) ??

If $\mathbf{\Lambda} = \text{rank } 1$, then $\mathbf{X}_{|r,c} \sim \text{Multivariate Hypergeometric Dist.}$

For general $\mathbf{\Lambda}$, no explicit form, no approximation scheme, etc.

Thm. (Duan, Powell, L., Mukherjee 26+)

$\mathbf{X}_{|r,c} =_d$ certain Gibbs random permutation of length $N = \text{total sum}$

$\mathbf{X}_{|r,c} \rightarrow \pi^* = \underset{\pi \in \Pi(\mathbf{r}^\infty, \mathbf{c}^\infty)}{\text{argmin}} D_{KL}(\pi \parallel \mathbb{E}\mathbf{X}^\infty)$ (weak convergence w/ explicit rates)

Empirical Singular Value Dist. $(\frac{mn}{N\sqrt{n}} (\mathbf{X}_{|r,c} - \mathbf{\Lambda}^{r,c})) \rightarrow \text{limiting distribution}$

► $\mathbf{X}_{|r,c} = \mathbf{X} \sim \text{Poisson}(\mathbf{\Lambda})$ conditioned on having margin (\mathbf{r}, \mathbf{c}) ??

If $\mathbf{\Lambda} = \text{rank } 1$, then $\mathbf{X}_{|r,c} \sim \text{Multivariate Hypergeometric Dist.}$

For general $\mathbf{\Lambda}$, no explicit form, no approximation scheme, etc.

Thm. (Duan, Powell, L., Mukherjee 26+)

$\mathbf{X}_{|r,c} =_d$ certain Gibbs random permutation of length $N = \text{total sum}$

$$\mathbf{X}_{|r,c} \rightarrow \pi^* = \underset{\pi \in \Pi(\mathbf{r}^\infty, \mathbf{c}^\infty)}{\text{argmin}} D_{KL}(\pi \parallel \mathbb{E}\mathbf{X}^\infty) \quad (\text{weak convergence w/ explicit rates})$$

Empirical Singular Value Dist. $(\frac{mn}{N\sqrt{n}} (\mathbf{X}_{|r,c} - \mathbf{\Lambda}^{r,c})) \rightarrow \text{limiting distribution}$

Proof uses: Transference principle, Varadhan's lemma, Barvinok's permanent identity

Thank you very much!