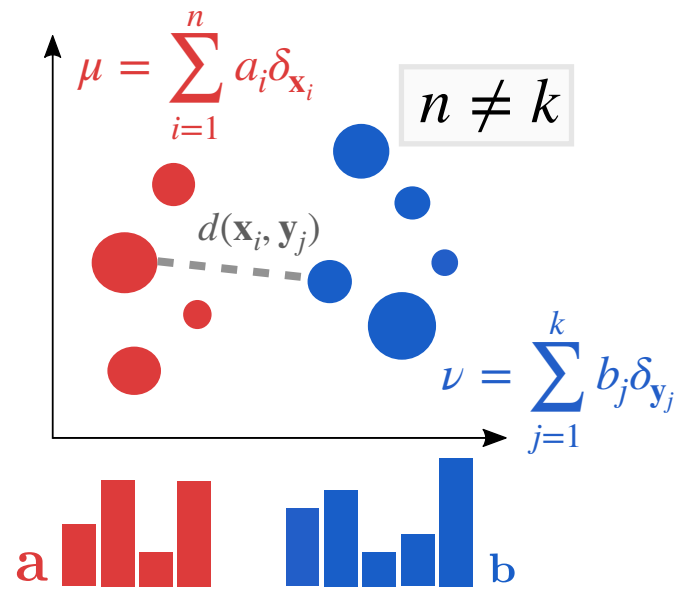


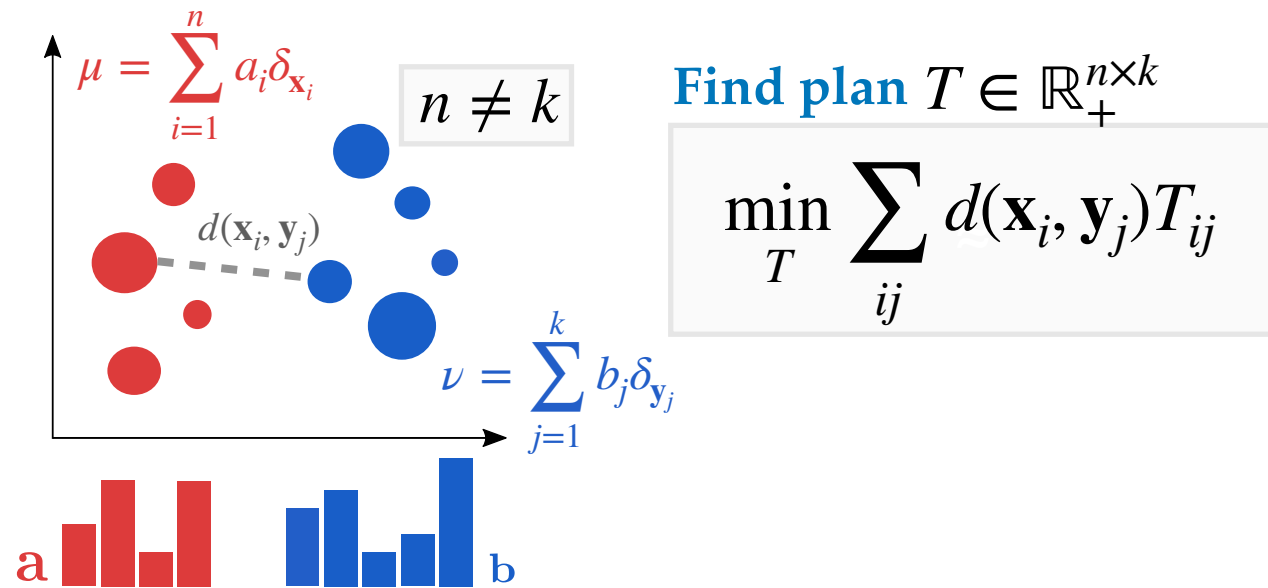
From Wasserstein to Sinkhorn

♦ Classical optimal transport (in a nutshell)



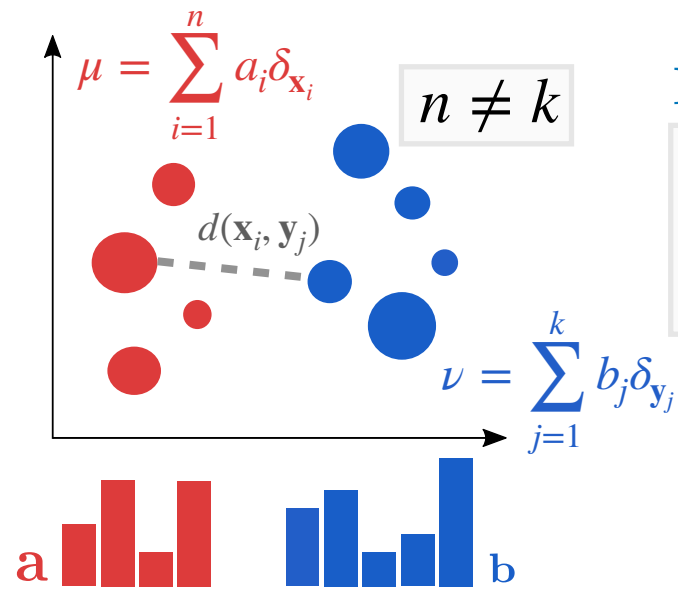
From Wasserstein to Sinkhorn

♦ Classical optimal transport (in a nutshell)



From Wasserstein to Sinkhorn

♦ Classical optimal transport (in a nutshell)



Find plan $T \in \mathbb{R}_+^{n \times k}$

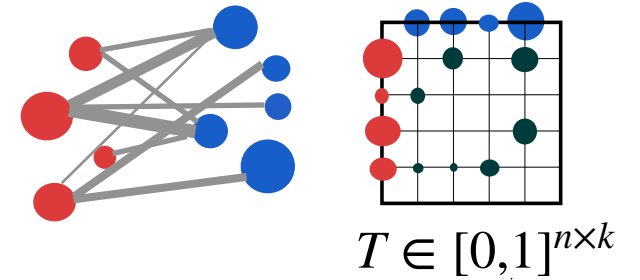
$$\min_T \sum_{ij} d(\mathbf{x}_i, \mathbf{y}_j) T_{ij}$$

which constraints ?

Coupling
 $\Pi(\mathbf{a}, \mathbf{b})$

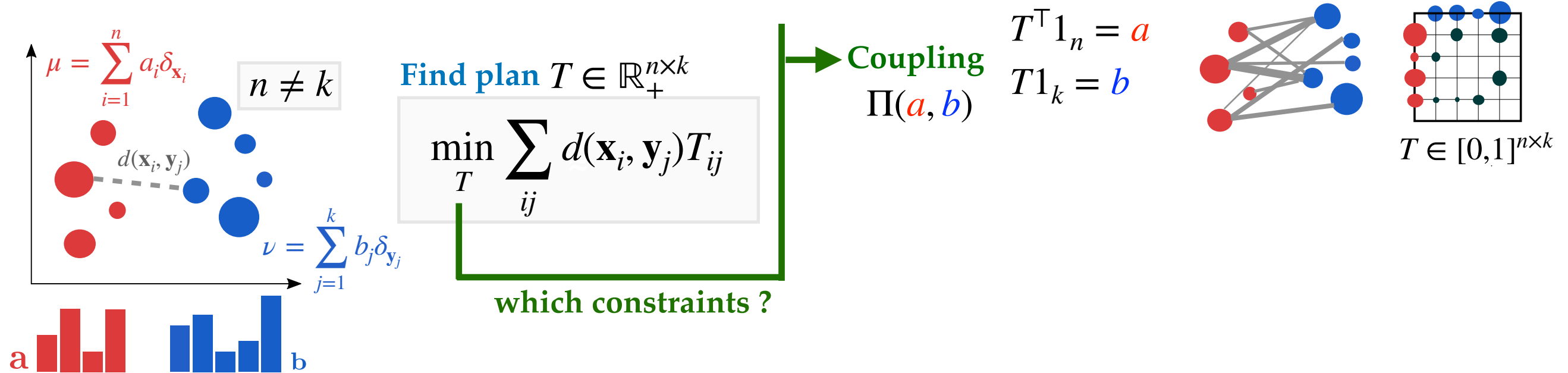
$$T^\top \mathbf{1}_n = \mathbf{a}$$

$$T \mathbf{1}_k = \mathbf{b}$$



From Wasserstein to Sinkhorn

♦ Classical optimal transport (in a nutshell)



♦ Wasserstein distance

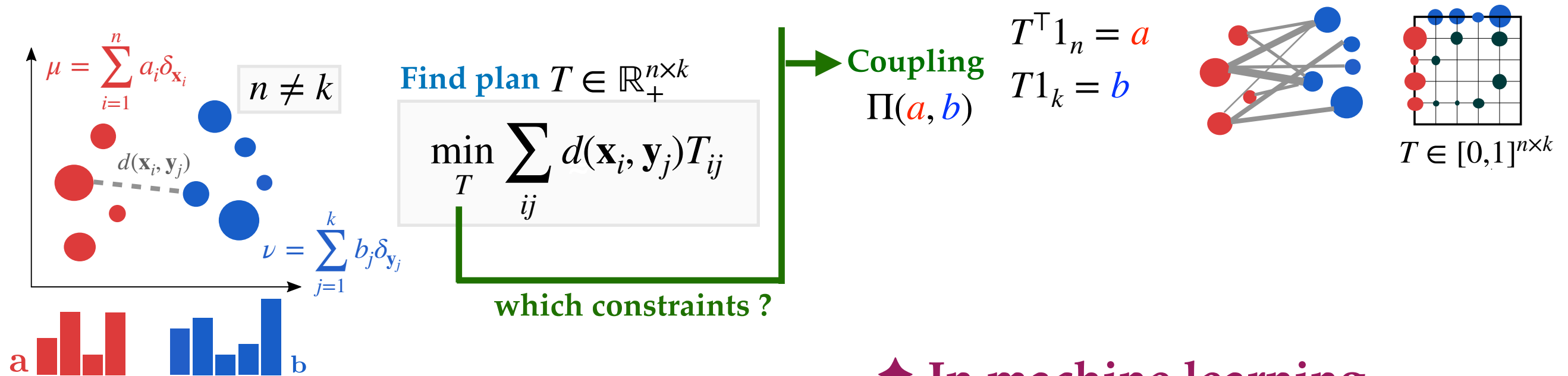
$$\begin{aligned} \mu &\in \mathcal{P}(X) \\ \nu &\in \mathcal{P}(X) \end{aligned}$$

$$W_p(\mu, \nu) = \left(\min_T \int_{X \times X} d(x, y)^p dT(x, y) \right)^{1/p}$$

- ♦ It is always **well-defined**
- ♦ It is a proper distance on $\mathcal{P}(X)$
- ♦ Lifts the geometry of $X \rightarrow \mathcal{P}(X)$

From Wasserstein to Sinkhorn

♦ Classical optimal transport (in a nutshell)



♦ Wasserstein distance

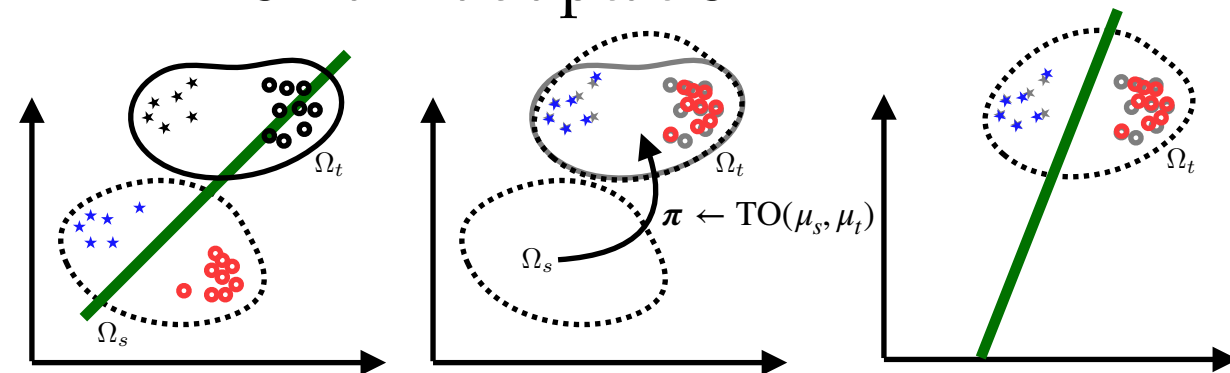
$$W_p(\mu, \nu) = \left(\min_T \int_{X \times X} d(x, y)^p dT(x, y) \right)^{1/p}$$

$\mu \in \mathcal{P}(X)$
 $\nu \in \mathcal{P}(X)$

- ♦ It is always **well-defined**
- ♦ It is a proper distance on $\mathcal{P}(X)$
- ♦ Lifts the geometry of $X \rightarrow \mathcal{P}(X)$

♦ In machine learning

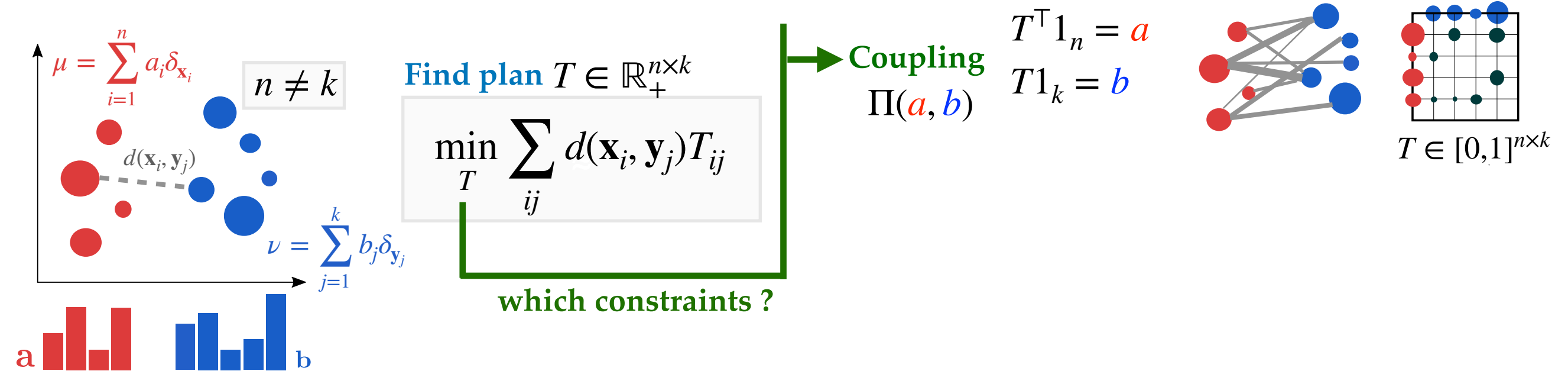
♦ Domain adaptation



- ♦ Generative modeling
- ♦ Analysis of NN convergence
- ♦ ML on graphs, fairness
- ♦ And many other ...

From Wasserstein to Sinkhorn

♦ Classical optimal transport (in a nutshell)



♦ Algorithmic foundations

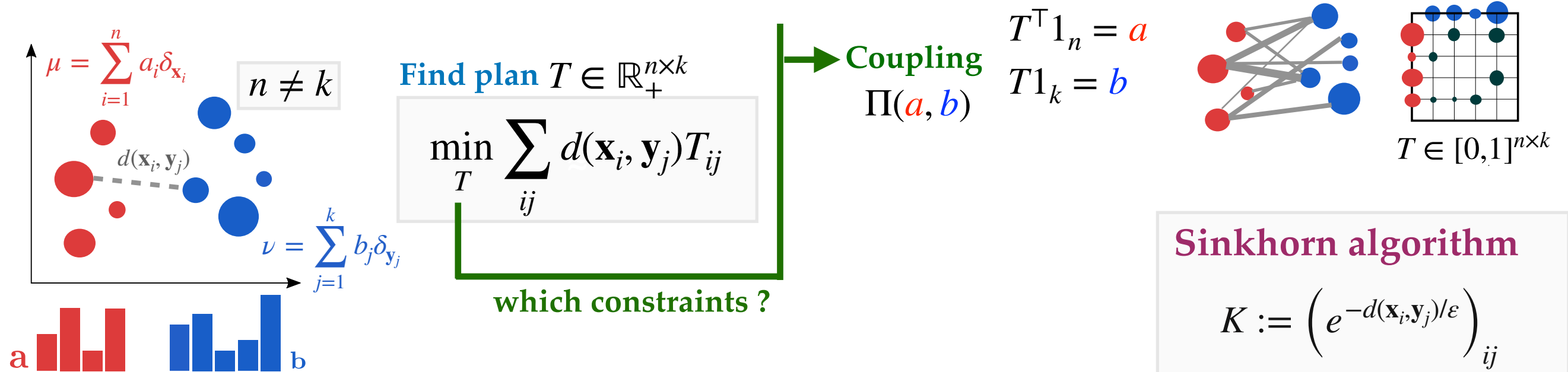
Unregularized problem

♦ Simplex, Network flow

$$\mathcal{O}(n^3 \log(n)^2)$$

From Wasserstein to Sinkhorn

♦ Classical optimal transport (in a nutshell)



♦ Algorithmic foundations

Unregularized problem

♦ Simplex, Network flow

$$\mathcal{O}(n^3 \log(n)^2)$$

Entropic regularization

$$\min_T \sum_{ij} d(\mathbf{x}_i, \mathbf{y}_j) T_{ij} - \varepsilon H(T)$$

Sinkhorn algorithm

$$K := \left(e^{-d(\mathbf{x}_i, \mathbf{y}_j)/\varepsilon} \right)_{ij}$$

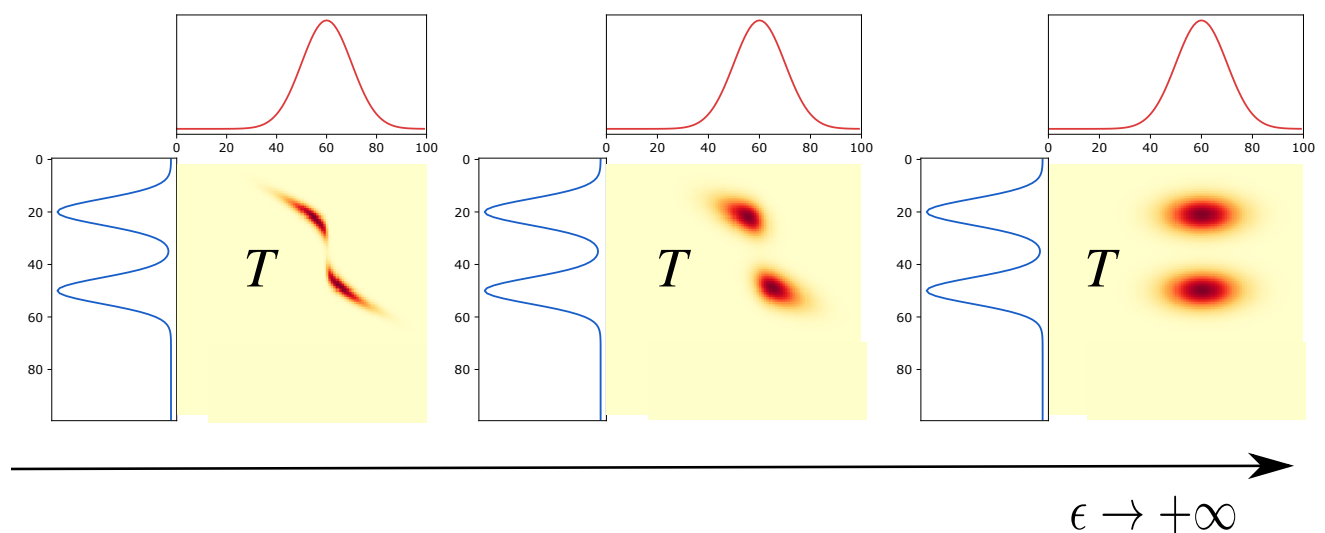
while not converged:

$$u = a \oslash K^\top v$$

$$v = b \oslash K u$$

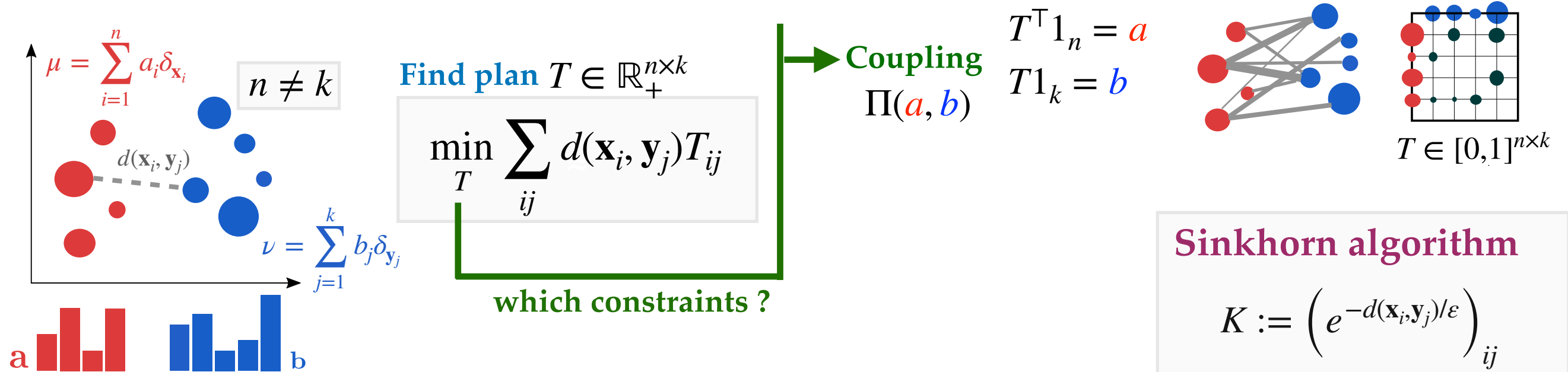
output

$$T = \text{diag}(u) K \text{diag}(v)$$



From Wasserstein to Sinkhorn

♦ Classical optimal transport (in a nutshell)



♦ Algorithmic foundations

Unregularized problem

♦ Simplex, Network flow

$$\mathcal{O}(n^3 \log(n)^2)$$

Entropic regularization

$$\min_T \sum_{ij} d(\mathbf{x}_i, \mathbf{y}_j) T_{ij} - \varepsilon H(T)$$

Regularized problem

$$\mathcal{O}(n^2)$$

Sinkhorn algorithm

$$K := \left(e^{-d(\mathbf{x}_i, \mathbf{y}_j)/\varepsilon} \right)_{ij}$$

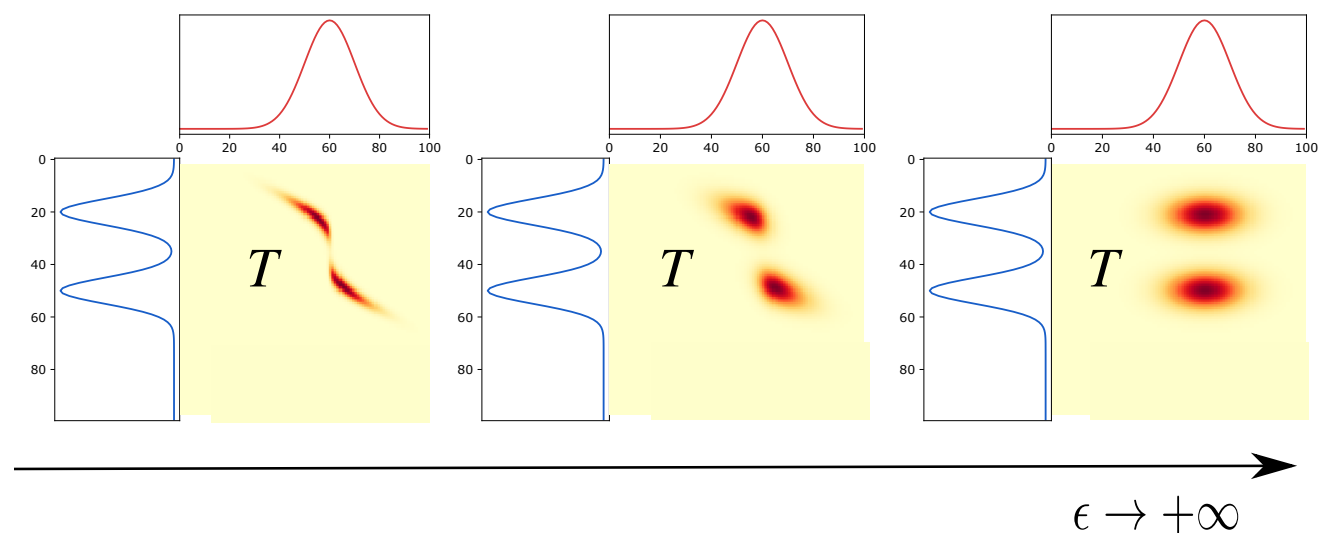
while not converged:

$$u = a \oslash K^T v$$

$$v = b \oslash K u$$

output

$$T = \text{diag}(u) K \text{diag}(v)$$



Accelerating Sinkhorn algorithm

♦ **Goal: fast approximation of** $u \rightarrow Ku$

Sinkhorn algorithm

$$K := \left(e^{-d(\mathbf{x}_i, \mathbf{y}_j)/\varepsilon} \right)_{ij}$$

while not converged:

$$u = \textcolor{red}{a} \oslash K^\top v$$

$$v = \textcolor{blue}{b} \oslash Ku$$

output

$$T = \text{diag}(u)K \text{diag}(v)$$

Accelerating Sinkhorn algorithm

♦ **Goal: fast approximation of $u \rightarrow Ku$**

Sinkhorn algorithm

$$K := \left(e^{-d(\mathbf{x}_i, \mathbf{y}_j)/\varepsilon} \right)_{ij}$$

while not converged:

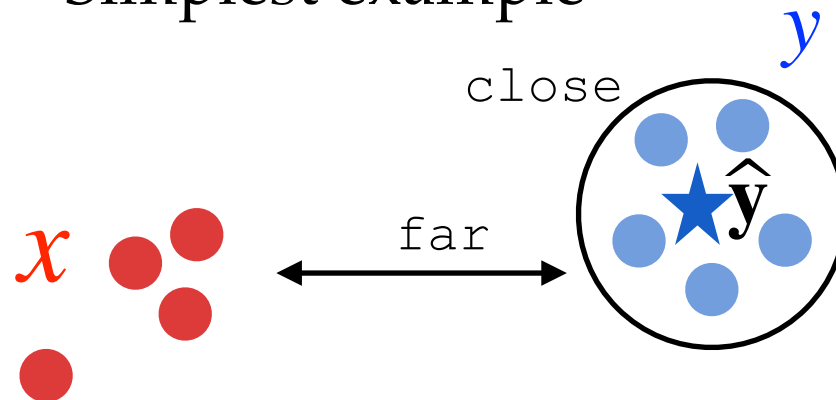
$$u = \textcolor{red}{a} \oslash K^\top v$$

$$v = \textcolor{blue}{b} \oslash Ku$$

output

$$T = \text{diag}(u)K \text{diag}(v)$$

♦ Simplest example



$$K \approx \mathbf{p}\mathbf{q}^\top = \begin{pmatrix} k(\mathbf{x}_1, \hat{\mathbf{y}}) \\ k(\mathbf{x}_2, \hat{\mathbf{y}}) \\ \vdots \\ k(\mathbf{x}_n, \hat{\mathbf{y}}) \end{pmatrix} \mathbf{1}^\top$$

$$u \rightarrow \tilde{K}u \text{ in } \mathcal{O}(n)$$

Accelerating Sinkhorn algorithm

◆ **Goal: fast approximation of $u \rightarrow Ku$**

Sinkhorn algorithm

$$K := \left(e^{-d(\mathbf{x}_i, \mathbf{y}_j)/\varepsilon} \right)_{ij}$$

while not converged:

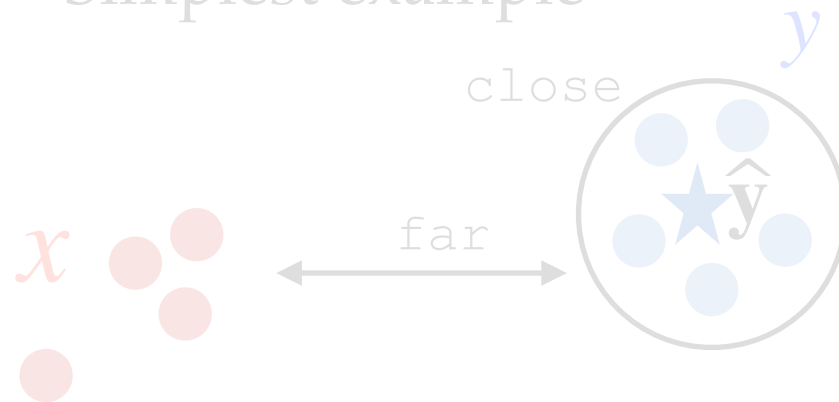
$$u = \tilde{a} \oslash K^\top v$$

$$v = \tilde{b} \oslash Ku$$

output

$$T = \text{diag}(u)K \text{diag}(v)$$

◆ Simplest example



$$K \approx \mathbf{p}\mathbf{q}^\top = \begin{pmatrix} k(\mathbf{x}_1, \hat{\mathbf{y}}) \\ k(\mathbf{x}_2, \hat{\mathbf{y}}) \\ \vdots \\ k(\mathbf{x}_n, \hat{\mathbf{y}}) \end{pmatrix} \mathbf{1}^\top$$

$$u \rightarrow \tilde{K}u \text{ in } \mathcal{O}(n)$$

◆ Generally if r clusters in y far enough from all the x

$$K \approx \mathbf{P}\mathbf{Q}^\top \quad u \rightarrow \tilde{K}u \text{ in } \mathcal{O}(rn)$$

◆ But unknown clusters + crude approximation !

Accelerating Sinkhorn algorithm

♦ **Goal: fast approximation of $u \rightarrow Ku$**

Sinkhorn algorithm

$$K := \left(e^{-d(\mathbf{x}_i, \mathbf{y}_j)/\varepsilon} \right)_{ij}$$

while not converged:

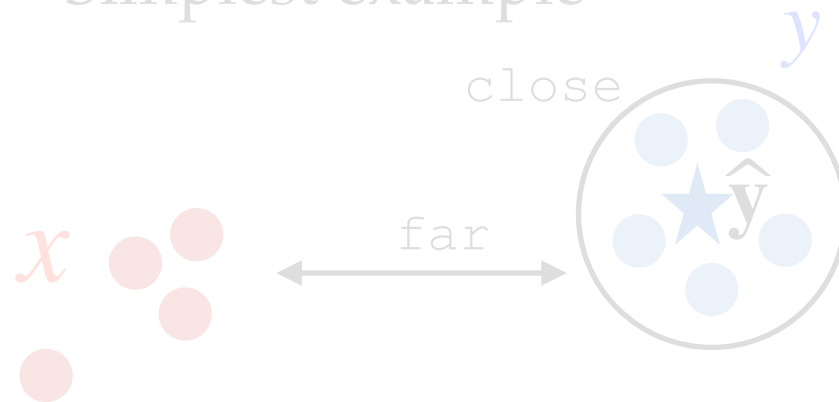
$$u = \tilde{a} \oslash K^\top v$$

$$v = \tilde{b} \oslash Ku$$

output

$$T = \text{diag}(u)K \text{diag}(v)$$

♦ Simplest example



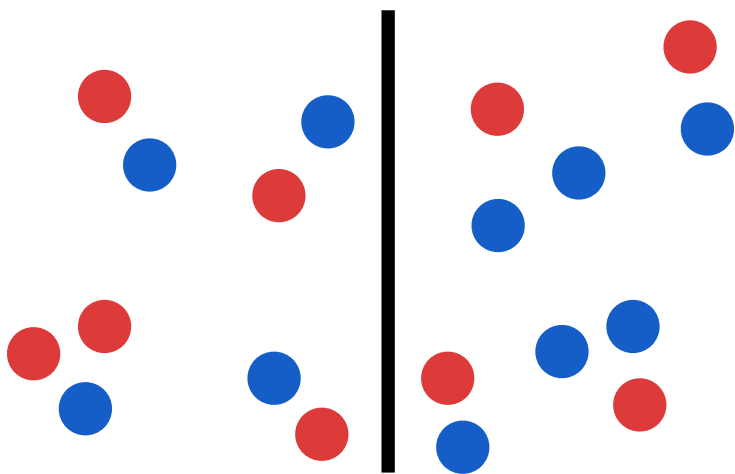
$$K \approx \mathbf{p}\mathbf{q}^\top = \begin{pmatrix} k(\mathbf{x}_1, \hat{\mathbf{y}}) \\ k(\mathbf{x}_2, \hat{\mathbf{y}}) \\ \vdots \\ k(\mathbf{x}_n, \hat{\mathbf{y}}) \end{pmatrix} \mathbf{1}^\top$$

$$u \rightarrow \tilde{K}u \text{ in } \mathcal{O}(n)$$

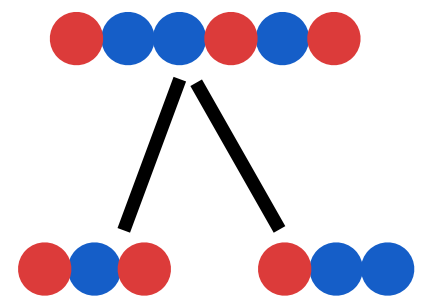
♦ Generally if r clusters in y far enough from all the x

$$K \approx \mathbf{P}\mathbf{Q}^\top \quad u \rightarrow \tilde{K}u \text{ in } \mathcal{O}(rn)$$

♦ **Idea: hierarchical clustering**



$$K = \begin{bmatrix} & \mathbf{x} \in C, \mathbf{y} \in \bar{C} \\ \mathbf{x} \in \bar{C}, \mathbf{y} \in C & \end{bmatrix}$$



Accelerating Sinkhorn algorithm

◆ **Goal: fast approximation of $u \rightarrow Ku$**

Sinkhorn algorithm

$$K := \left(e^{-d(\mathbf{x}_i, \mathbf{y}_j)/\varepsilon} \right)_{ij}$$

while not converged:

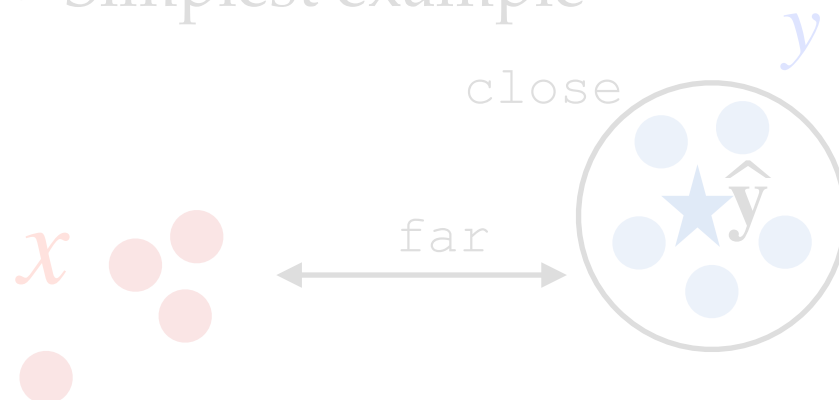
$$u = \tilde{a} \oslash K^\top v$$

$$v = \tilde{b} \oslash Ku$$

output

$$T = \text{diag}(u)K \text{diag}(v)$$

◆ Simplest example



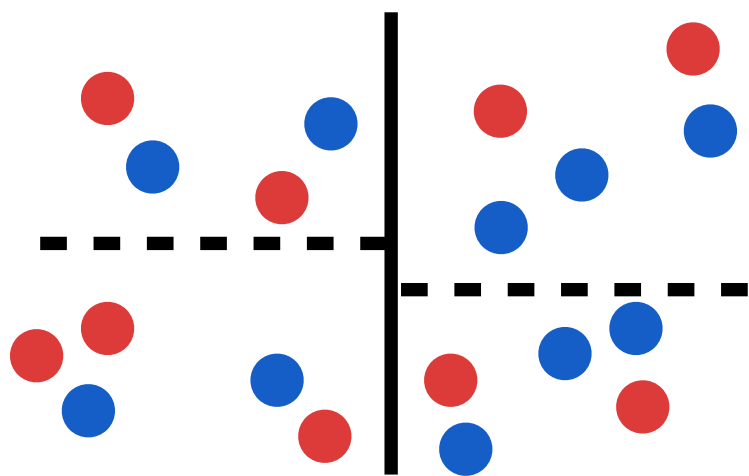
$$K \approx \mathbf{p}\mathbf{q}^\top = \begin{pmatrix} k(\mathbf{x}_1, \hat{\mathbf{y}}) \\ k(\mathbf{x}_2, \hat{\mathbf{y}}) \\ \vdots \\ k(\mathbf{x}_n, \hat{\mathbf{y}}) \end{pmatrix} \mathbf{1}^\top$$

$$u \rightarrow \tilde{K}u \text{ in } \mathcal{O}(n)$$

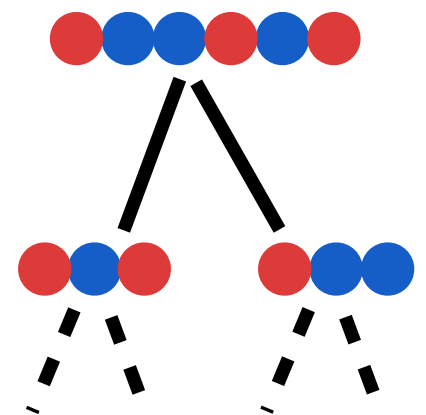
◆ Generally if r clusters in y far enough from all the x

$$K \approx \mathbf{P}\mathbf{Q}^\top \quad u \rightarrow \tilde{K}u \text{ in } \mathcal{O}(rn)$$

◆ **Idea: hierarchical clustering**



$$K = \begin{bmatrix} & \text{---} & & \\ \text{---} & & & \\ & & \text{---} & \\ \text{---} & & & \end{bmatrix} \begin{matrix} \\ \mathbf{x} \in C, \mathbf{y} \in \bar{C} \\ \\ \mathbf{x} \in \bar{C}, \mathbf{y} \in C \end{matrix}$$



Accelerating Sinkhorn algorithm

♦ **Goal: fast approximation of $u \rightarrow Ku$**

Sinkhorn algorithm

$$K := \left(e^{-d(\mathbf{x}_i, \mathbf{y}_j)/\varepsilon} \right)_{ij}$$

while not converged:

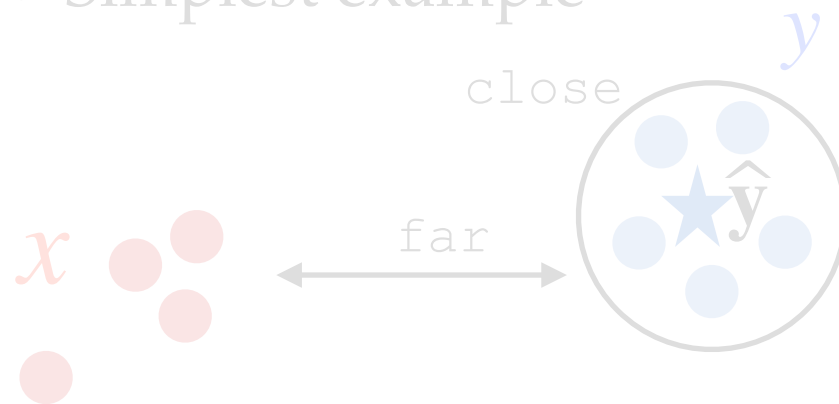
$$u = \tilde{a} \oslash K^\top v$$

$$v = \tilde{b} \oslash Ku$$

output

$$T = \text{diag}(u)K \text{diag}(v)$$

♦ Simplest example



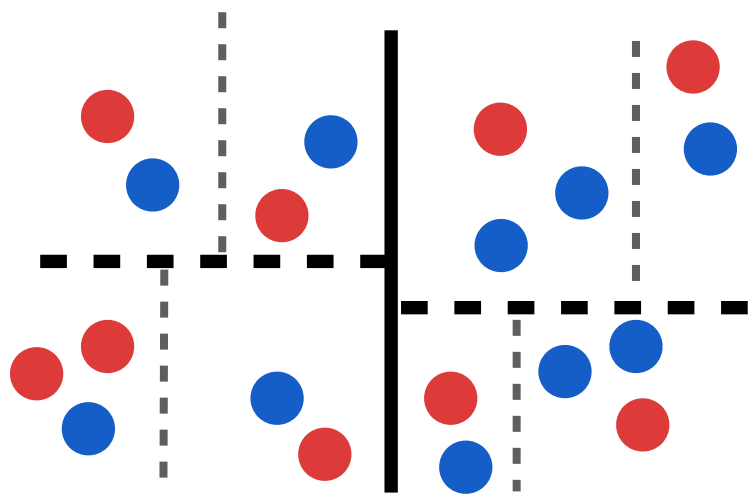
$$K \approx \mathbf{p}\mathbf{q}^\top = \begin{pmatrix} k(\mathbf{x}_1, \hat{\mathbf{y}}) \\ k(\mathbf{x}_2, \hat{\mathbf{y}}) \\ \vdots \\ k(\mathbf{x}_n, \hat{\mathbf{y}}) \end{pmatrix} \mathbf{1}^\top$$

$$u \rightarrow \tilde{K}u \text{ in } \mathcal{O}(n)$$

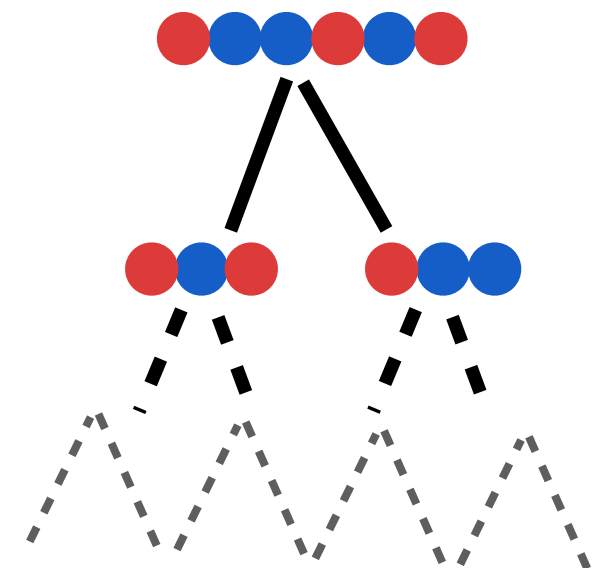
♦ Generally if r clusters in y far enough from all the x

$$K \approx \mathbf{P}\mathbf{Q}^\top \quad u \rightarrow \tilde{K}u \text{ in } \mathcal{O}(rn)$$

♦ **Idea: hierarchical clustering**



$$K = \begin{bmatrix} \text{---} & \text{---} \\ \text{---} & \text{---} \\ \text{---} & \text{---} \\ \text{---} & \text{---} \end{bmatrix} \quad \begin{matrix} \mathbf{x} \in C, \mathbf{y} \in \bar{C} \\ \mathbf{x} \in \bar{C}, \mathbf{y} \in C \end{matrix}$$



Accelerating Sinkhorn algorithm

♦ **Goal: fast approximation of $u \rightarrow Ku$**

Sinkhorn algorithm

$$K := \left(e^{-d(\mathbf{x}_i, \mathbf{y}_j)/\varepsilon} \right)_{ij}$$

while not converged:

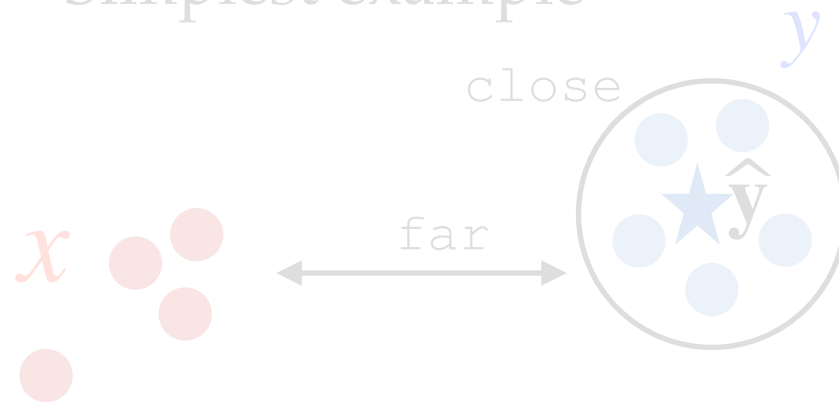
$$u = \tilde{a} \oslash K^\top v$$

$$v = \tilde{b} \oslash Ku$$

output

$$T = \text{diag}(u)K \text{diag}(v)$$

♦ Simplest example



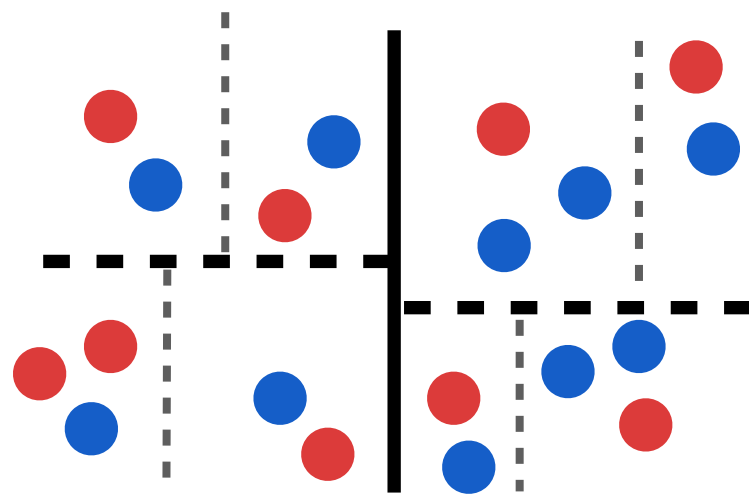
$$K \approx \mathbf{p}\mathbf{q}^\top = \begin{pmatrix} k(\mathbf{x}_1, \hat{\mathbf{y}}) \\ k(\mathbf{x}_2, \hat{\mathbf{y}}) \\ \vdots \\ k(\mathbf{x}_n, \hat{\mathbf{y}}) \end{pmatrix} \mathbf{1}^\top$$

$$u \rightarrow \tilde{K}u \text{ in } \mathcal{O}(n)$$

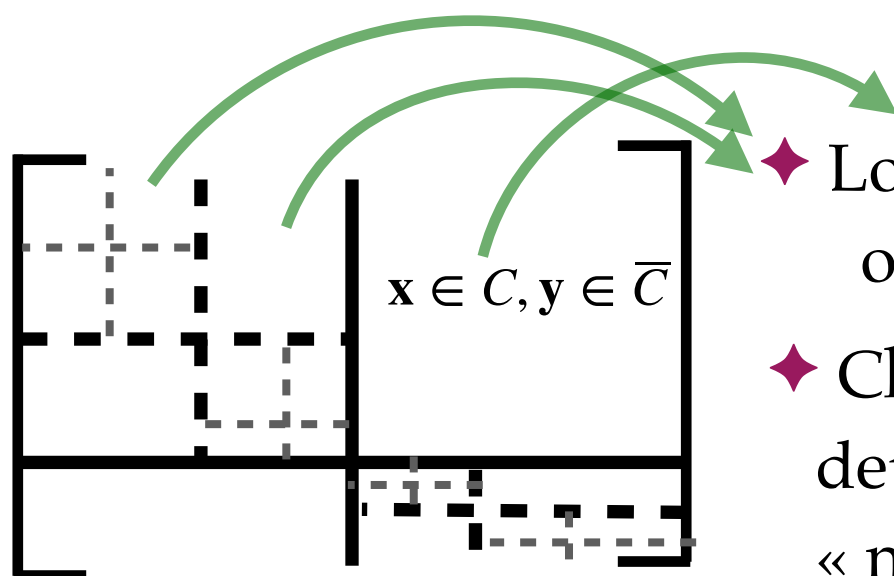
♦ Generally if r clusters in y far enough from all the x

$$K \approx \mathbf{P}\mathbf{Q}^\top \quad u \rightarrow \tilde{K}u \text{ in } \mathcal{O}(rn)$$

♦ **Idea: hierarchical clustering**



$$K =$$



- ♦ Low-rank approx of « far interactions »
- ♦ Cheap tree traversal to detect « far » and « near » interactions

♦ Fast multipole methods, Barnes-Hut algorithm, H-matrices