# A Review of Regularized Optimal Transport

## Marco Cuturi

ENSAE ParisTech

École nationale de la statistique et de l'administration économique

université PARIS-SACLAY

# What is Optimal Transport?

A geometric toolbox to
compare probability measures
supported on a metric space.



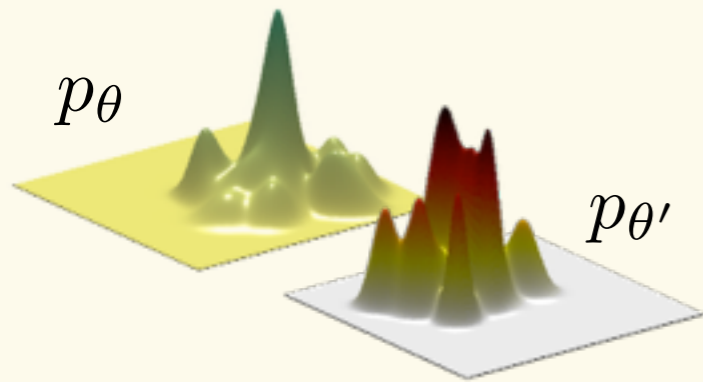Monge    Kantorovich    Dantzig    Wasserstein    Brenier    Otto    McCann    Villani
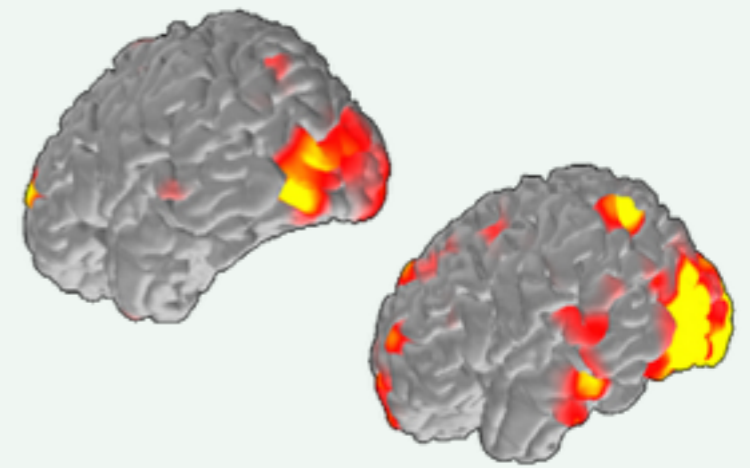
# What is Optimal Transport?

A geometric toolbox to
compare **probability measures**
supported on a metric space.



$p_\theta$
$p_{\theta'}$

*Statistical Models*



$h_1$
$h_2$

*Bags
of features*



*Brain Activation Maps*



$\nu$
$\mu$

*Empirical
Measures*



*Color Histograms*

# What is Optimal Transport?

A geometric toolbox to
compare probability measures
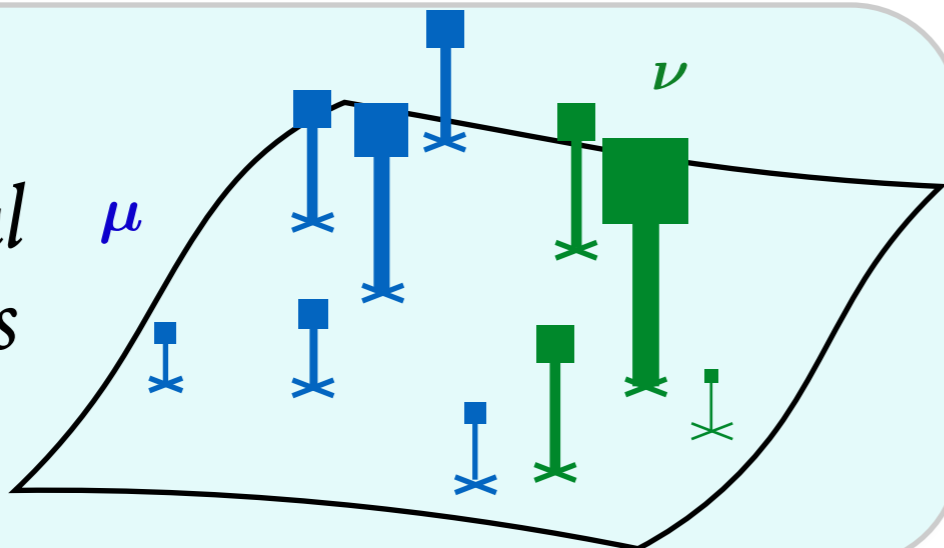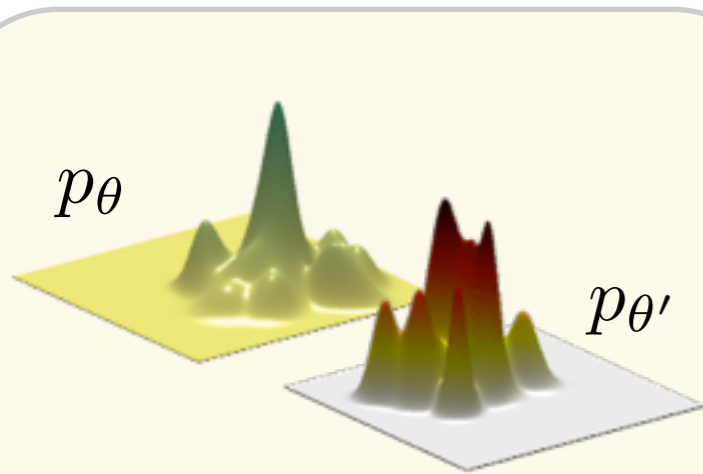supported on a **metric space**.



*Statistical Models*

$p_\theta$

$p_{\theta'}$



*Bags
of features*

$h_2$



*Brain Activation Maps*



*Empirical
Measures*

$\mu$

$\nu$



*Color Histograms*

# What is Optimal Transport?

A **geometric toolbox** to
compare probability measures
supported on a metric space.



$\mathcal{P}(\Omega)$

$p_\theta$

$p_{\theta'}$

# What is Optimal Transport?

A **geometric toolbox** to
compare probability measures
supported on a metric space.



$$W(p_\theta, p_{\theta'})$$

$\mathcal{P}(\Omega)$

$p_{\theta'}$

Wasserstein Distance

$p_\theta$

# What is Optimal Transport?

A **geometric toolbox** to
compare probability measures
supported on a metric space.



$\mathcal{P}(\Omega)$

$p_\theta$

$p_{\theta'}$

**[McCann'95]**
Interpolant

# What is Optimal Transport?

A **geometric toolbox** to compare probability measures supported on a metric space.

# What is Optimal Transport?
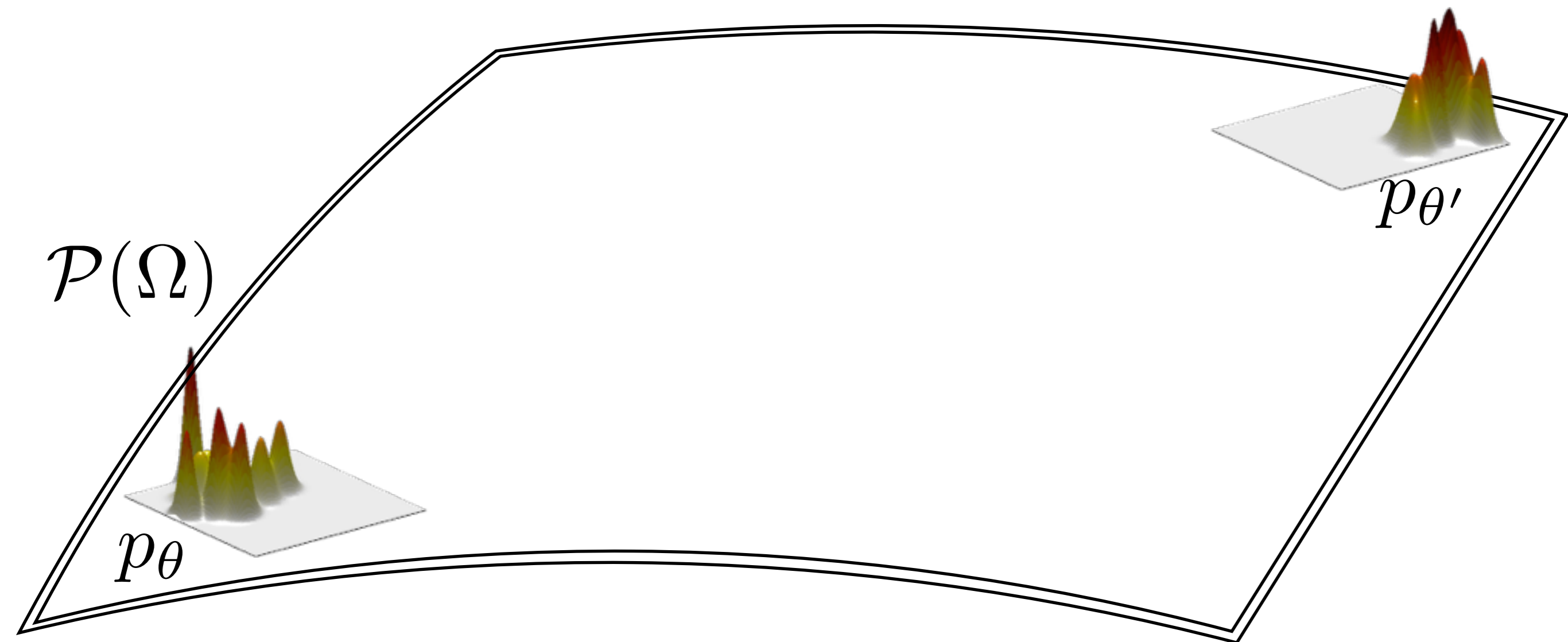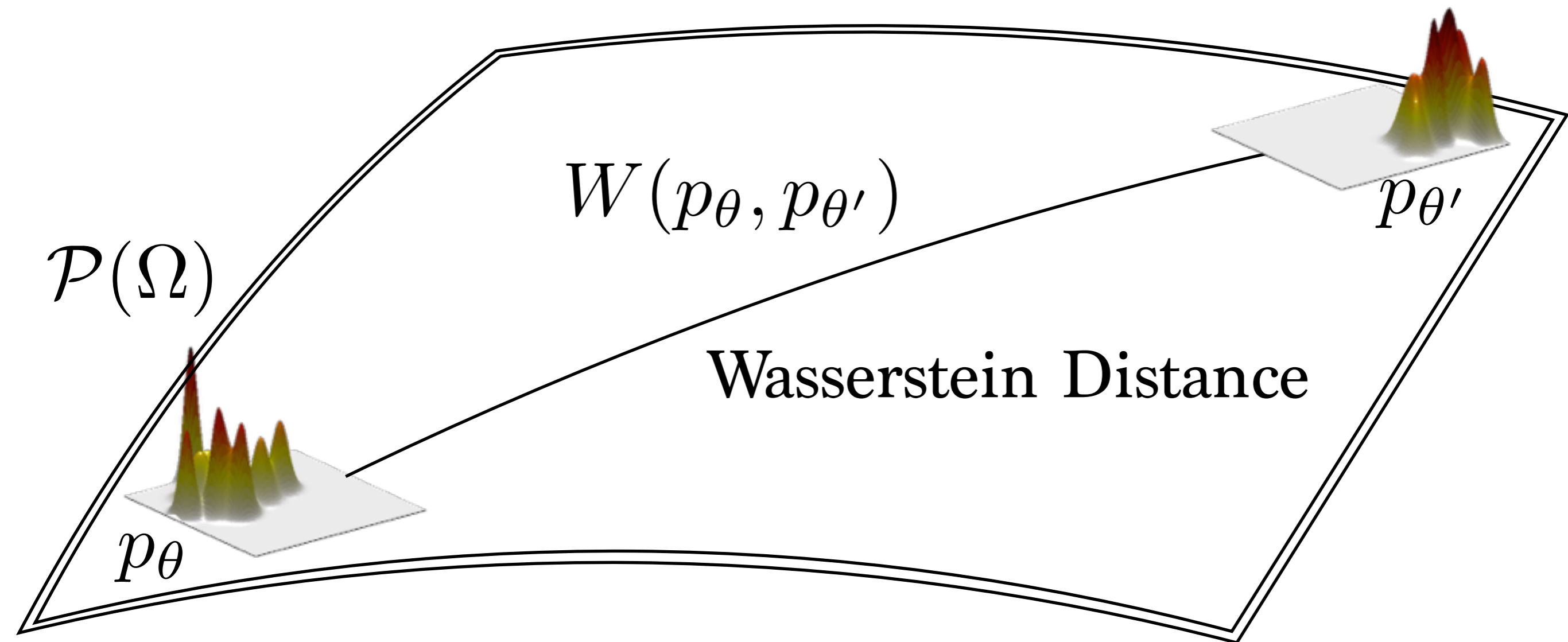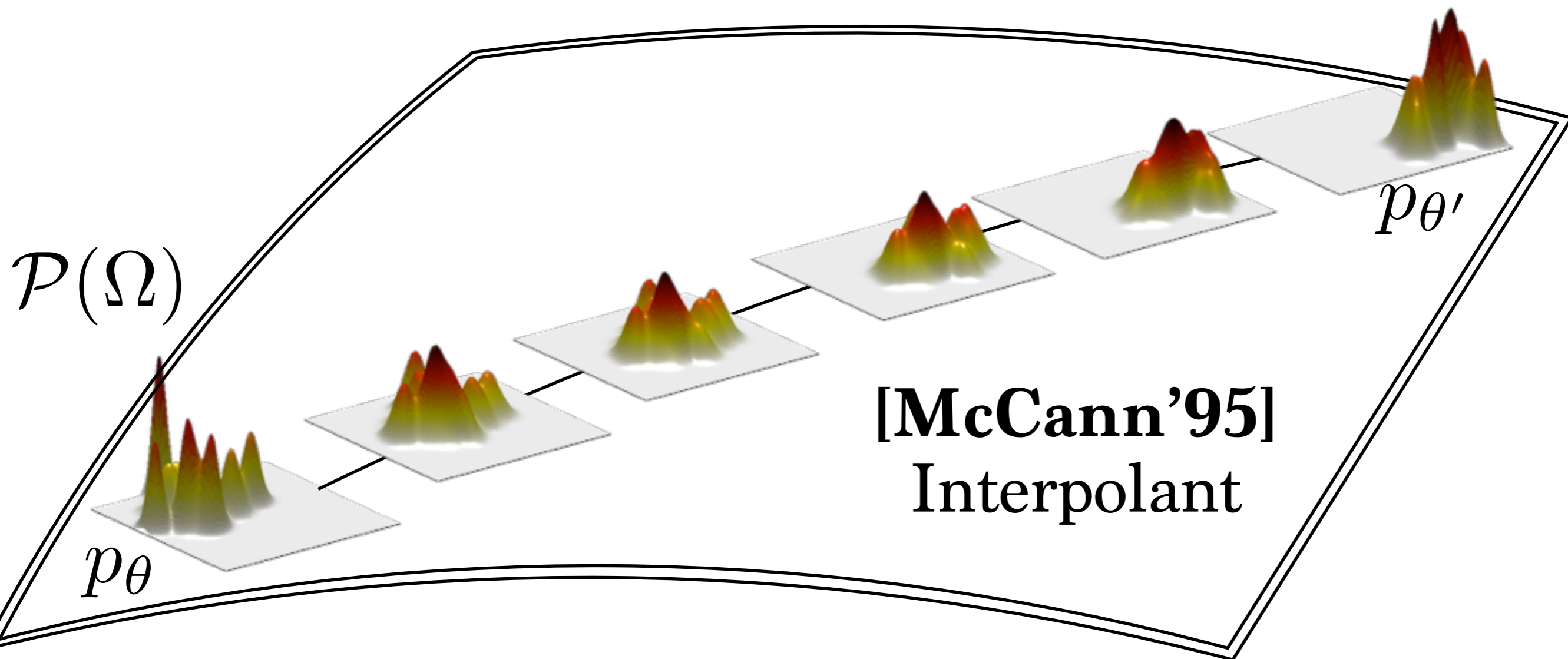
A **geometric toolbox** to compare probability measures supported on a metric space.



$P(\Omega)$

$p_{\theta'}$

$p_\theta$

$p_{\theta''}$

Wasserstein Barycenter [**Agueh'11**]

# OT and data-analysis

- Key developments in (applied) maths ~'90s
  **[McCann'95], [JKO'98], [Benamou'98], [Gangbo'98], [Ambrosio'06], [Villani'03/'09].**

- Key developments in TCS / graphics since '00s
  **[Rubner'98], [Indyk'03], [Naor'07], [Andoni'15].**

◉ Small to *no-impact* in large-scale data analysis:

  ✦ **computationally heavy;**

  ✦ **Wasserstein distance is not differentiable**

# OT and data-analysis

*Today's talk: Entropy Regularized OT*

- **Very fast** compared to usual approaches, <u>GPGPU parallel</u>.

- **Differentiable**, important if we want to use OT distances as **loss functions**.

- Can be **automatically differentiated**, simple iterative process, *DL*-toolboxes compatible.

- OT can become a building block in ML.

✦ Wasserstein distance is not differentiable

# *Background: OT Geometry*

**Consider** $(\Omega, {\color{green}D})$, a metric probability space. **Let** ${\color{red}\boldsymbol{\mu}}, {\color{blue}\boldsymbol{\nu}}$ be probability measures in $\mathcal{P}(\Omega)$.

- **[Monge'81]** problem: find a map ${\color{red}\boldsymbol{T}} : \Omega \to \Omega$

$$\inf_{{\color{red}\boldsymbol{T}}\#{\color{red}\boldsymbol{\mu}}={\color{blue}\boldsymbol{\nu}}} \int_{\Omega} {\color{green}D}(x, {\color{red}\boldsymbol{T}}(x)){\color{red}\boldsymbol{\mu}}(dx)$$



$x$    ${\color{red}\boldsymbol{T}}(x)$

# *Background: OT Geometry*

> **Consider** $(\Omega, D)$, a metric probability space. **Let** $\boldsymbol{\mu}, \boldsymbol{\nu}$ be probability measures in $\mathcal{P}(\Omega)$.

- **[Monge'81]** problem: find a map $\boldsymbol{T} : \Omega \to \Omega$

$$\inf_{\boldsymbol{T}\#\boldsymbol{\mu}=\boldsymbol{\nu}} \int_{\Omega} \boldsymbol{D}(x, \boldsymbol{T}(x))\boldsymbol{\mu}(dx)$$

$\boldsymbol{\delta_x}$

# [Kantorovich'42] Relaxation

- Instead of maps $T : \Omega \to \Omega$, consider probabilistic maps, i.e. **couplings** $P \in \mathcal{P}(\Omega \times \Omega)$:

$$\Pi(\mu, \nu) \overset{\mathrm{def}}{=} \{ P \in \mathcal{P}(\Omega \times \Omega) \mid \forall A, B \subset \Omega,$$
$$P(A \times \Omega) = \mu(A),$$
$$P(\Omega \times B) = \nu(B) \}$$

# [Kantorovich'42] Relaxation

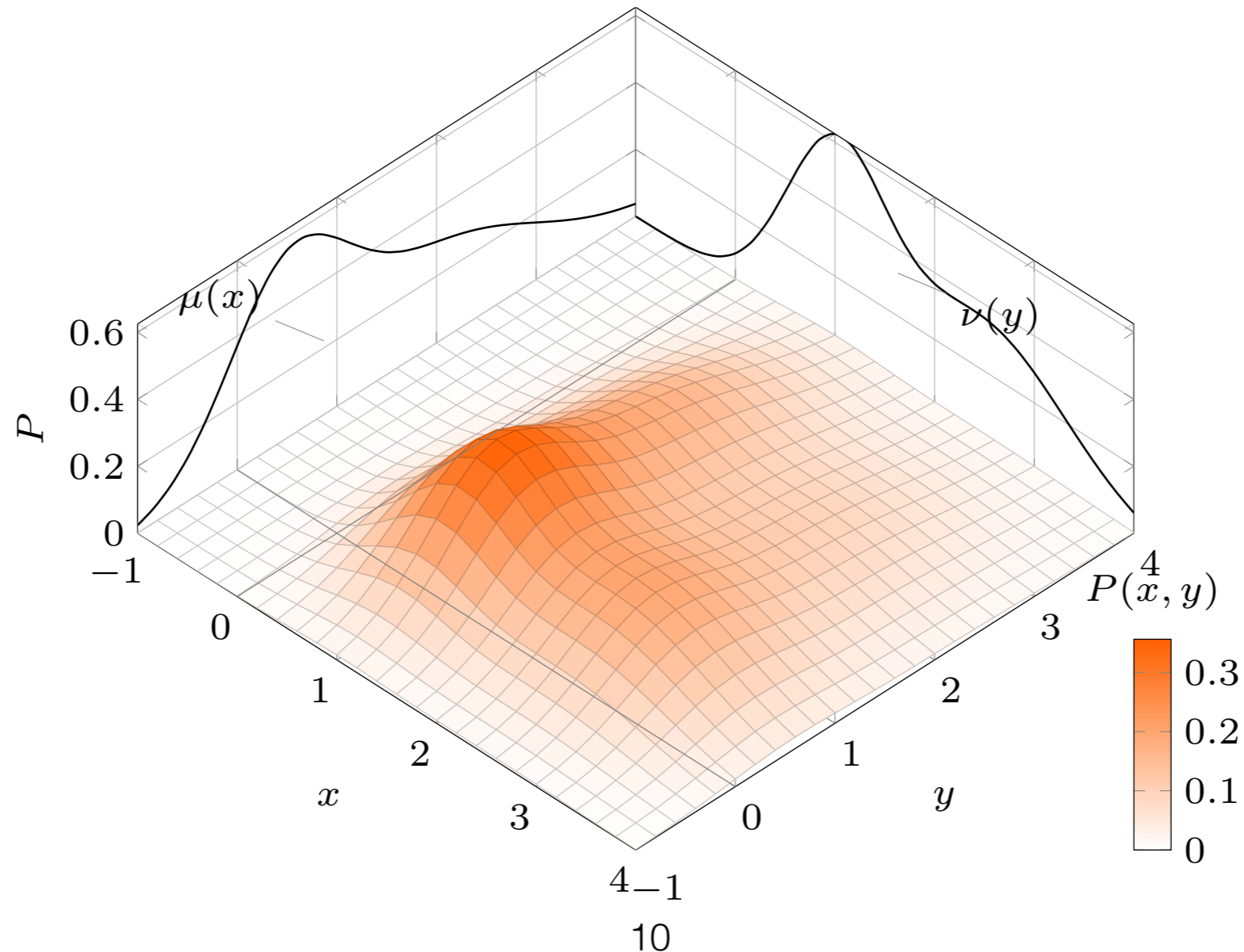$$\Pi(\textcolor{red}{\boldsymbol{\mu}}, \textcolor{blue}{\boldsymbol{\nu}}) \overset{\text{def}}{=} \{\textcolor{darkred}{\boldsymbol{P}} \in \mathcal{P}(\Omega \times \Omega) \,|\, \forall \textcolor{red}{\boldsymbol{A}}, \textcolor{blue}{\boldsymbol{B}} \subset \Omega,$$
$$\textcolor{darkred}{\boldsymbol{P}}(\textcolor{red}{\boldsymbol{A}} \times \Omega) = \textcolor{red}{\boldsymbol{\mu}}(\textcolor{red}{\boldsymbol{A}}), \textcolor{darkred}{\boldsymbol{P}}(\Omega \times \textcolor{blue}{\boldsymbol{B}}) = \textcolor{blue}{\boldsymbol{\nu}}(\textcolor{blue}{\boldsymbol{B}})\}$$

# [Kantorovich'42] Relaxation

$$\Pi(\textcolor{red}{\boldsymbol{\mu}}, \textcolor{blue}{\boldsymbol{\nu}}) \overset{\text{def}}{=} \{\textcolor{red}{\boldsymbol{P}} \in \mathcal{P}(\Omega \times \Omega) \,|\, \forall \textcolor{red}{\boldsymbol{A}}, \textcolor{blue}{\boldsymbol{B}} \subset \Omega,$$
$$\textcolor{red}{\boldsymbol{P}}(\textcolor{red}{\boldsymbol{A}} \times \Omega) = \textcolor{red}{\boldsymbol{\mu}}(\textcolor{red}{\boldsymbol{A}}), \textcolor{red}{\boldsymbol{P}}(\Omega \times \textcolor{blue}{\boldsymbol{B}}) = \textcolor{blue}{\boldsymbol{\nu}}(\textcolor{blue}{\boldsymbol{B}})\}$$

# Couplings

# Couplings

# Wasserstein Distance

**Def.** For $p \geq 1$, the $p$-Wasserstein distance between $\textcolor{red}{\boldsymbol{\mu}}, \textcolor{blue}{\boldsymbol{\nu}}$ in $\mathcal{P}(\Omega)$ is

$$W_p(\textcolor{red}{\boldsymbol{\mu}}, \textcolor{blue}{\boldsymbol{\nu}}) \stackrel{\text{def}}{=} \left( \inf_{\textcolor{brown}{\boldsymbol{P}} \in \Pi(\textcolor{red}{\boldsymbol{\mu}}, \textcolor{blue}{\boldsymbol{\nu}})} \mathbb{E}_{\textcolor{brown}{\boldsymbol{P}}}[D(X, Y)^p] \right)^{1/p}.$$

# Wasserstein between 2 Diracs



$$W_p^p(\delta_{\boldsymbol{x}}, \delta_{\boldsymbol{y}}) = \boldsymbol{D}(\boldsymbol{x}, \boldsymbol{y})$$

# Wasserstein on Uniform Measures

$$\mu = \sum_{i=1}^{n} \frac{1}{n} \delta_{x_i}$$

$$(\Omega, D)$$

$$\nu = \sum_{j=1}^{n} \frac{1}{n} \delta_{y_j}$$

# Wasserstein on Uniform Measures
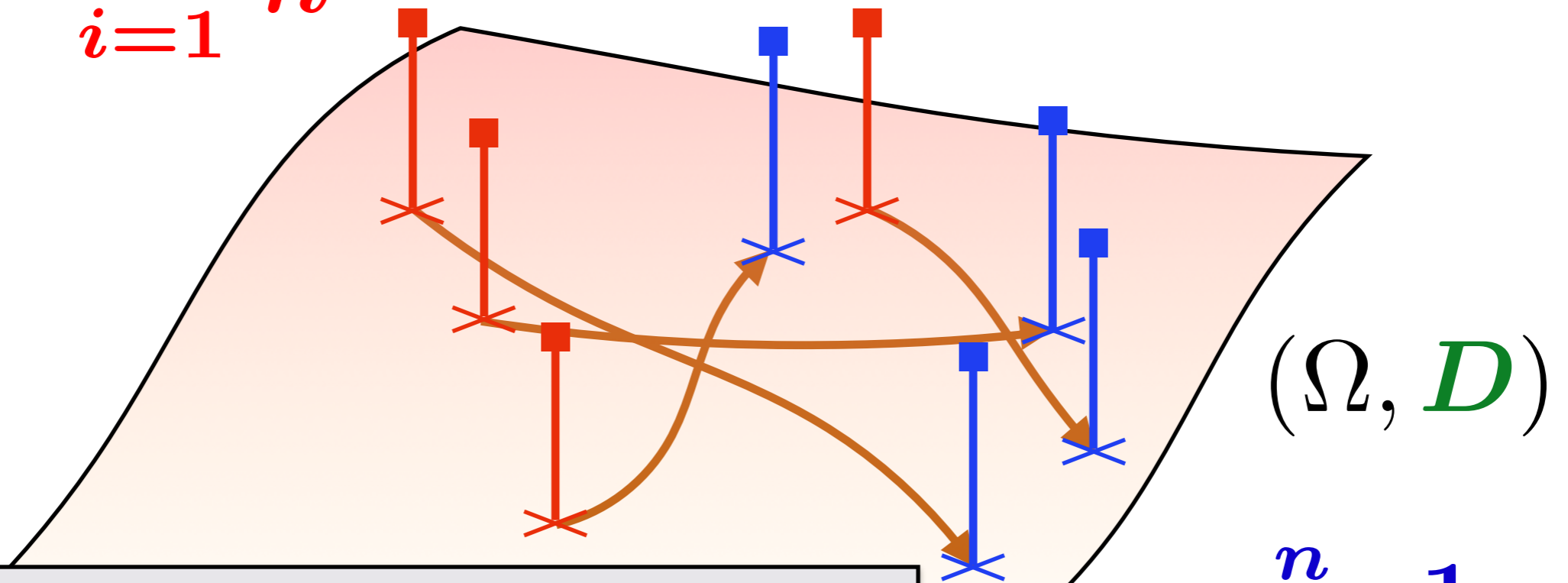
$$\mu = \sum_{i=1}^{n} \frac{1}{n} \delta_{x_i}$$



$(\Omega, D)$

$$\nu = \sum_{j=1}^{n} \frac{1}{n} \delta_{y_j}$$

$$C(\sigma) = \frac{1}{n} \sum_{i=1}^{n} D(x_i, y_{\sigma_i})^p$$

# Optimal Assignment ⊂ Wasserstein

$$\mu = \sum_{i=1}^{n} \frac{1}{n} \delta_{x_i}$$

$(\Omega, D)$

$$W_p^p(\mu, \nu) = \min_{\sigma \in S_n} C(\sigma)$$

$$\nu = \sum_{j=1}^{n} \frac{1}{n} \delta_{y_j}$$

# Wasserstein on Empirical Measures

$$\mu = \sum_{i=1}^{n} a_i \delta_{x_i}$$



$$(\Omega, D)$$

$$\nu = \sum_{j=1}^{m} b_j \delta_{y_j}$$

# Wasserstein on Empirical Measures

**Consider** $\boldsymbol{\mu} = \sum_{i=1}^{n} \boldsymbol{a_i} \boldsymbol{\delta_{x_i}}$ and $\boldsymbol{\nu} = \sum_{j=1}^{m} \boldsymbol{b_j} \boldsymbol{\delta_{y_j}}$.

$M_{\boldsymbol{XY}} \stackrel{\text{def}}{=} [D(\boldsymbol{x_i}, \boldsymbol{y_j})^p]_{ij}$

$U(\boldsymbol{a}, \boldsymbol{b}) \stackrel{\text{def}}{=} \{\boldsymbol{P} \in \mathbb{R}_+^{n \times m} | \boldsymbol{P}\mathbf{1}_m = \boldsymbol{a}, \boldsymbol{P}^T\mathbf{1}_n = \boldsymbol{b}\}$

# Wasserstein on Empirical Measures

Consider $\mu = \sum_{i=1}^{n} a_i \delta_{x_i}$ and $\nu = \sum_{j=1}^{m} b_j \delta_{y_j}$.

$$M_{XY} \stackrel{\text{def}}{=} [D(x_i, y_j)^p]_{ij}$$

$$U(a, b) \stackrel{\text{def}}{=} \{P \in \mathbb{R}_+^{n \times m} \mid P\mathbf{1}_m = a, P^T\mathbf{1}_n = b\}$$

$$
\begin{array}{c}
\quad y_1 \quad\quad \ldots \quad\quad y_m \\
\begin{array}{c} x_1 \\ \vdots \\ x_n \end{array}
\left[
\begin{array}{ccc}
\cdot & \cdot & \cdot \\
\cdot & D(x_i, y_j)^p & \cdot \\
\cdot & \cdot & \cdot
\end{array}
\right]
\end{array}
$$

$$
\begin{array}{c}
\quad b_1 \quad\quad \ldots \quad\quad b_m \\
\begin{array}{c} a_1 \\ \vdots \\ a_n \end{array}
\left[
\begin{array}{ccc}
\vdots & \vdots & \vdots \\
\vdots & P^T\mathbf{1}_n = b & \vdots \\
\vdots & \vdots & \vdots
\end{array}
\right]
\end{array}
$$

18

# Wasserstein on Empirical Measures

**Consider** $\boldsymbol{\mu} = \sum_{i=1}^{n} a_i \boldsymbol{\delta}_{x_i}$ and $\boldsymbol{\nu} = \sum_{j=1}^{m} b_j \boldsymbol{\delta}_{y_j}$.

$$M_{\boldsymbol{XY}} \overset{\text{def}}{=} [D(\boldsymbol{x}_i, \boldsymbol{y}_j)^p]_{ij}$$

$$U(\boldsymbol{a}, \boldsymbol{b}) \overset{\text{def}}{=} \{\boldsymbol{P} \in \mathbb{R}_+^{n \times m} | \boldsymbol{P}\mathbf{1}_m = \boldsymbol{a}, \boldsymbol{P}^T \mathbf{1}_n = \boldsymbol{b}\}$$

**Def.** Optimal Transport Problem

$$W_p^p(\boldsymbol{\mu}, \boldsymbol{\nu}) = \min_{\boldsymbol{P} \in U(\boldsymbol{a}, \boldsymbol{b})} \langle \boldsymbol{P}, M_{\boldsymbol{XY}} \rangle$$

# Discrete OT Problem

# Discrete OT Problem

# Discrete OT Problem



$M_{\boldsymbol{XY}}$

$U(\boldsymbol{a}, \boldsymbol{b})$

$P^\star$

**Def.** Dual OT problem

$$W_p^p(\boldsymbol{\mu}, \boldsymbol{\nu}) = \max_{\substack{\boldsymbol{\alpha} \in \mathbb{R}^n, \boldsymbol{\beta} \in \mathbb{R}^m \\ \boldsymbol{\alpha_i} + \boldsymbol{\beta_j} \leq D(\boldsymbol{x_i}, \boldsymbol{y_j})^p}} \alpha^T \boldsymbol{a} + \beta^T \boldsymbol{b}$$

# Discrete OT Problem



$M_{\textbf{X}\textbf{Y}}$

$P^{\star}$

$U(\textbf{\textit{a}}, \textbf{\textit{b}})$

*network flow solver
used in practice.*
$O(n^3 \log(n))$

**Note:** flow/PDE formulations **[Beckman'61]/[Benamou'98]** can be used for *p=1/p=2* for a sparse-graph metric/Euclidean metric.

# Discrete OT Problem

```c
/*
    emd.c

    Last update: 3/14/98

    An implementation of the Earth Movers Distance.
    Based of the solution for the Transportation problem as described in
    "Introduction to Mathematical Programming" by F. S. Hillier and
    G. J. Lieberman, McGraw-Hill, 1990.

    Copyright (C) 1998 Yossi Rubner
    Computer Science Department, Stanford University
    E-Mail: rubner@cs.stanford.edu    URL: http://vision.stanford.edu/~rubner
*/

/*#include <stdio.h>
#include <stdlib.h>*/
#include <math.h>

#include "emd.h"

#define DEBUG_LEVEL 0
/*
 DEBUG_LEVEL:
    0 = NO MESSAGES
    1 = PRINT THE NUMBER OF ITERATIONS AND THE FINAL RESULT
    2 = PRINT THE RESULT AFTER EVERY ITERATION
    3 = PRINT ALSO THE FLOW AFTER EVERY ITERATION
    4 = PRINT A LOT OF INFORMATION (PROBABLY USEFUL ONLY FOR THE AUTHOR)
*/


#define MAX_SIG_SIZE1  (MAX_SIG_SIZE+1)  /* FOR THE POSSIBLE DUMMY FEATURE */

/* NEW TYPES DEFINITION */

/* node1_t IS USED FOR SINGLE-LINKED LISTS */
typedef struct node1_t {
  int i;
  double val;
  struct node1_t *Next;
} node1_t;

/* node1_t IS USED FOR DOUBLE-LINKED LISTS */
typedef struct node2_t {
  int i, j;
  double val;
  struct node2_t *NextC;                /* NEXT COLUMN */
  struct node2_t *NextR;                /* NEXT ROW */
} node2_t;




/* GLOBAL VARIABLE DECLARATION */
static int _n1, _n2;                          /* SIGNATURES SIZES */
static float _C[MAX_SIG_SIZE1][MAX_SIG_SIZE1];/* THE COST MATRIX */
static node2_t _X[MAX_SIG_SIZE1*2];           /* THE BASIC VARIABLES VECTOR */
```

# Discrete OT Problem

```c
/*

    emd.c

    Last update: 3/14/98

    An implementation of the Earth Movers Distance.
    Based of the solution for the Transportation problem as described in
    "Introduction to Mathematical Programming" by F. S. Hillier and
    G. J. Lieberman, McGraw-Hill, 1990.

    Copyright (C) 1998 Yossi Rubner
    Computer Science Department, Stanford University
    E-Mail: rubner@cs.stanford.edu    URL: http://vision.stanford.edu/~rubner
*/

/*#include <stdio.h>
#include <stdlib.h>*/
#include <math.h>

#include "emd.h"

#define DEBUG_LEVEL 0
/*
 DEBUG_LEVEL:
    0 = NO MESSAGES
    1 = PRINT THE NUMBER OF ITERATIONS AND THE FINAL RESULT
    2 = PRINT THE RESULT AFTER EVERY ITERATION
    3 = PRINT ALSO THE FLOW AFTER EVERY ITERATION
    4 = PRINT A LOT OF INFORMATION (PROBABLY USEFUL ONLY FOR THE AUTHOR)
*/


#define MAX_SIG_SIZE1  (MAX_SIG_SIZE+1)  /* FOR THE POSIBLE DUMMY FEATURE */

/* NEW TYPES DEFINITION */

/* node1_t IS USED FOR SINGLE-LINKED LISTS */
typedef struct node1_t {
  int i;
  double val;
  struct node1_t *Next;
} node1_t;

/* node1_t IS USED FOR DOUBLE-LINKED LISTS */
typedef struct node2_t {
  int i, j;
  double val;
  struct node2_t *NextC;               /* NEXT COLUMN */
  struct node2_t *NextR;               /* NEXT ROW */
} node2_t;



/* GLOBAL VARIABLE DECLARATION */
static int _n1, _n2;                          /* SIGNATURES SIZES */
static float _C[MAX_SIG_SIZE1][MAX_SIG_SIZE1];/* THE COST MATRIX */
static node2_t _X[MAX_SIG_SIZE1*2];           /* THE BASIC VARIABLES VECTOR */
```
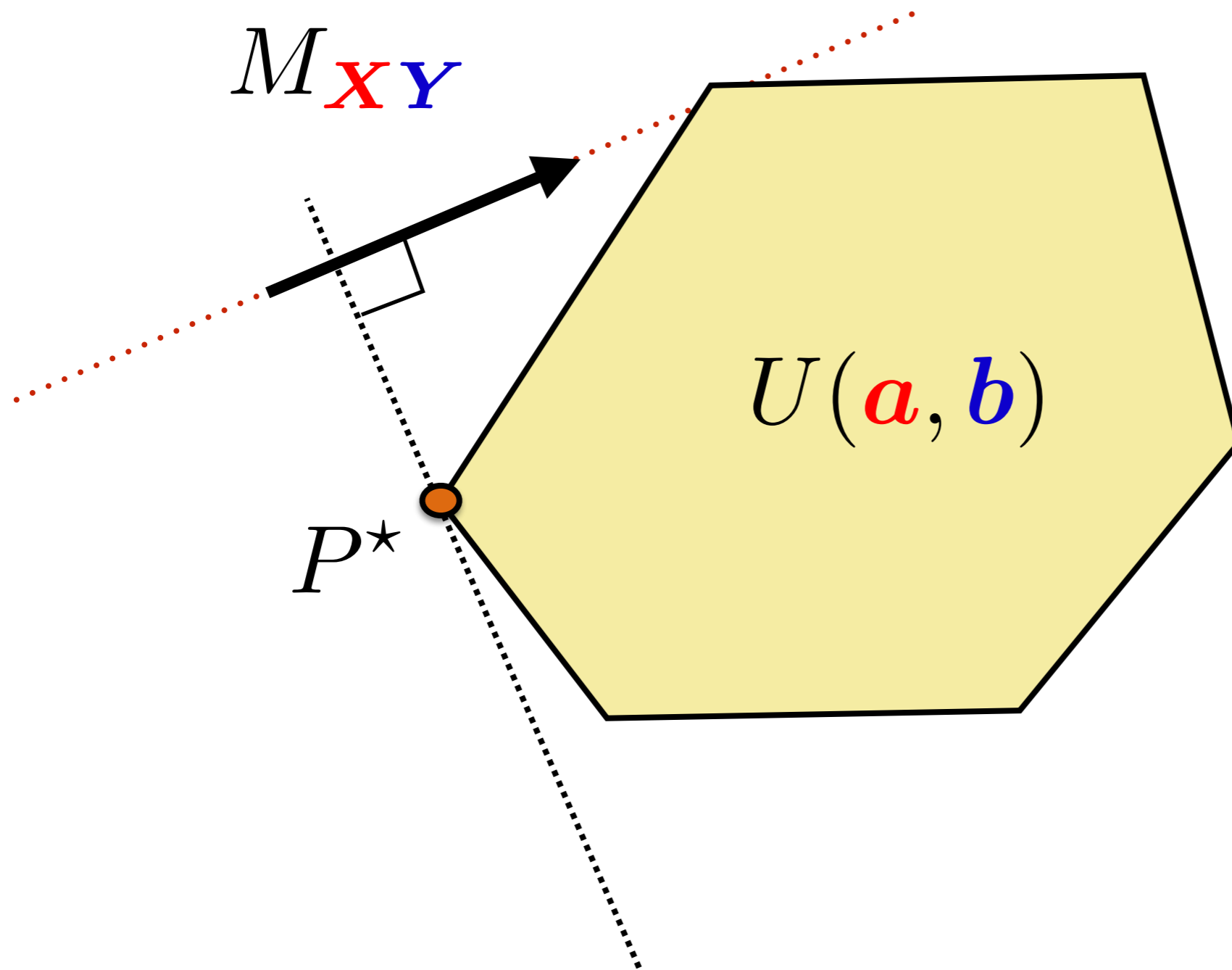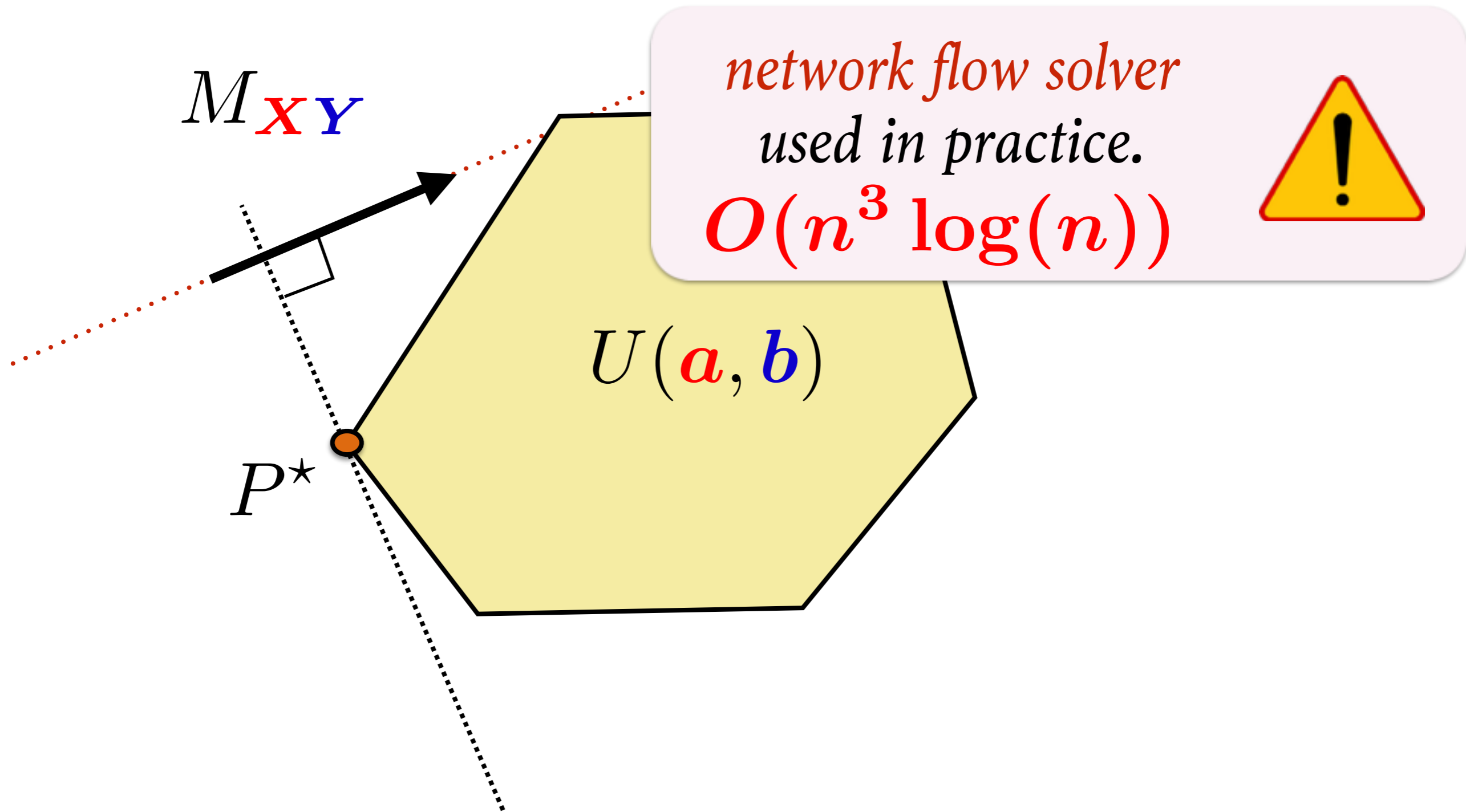
# Discrete OT Problem

# Discrete OT Problem

# Discrete OT Problem



$M_{\textcolor{red}{X}\textcolor{blue}{Y}}$

$P^\star$

$U(\textcolor{red}{a}, \textcolor{blue}{b})$

*network flow solver used in practice.*
$O(n^3 \log(n))$

# Discrete OT Problem

$M_{\mathbf{X}\mathbf{Y}}$

$U(\mathbf{a}, \mathbf{b})$

*network flow solver used in practice.*
$$O(n^3 \log(n))$$

*Solution $P^\star$ unstable and not always unique.*

$P^\star$

# Discrete OT Problem



$M_{\mathbf{X}\mathbf{Y}}$

$U(\mathbf{a}, \mathbf{b})$

$\{P^\star\}$

*network flow solver used in practice.*
$$O(n^3 \log(n))$$

*Solution $P^\star$ unstable and not always unique.*

# Discrete OT Problem



$M_{\textcolor{red}{X}\textcolor{blue}{Y}}$

$U(\textcolor{red}{a}, \textcolor{blue}{b})$

$\{P^\star\}$

*network flow solver used in practice.*
$O(n^3 \log(n))$

*Solution $P^\star$ unstable and not always unique.*

# Discrete OT Problem



$M_{\textbf{X}\textbf{Y}}$

$U(\textbf{a}, \textbf{b})$

*network flow solver used in practice.*
$O(n^3 \log(n))$

*Solution $P^\star$ unstable and not always unique.*

$P^\star$

# Discrete OT Problem



$M_{\textbf{X}\textbf{Y}}$

$U(\textbf{a}, \textbf{b})$

*network flow solver*
*used in practice.*
$O(n^3 \log(n))$

*Solution $P^\star$ unstable*
*and not always unique.*

$P^\star$

$W_p^p(\boldsymbol{\mu}, \boldsymbol{\nu})$ not differentiable.

# Entropic Regularization [**Wilson'62**]

**Def.** Regularized Wasserstein, $\gamma \geq 0$

$$W_\gamma(\textcolor{red}{\boldsymbol{\mu}}, \textcolor{blue}{\boldsymbol{\nu}}) \overset{\text{def}}{=} \min_{\textcolor{red}{\boldsymbol{P}} \in U(\textcolor{red}{\boldsymbol{a}}, \textcolor{blue}{\boldsymbol{b}})} \langle \textcolor{red}{\boldsymbol{P}}, M_{\textcolor{red}{\boldsymbol{X}}\textcolor{blue}{\boldsymbol{Y}}} \rangle - \gamma E(\textcolor{red}{\boldsymbol{P}})$$

$$E(P) \overset{\text{def}}{=} - \sum_{i,j=1}^{nm} P_{ij}(\log P_{ij})$$

**Note: Unique** optimal solution because of strong concavity of Entropy

# Entropic Regularization [**Wilson'62**]

**Def.** Regularized Wasserstein, $\gamma \geq 0$

$$W_\gamma(\textcolor{red}{\boldsymbol{\mu}}, \textcolor{blue}{\boldsymbol{\nu}}) \overset{\text{def}}{=} \min_{\textcolor{brown}{\boldsymbol{P}} \in U(\textcolor{red}{\boldsymbol{a}}, \textcolor{blue}{\boldsymbol{b}})} \langle \textcolor{brown}{\boldsymbol{P}}, M_{\textcolor{red}{\boldsymbol{X}}\textcolor{blue}{\boldsymbol{Y}}} \rangle - \gamma E(\textcolor{brown}{\boldsymbol{P}})$$



**Note: Unique** optimal solution because of strong concavity of Entropy

# Fast & Scalable Algorithm

**Prop.** If $P_\gamma \overset{\text{def}}{=} \underset{P \in U(a,b)}{\text{argmin}} \langle P, M_{XY} \rangle - \gamma E(P)$

then $\exists! u \in \mathbb{R}_+^n, v \in \mathbb{R}_+^m$, such that

$P_\gamma = \text{diag}(u) K \text{diag}(v), \quad K \overset{\text{def}}{=} e^{-M_{XY}/\gamma}$

# Fast & Scalable Algorithm

**Prop.** If $P_\gamma \overset{\text{def}}{=} \underset{P \in U(a,b)}{\text{argmin}} \langle P, M_{XY} \rangle - \gamma E(P)$

then $\exists! u \in \mathbb{R}^n_+, v \in \mathbb{R}^m_+$, such that

$$P_\gamma = \textbf{diag}(u) K \textbf{diag}(v), \quad K \overset{\text{def}}{=} e^{-M_{XY}/\gamma}$$

$$L(P, \alpha, \beta) = \sum_{ij} P_{ij} M_{ij} + \gamma P_{ij} \log P_{ij} + \alpha^T(P\mathbf{1} - a) + \beta^T(P^T\mathbf{1} - b)$$

$$\partial L / \partial P_{ij} = M_{ij} + \gamma(\log P_{ij} + 1) + \alpha_i + \beta_j$$

$$(\partial L / \partial P_{ij} = 0) \Rightarrow P_{ij} = e^{\frac{\alpha_i}{\gamma} + \frac{1}{2}} \ e^{-\frac{M_{ij}}{\gamma}} \ e^{\frac{\beta_j}{\gamma} + \frac{1}{2}} = u_i \ K_{ij} v_j$$

# Fast & Scalable Algorithm

**Prop.** If $P_\gamma \overset{\text{def}}{=} \underset{\boldsymbol{P} \in U(\boldsymbol{a}, \boldsymbol{b})}{\text{argmin}} \langle \boldsymbol{P}, M_{\boldsymbol{XY}} \rangle - \gamma E(\boldsymbol{P})$

then $\exists! \boldsymbol{u} \in \mathbb{R}_+^n, \boldsymbol{v} \in \mathbb{R}_+^m$, such that

$$P_\gamma = \mathbf{diag}(\boldsymbol{u}) K \mathbf{diag}(\boldsymbol{v}), \quad K \overset{\text{def}}{=} e^{-M_{\boldsymbol{XY}}/\gamma}$$

- **[Sinkhorn'64]** fixed-point iterations for $(\boldsymbol{u}, \boldsymbol{v})$

$$\boldsymbol{u} \leftarrow \boldsymbol{a}/K\boldsymbol{v}, \quad \boldsymbol{v} \leftarrow \boldsymbol{b}/K^T\boldsymbol{u}$$

- $O(nm)$ complexity, GPGPU parallel [C'13] .
- $O(n^{d+1})$ if $\Omega = \{1, \dots, n\}^d$ and $D^p$ separable.

[S..C..'15]

26

# Very Fast EMD Approx. Solver



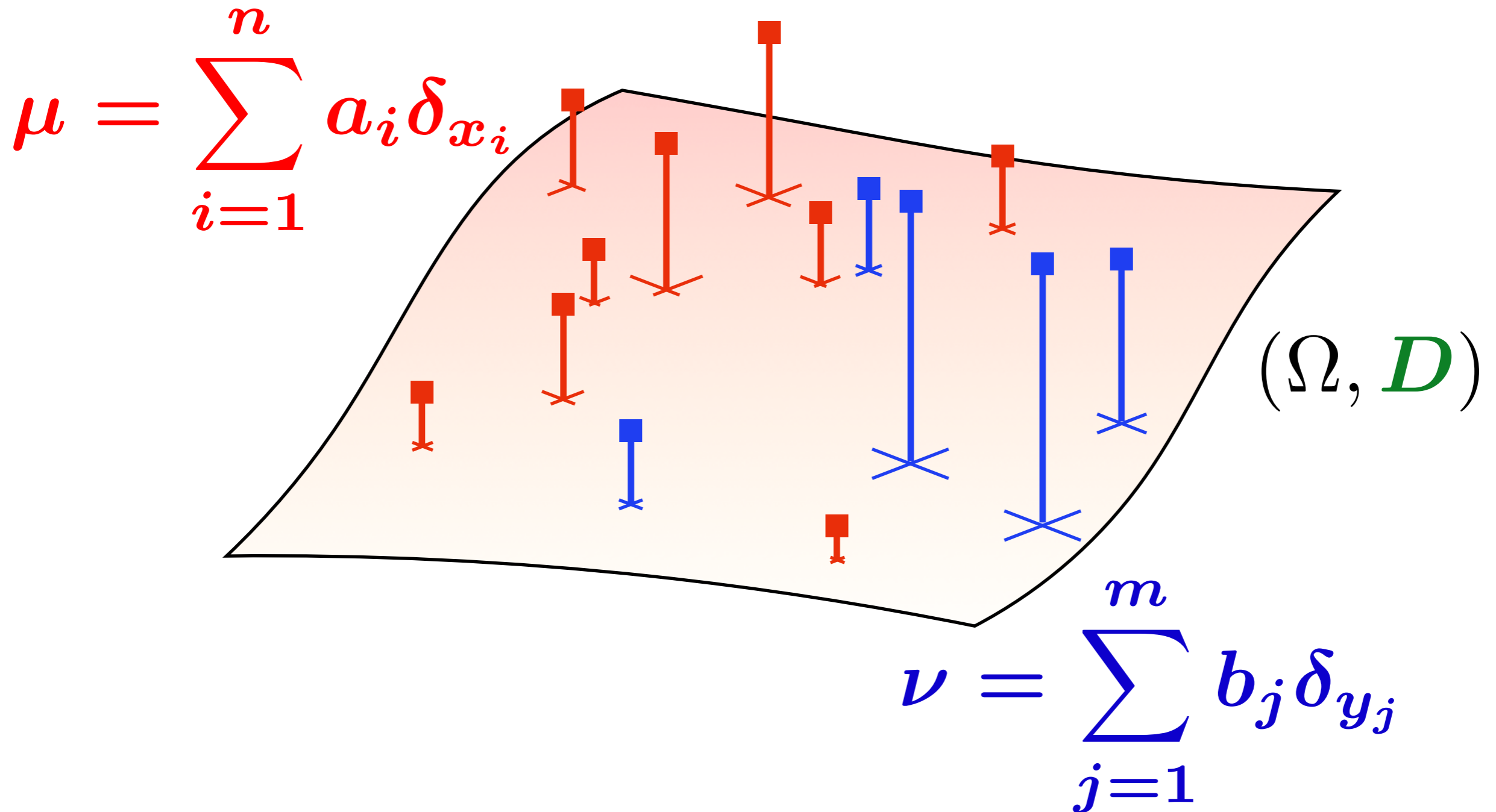**Note.** $(\Omega, D)$ is a random graph with shortest path metric, histograms sampled uniformly on simplex, Sinkhorn tolerance $10^{-2}$.

# Regularization ⤳ *Differentiability*

$$W_\gamma((\textcolor{red}{\boldsymbol{a}}, \textcolor{red}{\boldsymbol{X}}), (\textcolor{blue}{\boldsymbol{b}}, \textcolor{blue}{\boldsymbol{Y}})) = \min_{\textcolor{red}{\boldsymbol{P}} \in U(\textcolor{red}{\boldsymbol{a}}, \textcolor{blue}{\boldsymbol{b}})} \langle \textcolor{red}{\boldsymbol{P}}, M_{\textcolor{red}{\boldsymbol{X}}\textcolor{blue}{\boldsymbol{Y}}} \rangle - \gamma E(\textcolor{red}{\boldsymbol{P}})$$

$$\textcolor{red}{\mu = \sum_{i=1}^{n} a_i \delta_{x_i}}$$

$$(\Omega, \textcolor{green}{\boldsymbol{D}})$$

$$\textcolor{blue}{\nu = \sum_{j=1}^{m} b_j \delta_{y_j}}$$



28

# Regularization ⤳ *Differentiability*

$$W_\gamma((\textcolor{red}{a + \Delta a, X}), (\textcolor{blue}{b, Y})) = W_\gamma((\textcolor{red}{a, X}), (\textcolor{blue}{b, Y})) + ??$$

$$\textcolor{red}{\mu = \sum_{i=1}^{n} a_i \delta_{x_i}}$$

$$(\Omega, \textcolor{green}{D})$$

$$\textcolor{blue}{\nu = \sum_{j=1}^{m} b_j \delta_{y_j}}$$

# Regularization ⤳ *Differentiability*

$$W_\gamma\big((a + \Delta a, X), (b, Y)\big) = W_\gamma\big((a, X), (b, Y)\big) + ??$$

$$\mu = \sum_{i=1}^{n} a_i \delta_{x_i}$$



$(\Omega, D)$

$$a \leftarrow a + \Delta a$$

$$\nu = \sum_{j=1}^{m} b_j \delta_{y_j}$$

# Regularization ⤳ *Differentiability*

$$W_\gamma((\textcolor{red}{a, X + \Delta X}), \textcolor{blue}{(b, Y)}) = W_\gamma(\textcolor{red}{(a, X)}, \textcolor{blue}{(b, Y)}) + ??$$

$$\textcolor{red}{\mu = \sum_{i=1}^{n} a_i \delta_{x_i}}$$

$(\Omega, \textcolor{green}{D})$

$$\textcolor{blue}{\nu = \sum_{j=1}^{m} b_j \delta_{y_j}}$$

# Regularization ⤳ *Differentiability*

$$W_\gamma((\textcolor{red}{\boldsymbol{a}, \boldsymbol{X} + \boldsymbol{\Delta X}}), (\textcolor{blue}{\boldsymbol{b}, \boldsymbol{Y}})) = W_\gamma((\textcolor{red}{\boldsymbol{a}, \boldsymbol{X}}), (\textcolor{blue}{\boldsymbol{b}, \boldsymbol{Y}})) + ??$$

$$\textcolor{red}{\boldsymbol{\mu} = \sum_{i=1}^{n} \boldsymbol{a_i} \boldsymbol{\delta_{x_i}}}$$



$$(\Omega, \textcolor{green}{\boldsymbol{D}})$$

$$\textcolor{red}{\boldsymbol{X} \leftarrow \boldsymbol{X} + \boldsymbol{\Delta X}}$$

$$\textcolor{blue}{\boldsymbol{\nu} = \sum_{j=1}^{m} \boldsymbol{b_j} \boldsymbol{\delta_{y_j}}}$$

# Crucial for "min *data* + $W$" problems

- Quantization, $k$-means problem [**Lloyd'82**]

$$\min_{\substack{\boldsymbol{\mu} \in \mathcal{P}(\mathbb{R}^d) \\ |\operatorname{supp} \boldsymbol{\mu}| = k}} W_2^2(\boldsymbol{\mu}, \boldsymbol{\nu}_{\text{data}})$$

- [**McCann'95**] Interpolant

$$\min_{\boldsymbol{\mu} \in \mathcal{P}(\Omega)} (1 - t) W_2^2(\boldsymbol{\mu}, \boldsymbol{\nu_1}) + t W_2^2(\boldsymbol{\mu}, \boldsymbol{\nu_2})$$

- [**JKO'98**] PDE's as gradient flows in $(\mathcal{P}(\Omega), W)$.

$$\mu_{t+1} = \operatorname*{argmin}_{\boldsymbol{\mu} \in \mathcal{P}(\Omega)} J(\boldsymbol{\mu}) + \lambda_t W_p^p(\boldsymbol{\mu}, \mu_t)$$

# Crucial for "min *data* + *W* " problems

Any (ML) problem involving a **KL** or **L2** loss between (parameterized) histograms or probabilility measures can be easily *Wasserstein-ized* if we can differentiate *W* efficiently.

# 1. Differentiability of Regularized OT

**Def.** Dual regularized OT Problem

$$W_\gamma(\boldsymbol{\mu}, \boldsymbol{\nu}) = \max_{\alpha, \beta} \alpha^T \boldsymbol{a} + \beta^T \boldsymbol{b} - \frac{1}{\gamma}(e^{\alpha/\gamma})^T K e^{\beta/\gamma}$$

**Prop.** $W_\gamma(\boldsymbol{\mu}, \boldsymbol{\nu})$ is                    **[CD'14]**

1. convex w.r.t. $\boldsymbol{a}$ (Danskin),
$$\nabla_{\boldsymbol{a}} W_\gamma = \alpha^\star = \gamma \log(\boldsymbol{u}).$$

2. decreased, when $p = 2, \Omega = \mathbb{R}^d$, using
$$X \leftarrow Y P_\gamma^T \mathbf{D}(\boldsymbol{a}^{-1}).$$

# 2. Duality for Regularized OT's

**Prop.** Writing $H_{\boldsymbol{\nu}} : \boldsymbol{a} \mapsto W_{\gamma}(\boldsymbol{\mu}, \boldsymbol{\nu})$, **[CP'16]**

**1.** $H_{\boldsymbol{\nu}}$ has simple Legendre transform:

$$H_{\boldsymbol{\nu}}^{*} : \boldsymbol{g} \in \mathbb{R}^n \mapsto \gamma \left( E(\boldsymbol{b}) + \boldsymbol{b}^T \log(K e^{\boldsymbol{g}/\gamma}) \right)$$

**2.** If $A \in \mathbb{R}^{n \times d}$, $f$ convex on $\mathbb{R}^d$,

$$\min_{\boldsymbol{a} \in \Sigma_n} H_{\boldsymbol{\nu}}(\boldsymbol{a}) + f(A\boldsymbol{a}) = \max_{\boldsymbol{g} \in \mathbb{R}^d} -H_{\boldsymbol{\nu}}^{*}(A^T \boldsymbol{g}) - f^{*}(-\boldsymbol{g})$$

# 3. Stochastic Formulation

$$
\begin{aligned}
W_\gamma(\boldsymbol{\mu}, \boldsymbol{\nu}) &= \max_{\alpha, \beta} \alpha^T \boldsymbol{a} + \beta^T \boldsymbol{b} - \frac{1}{\gamma}(e^{\alpha/\gamma})^T K e^{\beta/\gamma} \\
&= \max_{\boldsymbol{\alpha}} \boldsymbol{\alpha}^T \boldsymbol{a} - \gamma(\log K e^{\boldsymbol{\alpha}/\gamma})^T \boldsymbol{b} \\
&= \max_{\boldsymbol{\alpha}} \sum_{j=1}^{m} \boldsymbol{b_j}\left(\boldsymbol{\alpha}^T \boldsymbol{a} - \gamma \log K_{.j}^T e^{\boldsymbol{\alpha}/\gamma}\right) \\
&= \max_{\boldsymbol{\alpha}} \sum_{j=1}^{m} f_j(\boldsymbol{\alpha})
\end{aligned}
$$

- **[GCPB'16]** shows how incremental gradient methods can be used to scale this further.

# 4. Algorithmic Formulation

**Def.** For $L \geq 1$, define

$$W_L(\boldsymbol{\mu}, \boldsymbol{\nu}) \overset{\text{def}}{=} \langle P_L, M_{\boldsymbol{XY}} \rangle,$$

where $P_L \overset{\text{def}}{=} \mathbf{diag}(\boldsymbol{u_L}) K \mathbf{diag}(\boldsymbol{v_L}),$

$\boldsymbol{v_0} = \mathbf{1}_m; l \geq 0, \boldsymbol{u_l} \overset{\text{def}}{=} \boldsymbol{a}/K\boldsymbol{v_l}, \boldsymbol{v_{l+1}} \overset{\text{def}}{=} \boldsymbol{b}/K^T\boldsymbol{u_l}.$

**Prop.** $\frac{\partial W_L}{\partial \boldsymbol{X}}, \frac{\partial W_L}{\partial \boldsymbol{a}}$ can be computed recursively, in $O(L)$ kernel $K \times$ vector products.

# Algorithmic Formulation of Reg. OT

**Example**: Differentiability w.r.t. $a$

$$\left(\frac{\partial \boldsymbol{v_0}}{\partial a}\right)^T = \boldsymbol{0}_{m \times n},$$

$$\left(\frac{\partial \boldsymbol{u_l}}{\partial a}\right)^T \boldsymbol{x} = \frac{\boldsymbol{x}}{K \boldsymbol{v_l}} - \left(\frac{\partial \boldsymbol{v_l}}{\partial a}\right)^T K^T \frac{\boldsymbol{x} \circ a}{(K \boldsymbol{v_l})^2},$$

$$\left(\frac{\partial \boldsymbol{v_{l+1}}}{\partial a}\right)^T \boldsymbol{y} = -\left(\frac{\partial \boldsymbol{u_l}}{\partial a}\right)^T K \frac{\boldsymbol{y} \circ b}{(K^T \boldsymbol{u_l})^2}.$$

# Algorithmic Formulation of Reg. OT

**Example**: Differentiability w.r.t. $a$

$$N = K \circ M_{\boldsymbol{XY}}$$

$$\nabla_{\boldsymbol{a}} W_L(\boldsymbol{\mu}, \boldsymbol{\nu}) = \left(\frac{\partial \boldsymbol{u_L}}{\partial a}\right)^T N \boldsymbol{v_L} + \left(\frac{\partial \boldsymbol{v_L}}{\partial a}\right)^T N^T \boldsymbol{u_L}$$

```matlab
function [d,grad_a,grad_b,hess_a,hess_b] = sinkhornObjGradHess(a,b,K,M,niter)

u_update = @(v,a) a./(K*v);
v_update = @(u,b) b./(K'*u);


% DuDa = @(eps,dvda,a,v)  (eps./(K*v))- (a./((K*v).^2)).*(K*dvda(eps));
%
% DvDa = @(eps,duda,b,u)  -(b./((K'*u).^2)).*(K'*duda(eps));
%
% DuDb = @(eps,dvdb,a,v)  -(a./((K*v).^2)).*(K*dvdb(eps));
%
% DvDb = @(eps,dudb,b,u)  (eps./(K'*u))-(b./((K'*u).^2)).*(K'*dudb(eps));



DuDat = @(x,dvdat,a,v)  bsxfun(@rdivide,x,K*v)... (x./(K*v))
    -dvdat(K'*( bsxfun(@times,x,(a./((K*v).^2))))));...-dvdat(K'*( (a./((K*v).^2)).*x));


DvDat = @(x,dudat,b,u)  -dudat(K*(bsxfun(@times,x,(b./((K'*u).^2))))); ...(b./((K'*u).^2)).*x))


JDuDat= @(x,Jdvdat,dvdat,a,v) -diag((x'*dvdat(K'))'./((K*v).^2)) ...(K*dvda(x))
    - Jdvdat(x)*K'*diag(a./((K*v).^2))...
    - dvdat(K'* ...
    ( diag(a.*( (-2*(x'*dvdat(K'))')./((K*v).^3)))+...
    diag(x./((K*v).^2))  ));          %1

JDvDat = @(x,Jdudat,dudat,b,u) ...
    -Jdudat(x)*K*diag(b./((K'*u).^2))...
    - dudat(K)* ( ...
    diag(b.*( (-2* (x'*dudat(K))')./((K'*u).^3)))) ;...
```

```matlab
DuDbt = @(x,dvdbt,a,v)  -dvdbt(K'*(bsxfun(@times,x,(a./((K*v).^2))))); ...(a./((K*v).^2)).*x));

DvDbt = @(x,dudbt,b,u)  bsxfun(@rdivide,x,K'*u) ... (x./(K'*u))...
    -dudbt(K*( bsxfun(@times,x,(b./((K'*u).^2))))));...( b./((K'*u).^2)) .*x));




JDvDbt= @(x,Jdudbt,dudbt,b,u) -diag((x'*dudbt(K))'./((K'*u).^2)) ...  (K'*dudb(x))
    - Jdudbt(x)*K*diag(b./((K'*u).^2))...
    - dudbt(K)* ( ...
    diag(b.*( (-2*(x'*dudbt(K))')./((K'*u).^3)))+...
    diag(x./((K'*u).^2))  ) ;

JDuDbt = @(x,Jdvdbt,dvdbt,a,v) ...
    -Jdvdbt(x)*K'*diag(a./((K*v).^2))...
    - dvdbt(K')* ( ...
    diag(a.*( (-2* (x'*dvdbt(K'))')./((K*v).^3)))) ;
```

```
n=size(a,1);
m=size(b,1);

DVDAT= @(eps) zeros(n,size(eps,2));
DVDBT= @(eps) zeros(m,size(eps,2));

JDVDAT= @(eps) zeros(n,m);
JDVDBT= @(eps) zeros(m,m);

v=ones(m,size(b,2));

for j=1:niter,
    u=u_update(v,a);
    DUDAT = @(x) DuDat(x,DVDAT,a,v);
    DUDBT = @(x) DuDbt(x,DVDBT,a,v);

    if nargout>3
        JDUDAT = @(x) JDuDat(x,JDVDAT,DVDAT,a,v);
        JDUDBT = @(x) JDuDbt(x,JDVDBT,DVDBT,a,v);
    end


    v=v_update(u,b);
    DVDAT = @(x) DvDat(x,DUDAT,b,u);
    DVDBT = @(x) DvDbt(x,DUDBT,b,u);

    if nargout>3
        JDVDAT = @(x) JDvDat(x,JDUDAT,DUDAT,b,u);
        JDVDBT = @(x) JDvDbt(x,JDUDBT,DUDBT,b,u);
    end
end
end
```

```
U=K.*M;
d=diag(u'*U*v);

grad_a=(DUDAT(U*v)+DVDAT(U'*u));
grad_b=(DUDBT(U*v)+DVDBT(U'*u));


if nargout>3
    hess_a= @(eps) JDUDAT(eps)*(U*v)+DUDAT((eps'*DVDAT(U'))')+...
        JDVDAT(eps)*(U'*u)+DVDAT((eps'*DUDAT(U))');
end
```

# Thanks to these tricks...

- **[Agueh'11]** Barycenters **[CD'14][BCCNP'15] [GCP'15][S..C..'15]**

- **[Burger'12]** TV gradient flow using duality **[CP'16]**

- Dictionary Learning / Latent Factors **[RCP'16]**

- **[Bigot'15]** *W*-PCA **[SC'15]**

- Density fitting / parameter estimation **[MMC'16]**

- Inverse problems / Wasserstein regression **[BPC'16]**

# *Wasserstein Barycenters*

$$\min_{\boldsymbol{\mu} \in \mathcal{P}(\Omega)} \sum_{i=1}^{N} \lambda_i W_p^p(\boldsymbol{\mu}, \boldsymbol{\nu_i})$$



$\mathcal{P}(\Omega)$

$\boldsymbol{\nu_1}$

Wasserstein
Barycenter
[**Agueh'11**]

$\boldsymbol{\nu_2}$

$\boldsymbol{\nu_3}$

# Multimarginal Formulation

- Exact solution $(W_2)$ using MM-OT. **[Agueh'11]**

# Multimarginal Formulation

- Exact solution $(W_2)$ using MM-OT. **[Agueh'11]**



If $\left|\operatorname{supp}\boldsymbol{\nu_i}\right| = \boldsymbol{n_i}$, LP of size $\left(\prod_i \boldsymbol{n_i}, \sum_i \boldsymbol{n_i}\right)$

# Finite Case, LP Formulation

- When $\Omega$ is a finite set, metric $M$, another LP.

$$\min_{\boldsymbol{\mu}} \sum_i \lambda_i W_p^p(\boldsymbol{\mu}, \boldsymbol{\nu_i})$$

# Finite Case, LP Formulation

- When $\Omega$ is a <span style="color:brown">finite set</span>, metric $M$, another LP.

$$\min_{P_1, \cdots, P_N, a} \sum_{i=1}^{N} \lambda_i \langle P_i, M \rangle$$

$$\text{s.t. } P_i{}^T \mathbf{1}_n = b_i, \forall i \leq N,$$

$$P_1 \mathbf{1}_n = \cdots = P_N \mathbf{1}_d = a.$$

If $|\Omega| = n$, LP of size $(Nn^2, (2N-1)n)$; unstable

# Primal Descent on Regularized $W$

$$\min_{\boldsymbol{\mu} \in \boldsymbol{Q} \subset \mathcal{P}(\Omega)} \sum_{i=1}^{N} \lambda_i W_{\boldsymbol{\gamma}}(\boldsymbol{\mu}, \boldsymbol{\nu_i})$$



*Fast Computation of Wasserstein Barycenters*
**International Conference on Machine Learning 2014**

**[CD'14]**

# Primal Descent on Regularized $W$

$$\min_{\boldsymbol{\mu} \in \boldsymbol{Q} \subset \mathcal{P}(\Omega)} \sum_{i=1}^{N} \lambda_i W_{\boldsymbol{\gamma}}(\boldsymbol{\mu}, \boldsymbol{\nu_i})$$



*Fast Computation of Wasserstein Barycenters*
**International Conference on Machine Learning 2014**

**[CD'14]**

# Primal Descent on Regularized $W$

$$\min_{\boldsymbol{\mu} \in \boldsymbol{Q} \subset \mathcal{P}(\Omega)} \sum_{i=1}^{N} \lambda_i W_{\boldsymbol{\gamma}}(\boldsymbol{\mu}, \boldsymbol{\nu_i})$$



*Fast Computation of Wasserstein Barycenters*
**International Conference on Machine Learning 2014**

**[CD'14]**

# *Wasserstein Barycenter* = KL Projections

$$\langle P, M_{\textcolor{red}{X}\textcolor{blue}{Y}} \rangle - \gamma E(P) = \gamma \mathbf{KL}(P|\,\textcolor{red}{K})$$

$$\min_{\boldsymbol{a}} \sum_{i=1}^{N} \lambda_i W_\gamma(\textcolor{red}{\boldsymbol{a}}, \textcolor{blue}{\boldsymbol{b_i}}) = \min_{\substack{\mathbf{P}=[\textcolor{brown}{\boldsymbol{P_1},...,\boldsymbol{P_N}}] \\ \mathbf{P}\in \textcolor{red}{\boldsymbol{C_1}}\cap \textcolor{blue}{\boldsymbol{C_2}}}} \sum_{i=1}^{N} \lambda_i \mathbf{KL}(\textcolor{brown}{\boldsymbol{P_i}}|\textcolor{red}{K})$$

$$\textcolor{red}{\boldsymbol{C_1}} = \{\mathbf{P}|\exists \textcolor{red}{\boldsymbol{a}}, \forall i, P_i \mathbf{1}_m = \textcolor{red}{\boldsymbol{a}}\}$$

$$\textcolor{blue}{\boldsymbol{C_2}} = \{\mathbf{P}|\forall i, P_i^T \mathbf{1}_n = \textcolor{blue}{\boldsymbol{b_i}}\}$$

# *Wasserstein Barycenter* = KL Projections

$$\min_{\boldsymbol{a}} \sum_{i=1}^{N} \lambda_i W_\gamma(\textcolor{red}{\boldsymbol{a}}, \textcolor{blue}{\boldsymbol{b_i}}) = \min_{\substack{\mathbf{P}=[P_1,...,P_N] \\ \mathbf{P} \in \textcolor{red}{C_1} \cap \textcolor{blue}{C_2}}} \sum_{i=1}^{N} \lambda_i \mathbf{KL}(\textcolor{brown}{P_i}|\textcolor{red}{K})$$

$$\textcolor{red}{C_1} = \{\mathbf{P} | \exists \textcolor{red}{\boldsymbol{a}}, \forall i, P_i \mathbf{1}_m = \textcolor{red}{\boldsymbol{a}}\}$$

$$\textcolor{blue}{C_2} = \{\mathbf{P} | \forall i, P_i^T \mathbf{1}_n = \textcolor{blue}{\boldsymbol{b_i}}\}$$

**[BCCNP'15]**

$$[\ \textcolor{red}{K} \cdots \textcolor{blue}{K}\ ] \qquad \mathbf{P}_\gamma$$

# *Wasserstein Barycenter* = KL Projections

$$\min_{\boldsymbol{a}} \sum_{i=1}^{N} \lambda_i W_\gamma(\boldsymbol{a}, \boldsymbol{b_i}) = \min_{\substack{\mathbf{P}=[P_1,...,P_N] \\ \mathbf{P}\in C_1\cap C_2}} \sum_{i=1}^{N} \lambda_i \mathbf{KL}(P_i | K)$$

$$C_1 = \left\{ \mathbf{P} | \exists \boldsymbol{a}, \forall i, P_i \mathbf{1}_m = \boldsymbol{a} \right\}$$

$$C_2 = \left\{ \mathbf{P} | \forall i, P_i^T \mathbf{1}_n = \boldsymbol{b_i} \right\}$$

```
u=ones(size(B)); % d x N matrix
while not converged
    v=u.*(K'*(B./(K*u))); % 2(Nd^2) cost
    u=bsxfun(@times,u,exp(log(v)*weights))./v;
end
a=mean(v,2);
```

**[BCCNP'15]**

*Iterative Bregman Projections for Regularized Transportation Problems*
**SIAM J. on Sci. Comp. 2015**

46

# Application: Graphics

**[S..C..'15]**

47

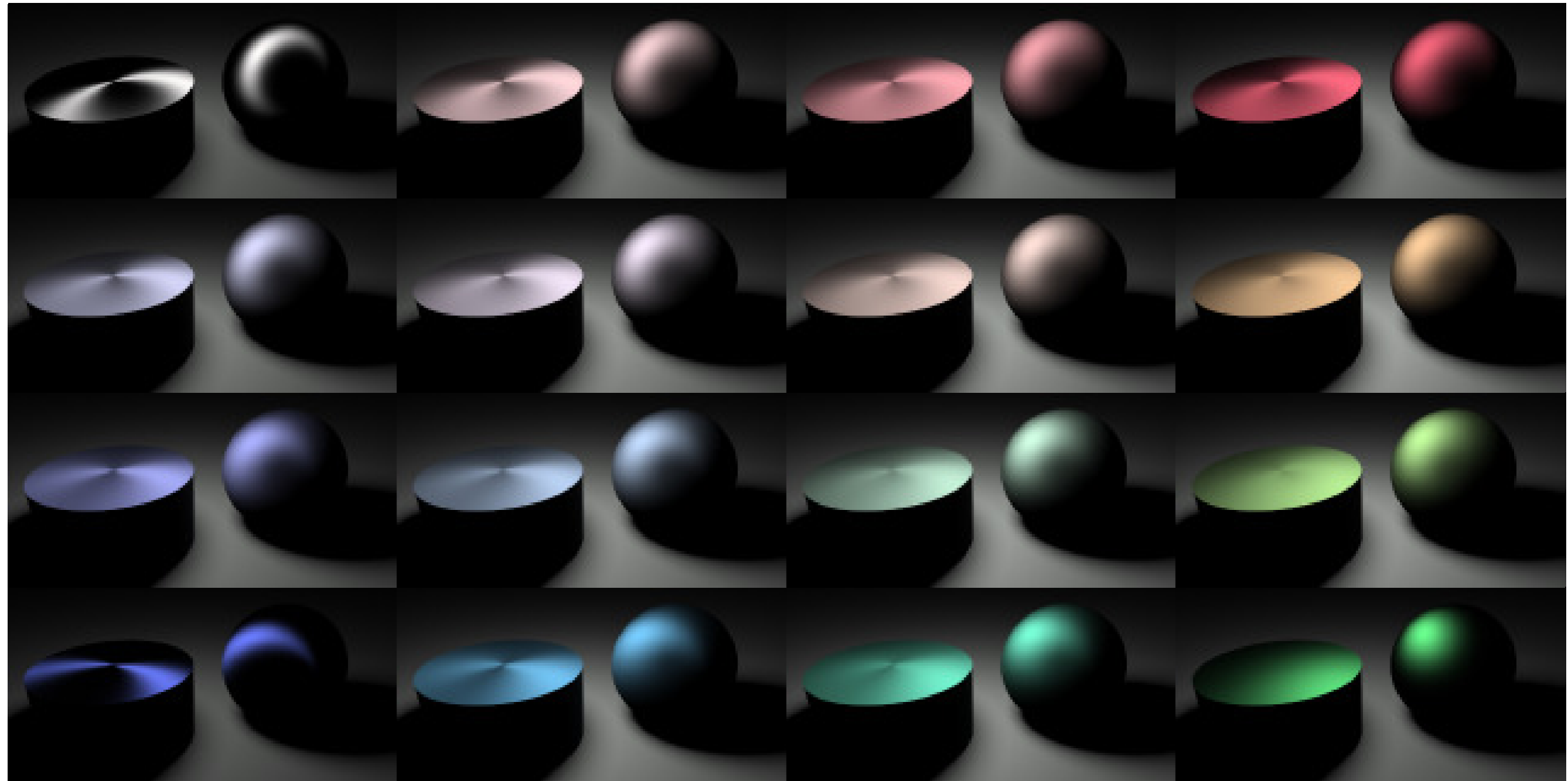# Application: Graphics

[S..C..'15]

# Application: Graphics



*Convolutional Wasserstein Distances: Efficient Optimal Transportation on Geometric Domains,*
**SIGGRAPH'15**

[S..C..'15]

# Application: Graphics
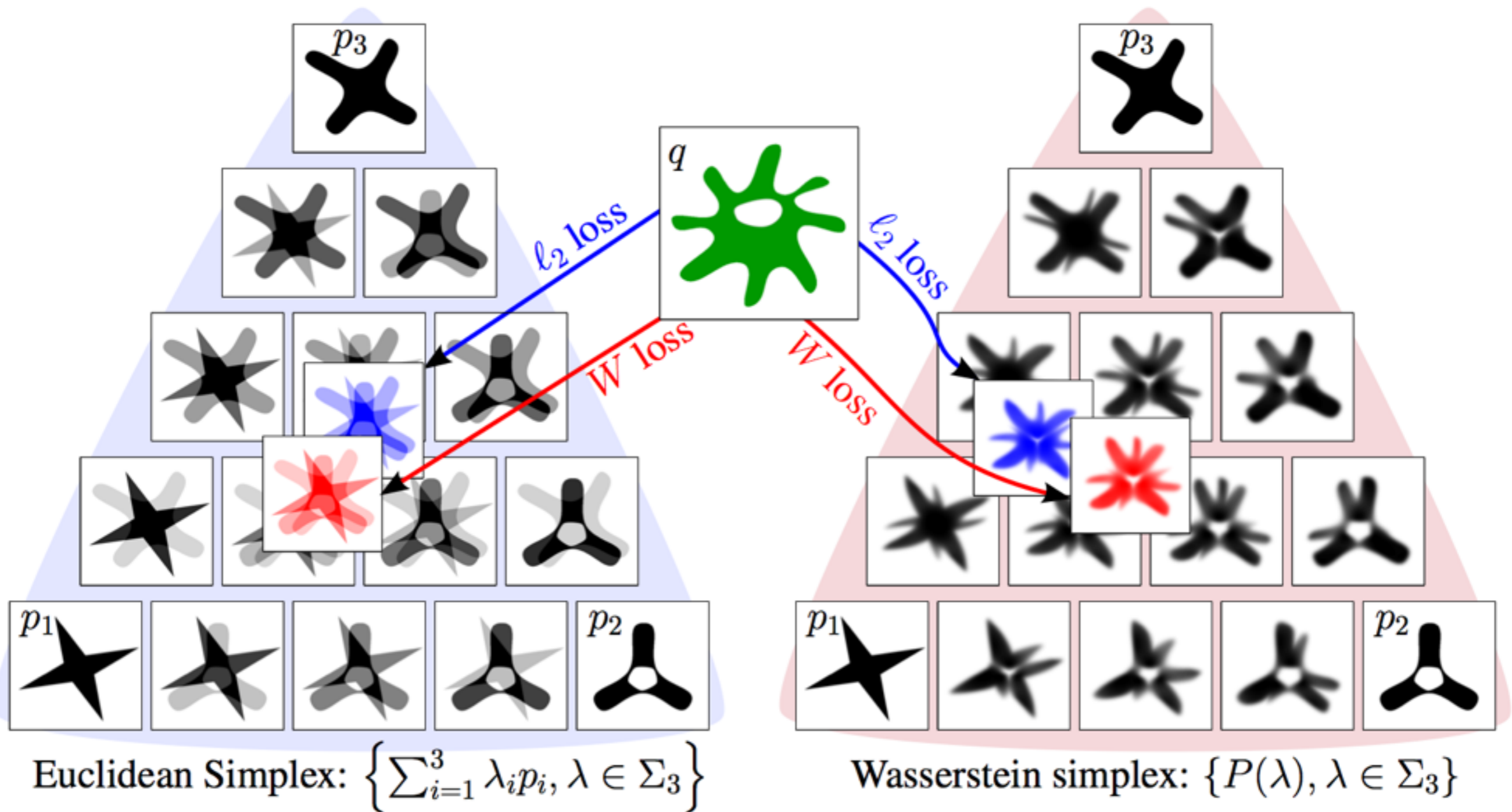


*Convolutional Wasserstein Distances: Efficient*
*Optimal Transportation on Geometric Domains,*
**SIGGRAPH'15**     **[S..C..'15]**

# Application: Graphics



*Convolutional Wasserstein Distances: Efficient*
*Optimal Transportation on Geometric Domains,*
**SIGGRAPH'15** **[S..C..'15]**

47

# *Inverse Wasserstein Problems*

- consider Barycenter operator:

$$b(\lambda) \stackrel{\text{def}}{=} \operatorname*{argmin}_{a} \sum_{i=1}^{N} \lambda_i W_\gamma(a, b_i)$$

- address now **Wasserstein inverse problems**:

$$\text{Given } a, \text{ find } \operatorname*{argmin}_{\lambda \in \Sigma_N} \mathcal{E}(\lambda) \stackrel{\text{def}}{=} \text{Loss}(a, b(\lambda))$$

# The Wasserstein Simplex



Euclidean Simplex: $\left\{\sum_{i=1}^{3} \lambda_i p_i, \lambda \in \Sigma_3\right\}$

Wasserstein simplex: $\{P(\lambda), \lambda \in \Sigma_3\}$

**Prop.** [**BCCNP'15**] Consider $\boldsymbol{B} \in \Sigma_d^N$ and let $\boldsymbol{U_0} = \mathbf{1}_{\boldsymbol{d \times N}}$, and then for $l \geq 0$:

$$\boldsymbol{b}^l \overset{\text{def}}{=} \exp\left(\log\left(K^T \boldsymbol{U_l}\right)\lambda\right) ; \begin{cases} \boldsymbol{V_{l+1}} \overset{\text{def}}{=} \dfrac{\boldsymbol{b}^l \mathbf{1}_N^T}{K^T \boldsymbol{U_l}}, \\ \boldsymbol{U_{l+1}} \overset{\text{def}}{=} \dfrac{\boldsymbol{B}}{K \boldsymbol{V_{l+1}}}. \end{cases}$$

# Using Truncated Barycenters

- instead of using the exact barycenter

$$\underset{\lambda \in \Sigma_N}{\operatorname{argmin}} \, \mathcal{E}(\lambda) \overset{\text{def}}{=} \operatorname{Loss}(\textcolor{green}{a}, \textcolor{blue}{b}(\lambda))$$

- use instead the L-iterate barycenter

$$\underset{\lambda \in \Sigma_N}{\operatorname{argmin}} \, \mathcal{E}^{(L)}(\lambda) \overset{\text{def}}{=} \operatorname{Loss}(\textcolor{green}{a}, \textcolor{blue}{b}^{(L)}(\lambda))$$

- Differente using **the chain rule.**

$$\nabla \mathcal{E}^{(L)}(\lambda) = [\partial \textcolor{blue}{b}^{(L)}]^T(\textcolor{red}{g}), \; \textcolor{red}{g} \overset{\text{def}}{=} \nabla \operatorname{Loss}(\textcolor{green}{a}, \cdot)|_{\textcolor{blue}{b}^{(L)}(\lambda)} \cdot$$

# Gradient / Barycenter Computation

**function** $\text{SINKHORN-DIFFERENTIATE}((p_s)_{s=1}^{S}, q, \lambda)$

$\quad \forall\, s, b_s^{(0)} \leftarrow \mathbb{1}$

$\quad (w, r) \leftarrow (0^S, 0^{S \times N})$

$\quad$**for** $\ell = 1, 2, \ldots, L \qquad$ // *Sinkhorn loop*

$\qquad \forall\, s, \varphi_s^{(\ell)} \leftarrow K^{\top} \dfrac{p_s}{K b_s^{(\ell-1)}}$

$\qquad p \leftarrow \prod_s \left( \varphi_s^{(\ell)} \right)^{\lambda_s}$

$\qquad \forall\, s, b_s^{(\ell)} \leftarrow \dfrac{p}{\varphi_s^{(\ell)}}$

$\quad g \leftarrow \nabla \mathcal{L}(p, q) \odot p$

$\quad$**for** $\ell = L, L-1, \ldots, 1 \qquad$ // *Reverse loop*

$\qquad \forall\, s, w_s \leftarrow w_s + \langle \log \varphi_s^{(\ell)}, g \rangle$

$\qquad \forall\, s, r_s \leftarrow -K^{\top} \left( K \left( \dfrac{\lambda_s g - r_s}{\varphi_s^{(\ell)}} \right) \odot \dfrac{p_s}{(K b_s^{(\ell-1)})^2} \right) \odot b_s^{(\ell-1)}$

$\qquad g \leftarrow \sum_s r_s$

$\quad$**return** $P^{(L)}(\lambda) \leftarrow p, \nabla \mathcal{E}_L(\lambda) \leftarrow w$

# Application: Volume Reconstruction



Shape database $(p_1, \ldots, p_5)$

Input shape $q$

Projection $P(\lambda)$

Iso-surface

*Wasserstein Barycentric Coordinates: Histogram Regression using Optimal Transport*, **SIGGRAPH'16**

[BPC'16]

53

# Application: Color Grading



$\lambda_0 = 0.03$

$\lambda_1 = 0.12$

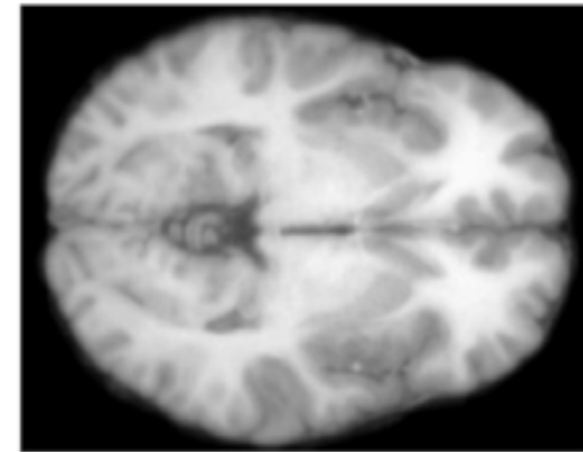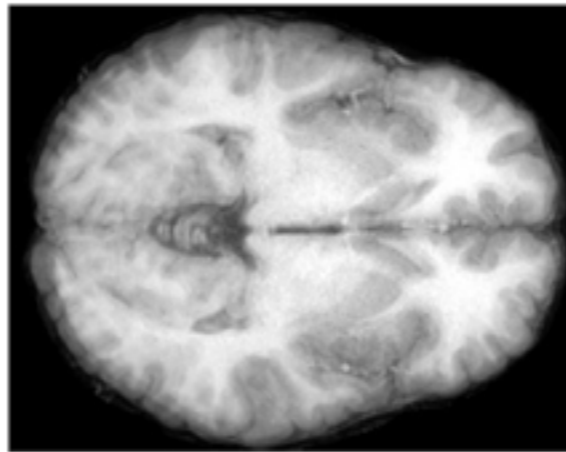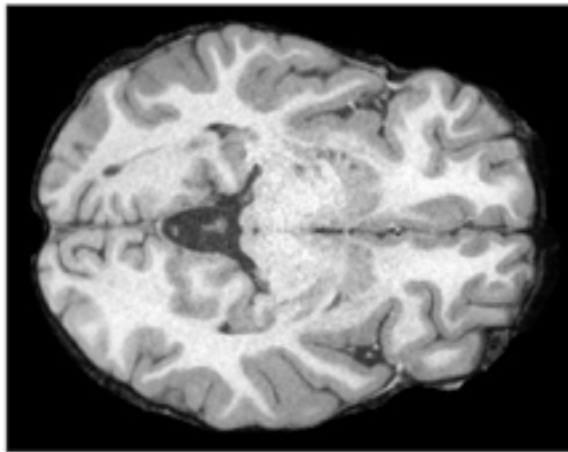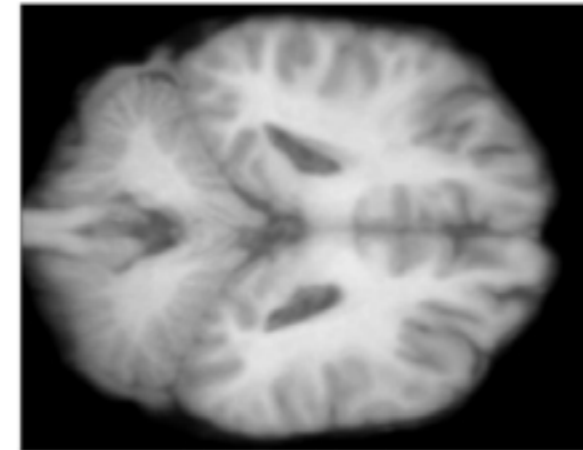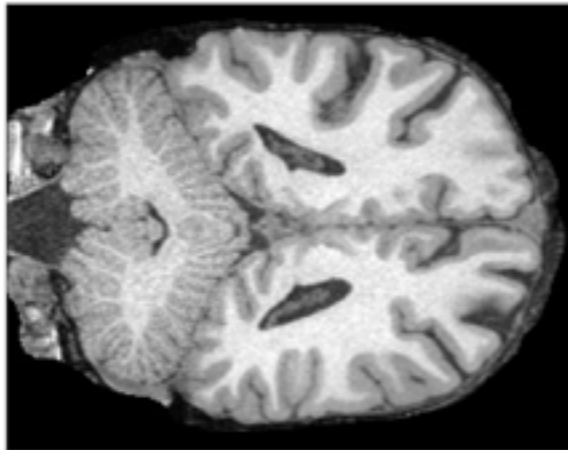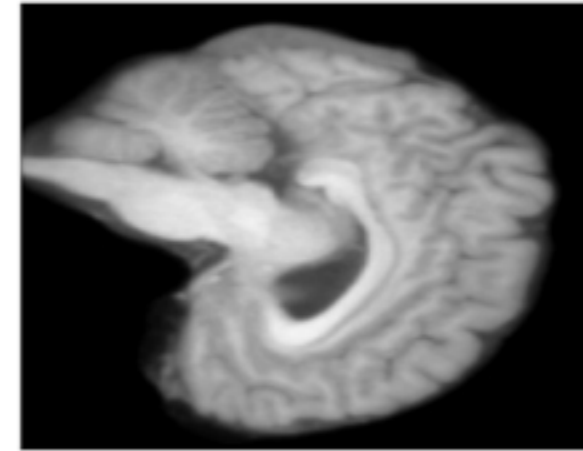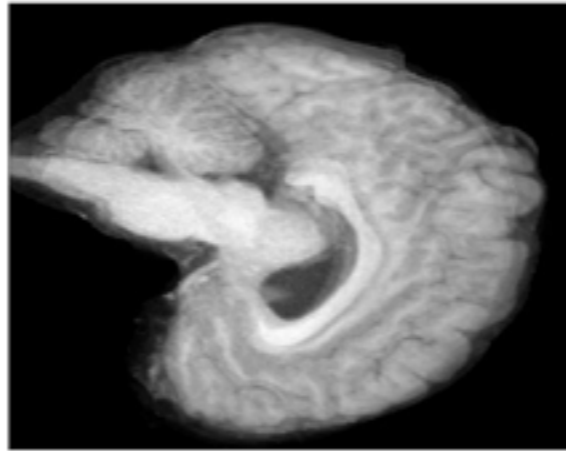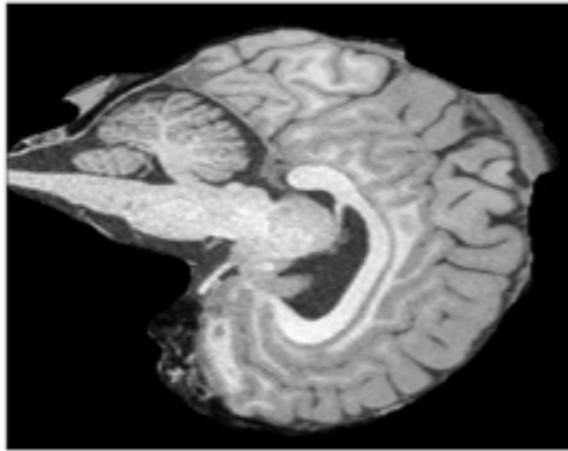$\lambda_2 = 0.40$

$\lambda_3 = 0.43$

# Application: Color Grading



*Wasserstein Barycentric Coordinates: Histogram Regression using Optimal Transport*, **SIGGRAPH'16**

**[BPC'16]**

# Application: Brain Mapping



Original | Euclidean projection | Wasserstein projection

# To conclude

- *Entropy* regularization is a very effective way to get OT to work as a generic loss.

- Many recent extensions:

  - **[Schmitzer'16]**: fast multiscale approaches

  - **[ZFMAP'15] [CSPV'16]**: Unbalanced transport

  - **[SPKS'16] [PCS'16]** extensions to *Gromov-W.*

  - **[FCTR'15]** Domain adaptation in ML