# Fitting Generative Models with Optimal Transport

## Marco Cuturi

ENSAE ParisTech

École nationale de la statistique et de l'administration économique

université PARIS-SACLAY

*Joint work / work in progress with*
G. Peyré, A. Genevay *(ENS)*, F. Bach *(INRIA)*,
G. Montavon, K-R Müller *(TU Berlin)*

# Maximum Likelihood Estimation



$$\mathcal{P}(\Omega)$$

$$\boldsymbol{\nu}_{\text{data}} = \frac{1}{N} \sum_{i=1}^{N} \delta_{\boldsymbol{x_i}}$$

$$\{p_\theta, \theta \in \Theta\}$$

$$p_{\boldsymbol{\theta}^\star}$$

MLE

$$\min_{\boldsymbol{\theta} \in \Theta} \text{KL}(\boldsymbol{\nu}_{\text{data}} \| \boldsymbol{p_\theta})$$

$$\min_{\boldsymbol{\theta} \in \Theta} -\frac{1}{N} \sum_{i=1}^{N} \log \boldsymbol{p_\theta}(\boldsymbol{x_i})$$
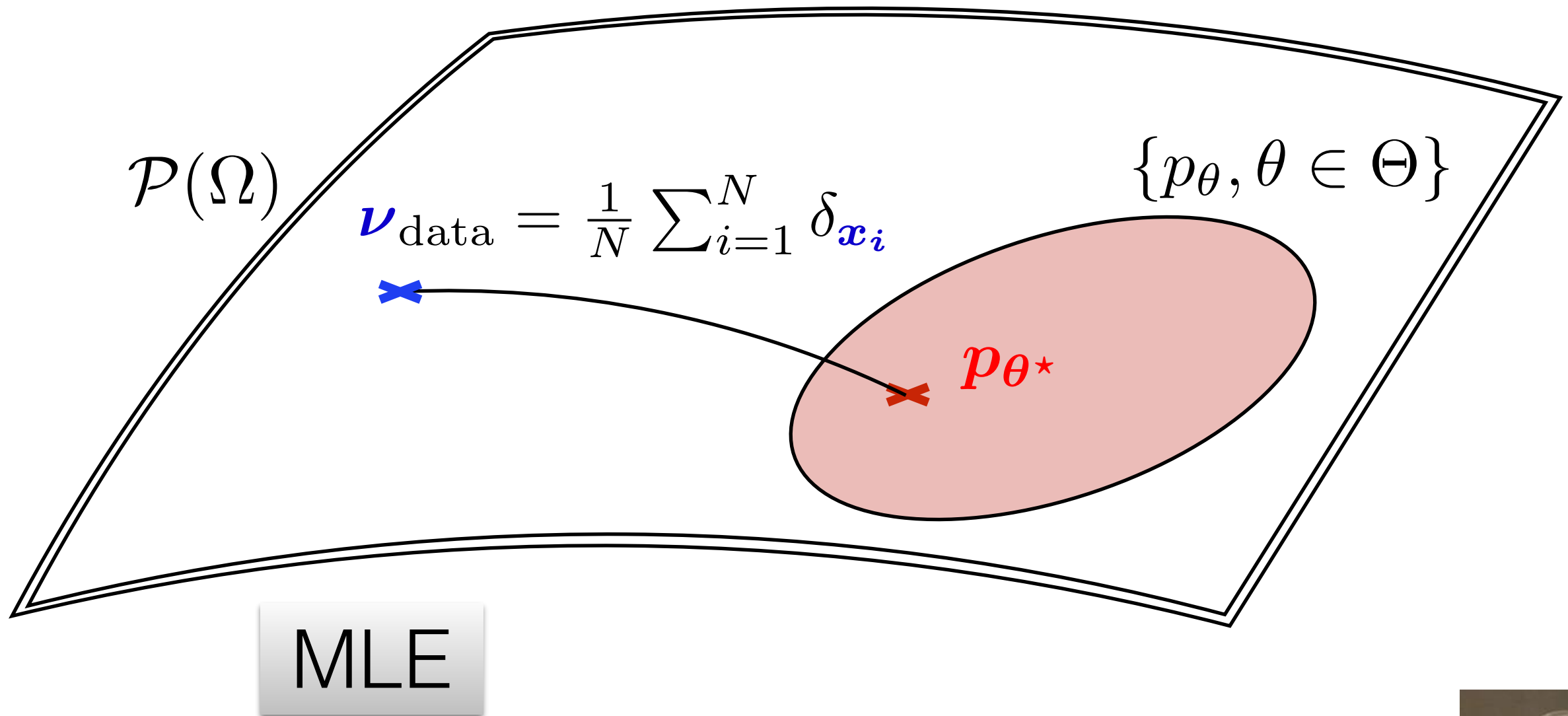
ON AN ABSOLUTE CRITERION
FOR FITTING FREQUENCY CURVES.

By *R. A. Fisher*, Gonville and Caius College, Cambridge.

1. IF we set ourselves the problem, in its essence one of frequent occurrence, of finding the arbitrary elements in a function of known form, which best suit a set of actual observations, we are met at the outset by an arbitrariness which appears to invalidate any results we may obtain. In

2

# Maximum Likelihood Estimation



$$\min_{\boldsymbol{\theta} \in \Theta} \mathrm{KL}(\boldsymbol{\nu}_{\mathrm{data}} \| \boldsymbol{p_\theta})$$

$$\min_{\boldsymbol{\theta} \in \Theta} -\frac{1}{N} \sum_{i=1}^{N} \log \boldsymbol{p_\theta}(\boldsymbol{x_i})$$

# Minimum * Estimation

## MINIMUM CHI-SQUARE, NOT MAXIMUM LIKELIHOOD!

BY JOSEPH BERKSON

Mayo Clinic, Rochester, Minnesota

ELSEVIER

COMPUTATIONAL
STATISTICS
& DATA ANALYSIS

Springer Series in Statistics

**Luc Devroye
Gábor Lugosi**

Combinatorial
Methods in
Density
Estimation

$l_1$

## Minimum Hellinger distance estimation for Poisson mixtures

Dimitris Karlis, Evdokia Xekalaki[*]

Department of Statistics, Athens University of Economics and Business, 76 Patission Str., 104 34 Athens, Greece

Available online at www.sciencedirect.com

SCIENCE ⓓ DIRECT®

ELSEVIER

STATISTICS &
PROBABILITY
LETTERS

www.elsevier.com/locate/stapro

## On minimum Kantorovich distance estimators

Federico Bassetti[a], Antonella Bodini[b], Eugenio Regazzini[a],[*]

# Statistical Estimation



$\mathcal{P}(\Omega)$

$\boldsymbol{\nu}_{\mathrm{data}} = \frac{1}{N} \sum_{i=1}^{N} \delta_{\boldsymbol{x_i}}$

$\{p_\theta, \theta \in \Theta\}$

$p_{\boldsymbol{\theta^\star}}$

# Statistical Estimation



$\mathcal{P}(\Omega)$

$\{p_\theta, \theta \in \Theta\}$

$\nu_{\text{data}} = \frac{1}{N} \sum_{i=1}^{N} \delta_{\boldsymbol{x_i}}$

$p_{\boldsymbol{\theta^\star}}$

$$\min_{\boldsymbol{\theta} \in \Theta} \mathrm{KL}(\boldsymbol{\nu}_{\text{data}} \| \boldsymbol{p_\theta} \|)$$ MLE

$$\min_{\boldsymbol{\theta} \in \Theta} W(\boldsymbol{\nu}_{\text{data}}, \boldsymbol{p_\theta})$$ MKE

**[Bassetti'06]**

# Model = positive densities



$$\boldsymbol{\nu}_{\mathrm{data}}$$

$$\Omega$$

$$\min_{\boldsymbol{\theta} \in \Theta} - \sum_{i=1}^{N} \log \boldsymbol{p_\theta}(\boldsymbol{x_i})$$

# Model = positive densities



$$\min_{\boldsymbol{\theta} \in \Theta} - \sum_{i=1}^{N} \log \boldsymbol{p_\theta}(\boldsymbol{x_i})$$

# Model = positive densities



$p_{\boldsymbol{\theta}}$

$\boldsymbol{\nu}_{\text{data}}$

$\Omega$

$$\min_{\boldsymbol{\theta} \in \Theta} - \sum_{i=1}^{N} \log \boldsymbol{p_{\theta}(x_i)}$$

$$\min_{\boldsymbol{\theta} \in \Theta} W(\boldsymbol{\nu}_{\text{data}}, \boldsymbol{p_{\theta}})$$

# Model = generative

# Model = generative



$\boldsymbol{\mu}$

latent
space

$\boldsymbol{\nu}_{\text{data}}$

data
space

$\Omega$

6

# Model = generative

$\boldsymbol{\mu}$

latent
space

$f_{\boldsymbol{\theta}}$ : latent space $\rightarrow$ data space

data
space

$\Omega$

$\boldsymbol{\nu}_{\text{data}}$

# Model = generative



$\boldsymbol{\mu}$

latent
space

$f_{\boldsymbol{\theta}} : \text{latent space} \to \text{data space}$

$f_{\boldsymbol{\theta}\sharp}\boldsymbol{\mu}$

$\boldsymbol{\nu}_{\text{data}}$

data
space

$\Omega$

# Model = generative

# Illustration - GAN



$\mu$

latent

[RMC'16]

100 z — Code — Project and reshape — 1024 — 4 / 4 — Deconv 1 — 512 — 8 / 8 — 5 / 5 — Stride 2 — 256 — 16 / 16 — Stride 2 — 5 / 5 — Deconv 2 — 128 — 32 / 32 — 5 / 5 — Stride 2 — Deconv 3 — 64 — 5 / 5 — Stride 2 — 3 — 64 — Image — Deconv 4

http://torch.ch/blog/2015/11/13/gan.html

# Illustration - GAN

$\mu$

latent



**[RMC'16]**

Code  Project and reshape

100 z

1024

4
4

Deconv 1

512

8

8

5
5

Stride 2

16

Stride 2

5
5

16

Deconv 2

32

5
5

32

Stride 2

Deconv 3

64

5

5

Stride 2

3

5
5

64

Deconv 4

Image

# Wasserstein Distances

**Def.** For $p \geq 1$, the $p$-Wasserstein distance between $\color{red}\boldsymbol{\mu}$, $\color{blue}\boldsymbol{\nu}$ in $\mathcal{P}(\Omega)$, defined by a metric $\color{green}\boldsymbol{D}$ on $\Omega$,

$$W_p^p(\color{red}\boldsymbol{\mu}\color{black}, \color{blue}\boldsymbol{\nu}\color{black}) \stackrel{\text{def}}{=} \inf_{\color{brown}\boldsymbol{P}\color{black}\in\Pi(\color{red}\boldsymbol{\mu}\color{black},\color{blue}\boldsymbol{\nu}\color{black})} \iint \color{green}\boldsymbol{D}(\boldsymbol{X}, \boldsymbol{Y})\color{black}^p d\color{brown}\boldsymbol{P}\color{black}(X, Y).$$

PRIMAL

$$W_p^p(\color{red}\boldsymbol{\mu}\color{black}, \color{blue}\boldsymbol{\nu}\color{black}) = \sup_{\substack{\color{red}\boldsymbol{\varphi}\color{black}\in L_1(\color{red}\boldsymbol{\mu}\color{black}),\color{blue}\boldsymbol{\psi}\color{black}\in L_1(\color{blue}\boldsymbol{\nu}\color{black}) \\ \color{red}\boldsymbol{\varphi}\color{black}(x)+\color{blue}\boldsymbol{\psi}\color{black}(y)\leq \color{green}\boldsymbol{D}\color{black}^p(x,y)}} \int \color{red}\boldsymbol{\varphi}\color{black} d\color{red}\boldsymbol{\mu}\color{black} + \int \color{blue}\boldsymbol{\psi}\color{black} d\color{blue}\boldsymbol{\nu}\color{black}.$$
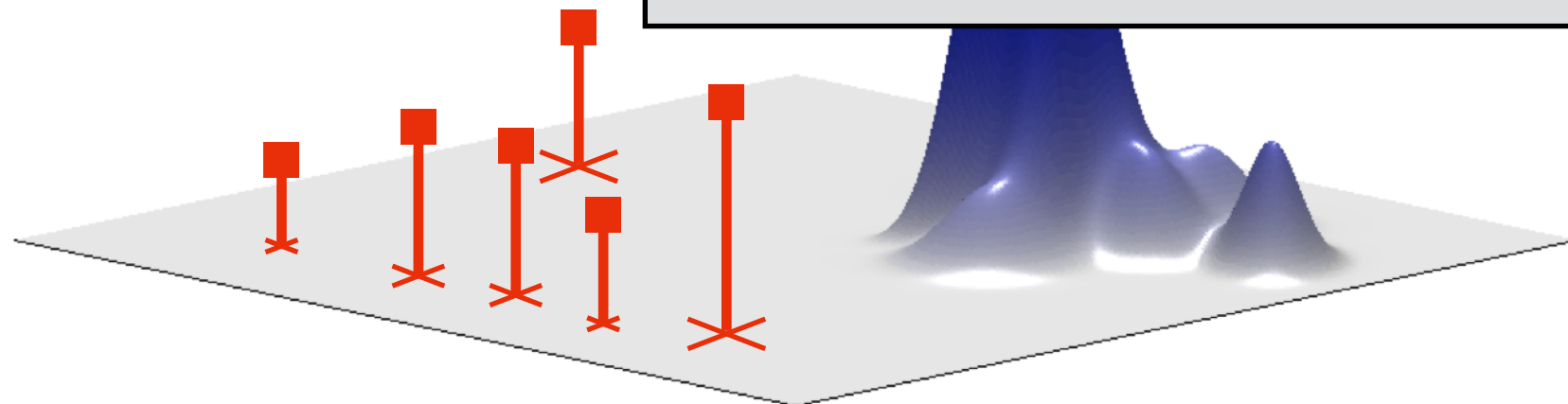
DUAL

8

# $W$ is versatile
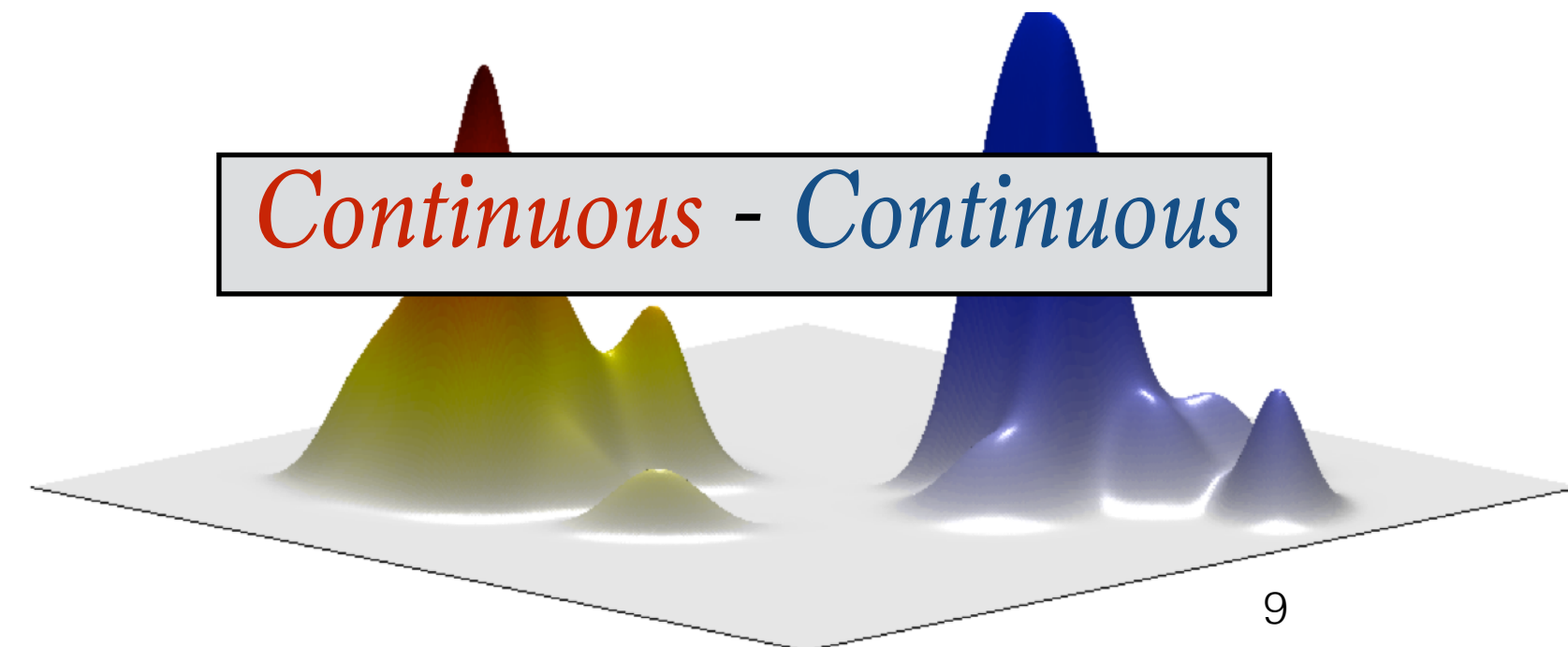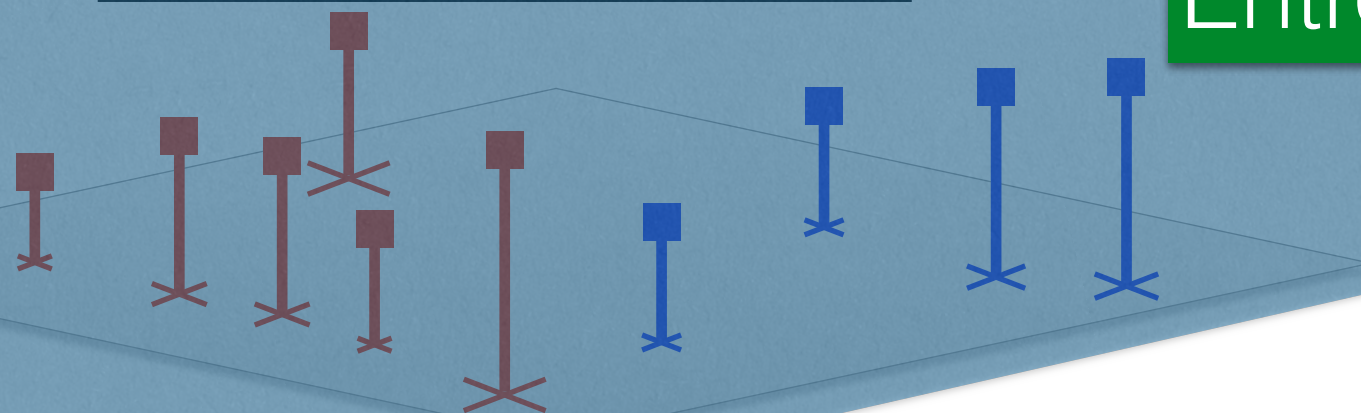


Discrete - Discrete

Discrete - Continuous
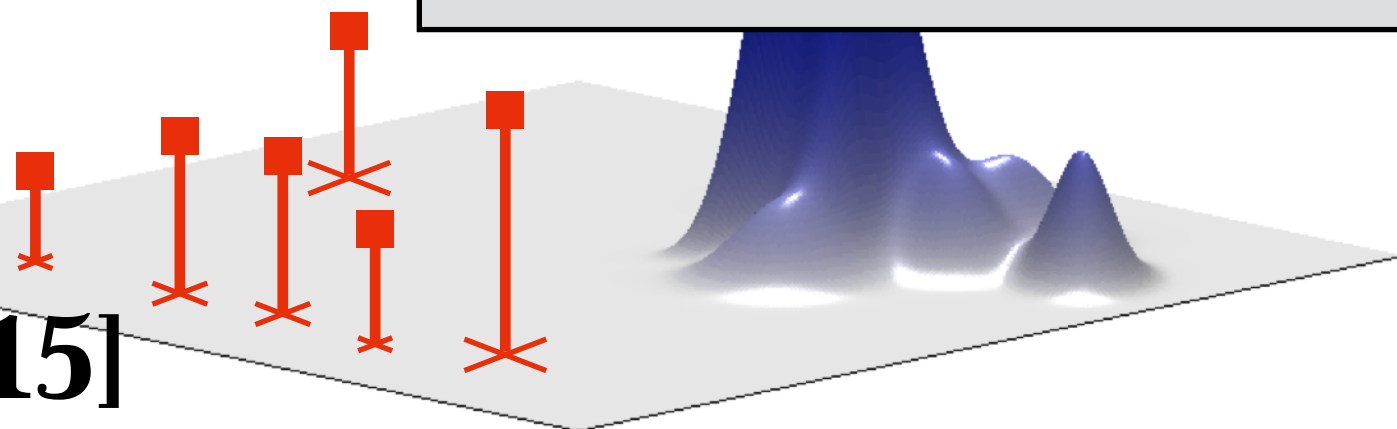
Continuous - Continuous

9

# $W$ is versatile

Discrete - Discrete

Network flow solver
Entropic regularization

Discrete - Continuous

low dim.

**[M'11][KMB'16] [L'15]**

Continuous - Continuous

Stochastic
Optimization

**[GCPB'16]**

# Dual regularization

$$W_p^p(\boldsymbol{\mu}, \boldsymbol{\nu}) = \sup_{\boldsymbol{\varphi}, \boldsymbol{\psi}} \int \boldsymbol{\varphi} d\boldsymbol{\mu} + \int \boldsymbol{\psi} d\boldsymbol{\nu} - \iota_C(\boldsymbol{\varphi}, \boldsymbol{\psi})$$

$$C = \{(\boldsymbol{\varphi}, \boldsymbol{\psi}) | \boldsymbol{\varphi} \oplus \boldsymbol{\psi} \leq \boldsymbol{D}^p\}$$

DUAL



$\gamma e^{-x/\gamma}$

Legend: $\gamma = .1$ (red), $\gamma = .01$ (blue), $\gamma = .001$ (black)

# Dual regularization

$$W_p^p(\boldsymbol{\mu}, \boldsymbol{\nu}) = \sup_{\boldsymbol{\varphi}, \boldsymbol{\psi}} \int \boldsymbol{\varphi} d\boldsymbol{\mu} + \int \boldsymbol{\psi} d\boldsymbol{\nu} - \iota_C(\boldsymbol{\varphi}, \boldsymbol{\psi})$$

$$C = \{(\boldsymbol{\varphi}, \boldsymbol{\psi}) | \boldsymbol{\varphi} \oplus \boldsymbol{\psi} \leq \boldsymbol{D}^p\}$$

DUAL



$\gamma e^{-x/\gamma}$

$\iota$

Legend: $\gamma = .1$ (red), $\gamma = .01$ (blue), $\gamma = .001$ (black)
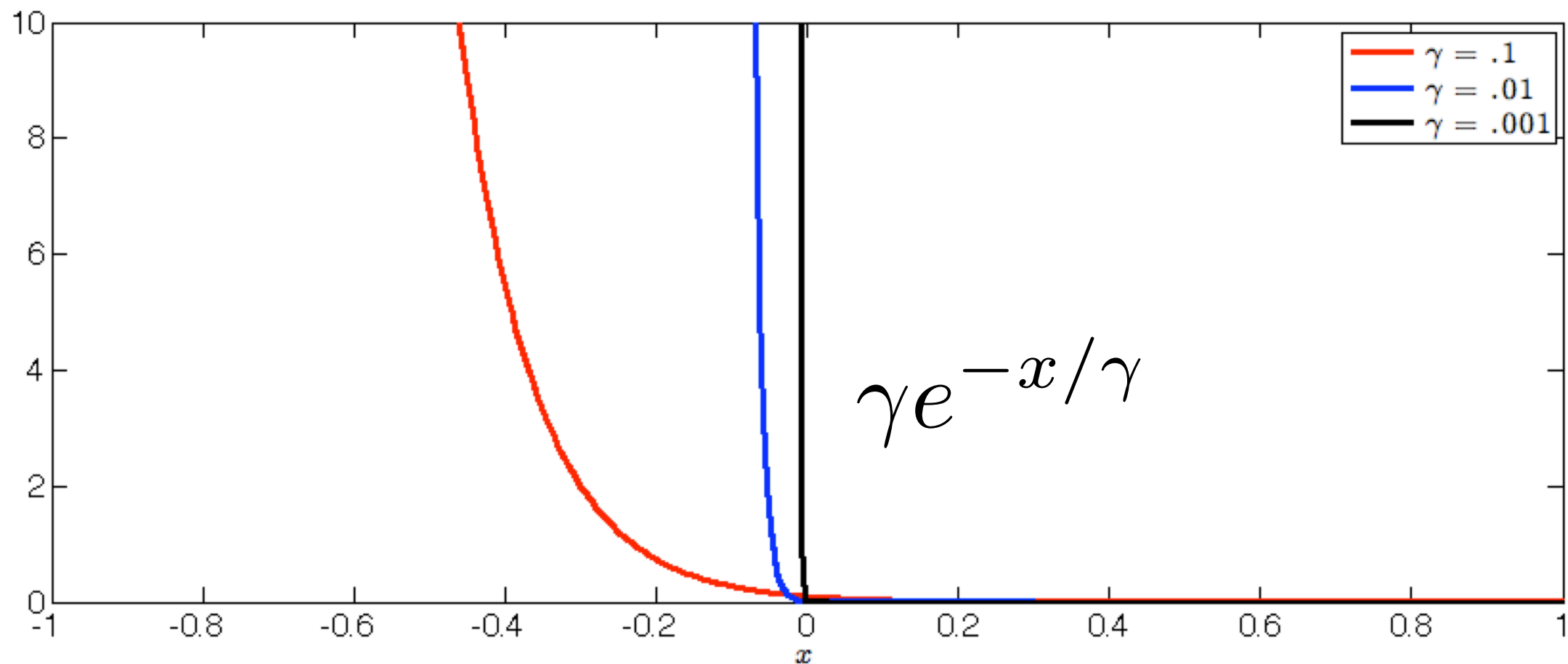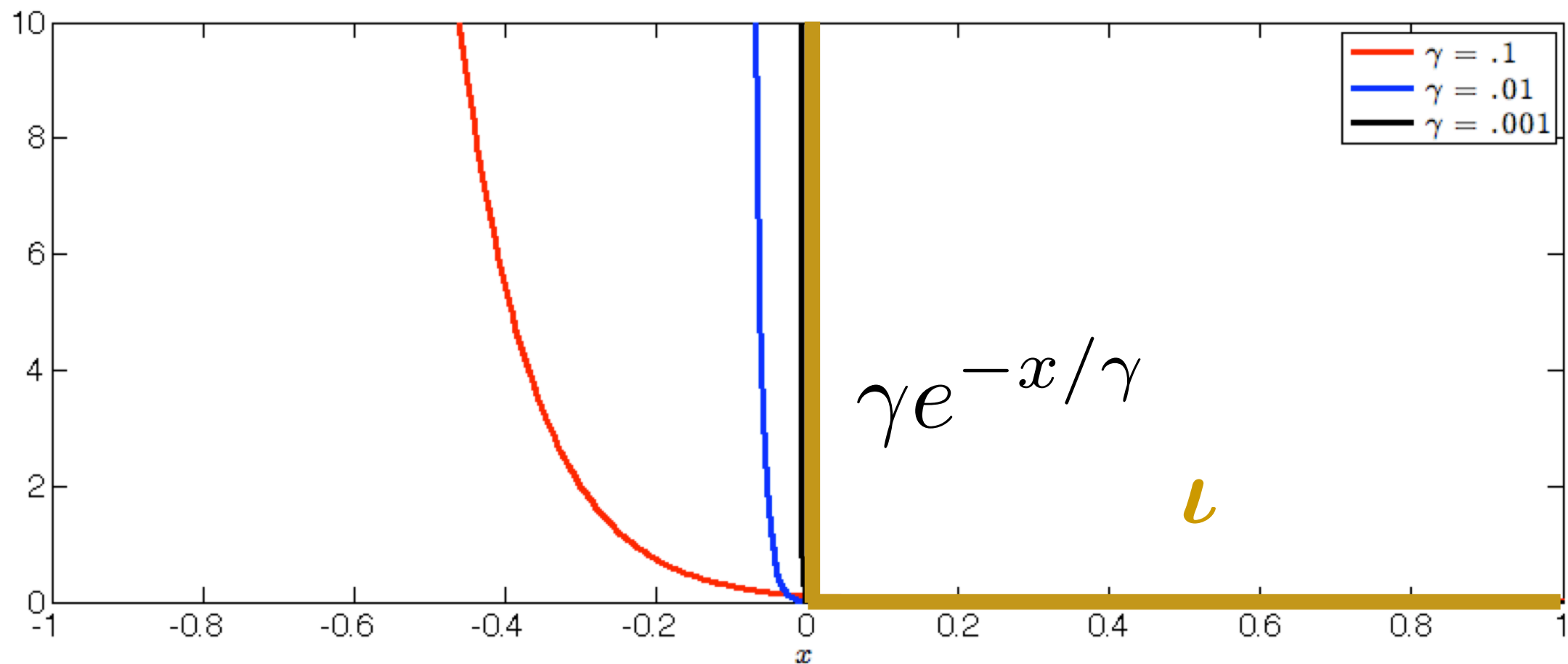
10

# Dual regularization

$$W_p^p(\boldsymbol{\mu}, \boldsymbol{\nu}) = \sup_{\boldsymbol{\varphi}, \boldsymbol{\psi}} \int \boldsymbol{\varphi} d\boldsymbol{\mu} + \int \boldsymbol{\psi} d\boldsymbol{\nu} - \iota_C(\boldsymbol{\varphi}, \boldsymbol{\psi})$$

$$C = \{(\boldsymbol{\varphi}, \boldsymbol{\psi}) | \boldsymbol{\varphi} \oplus \boldsymbol{\psi} \leq \boldsymbol{D}^p\}$$

DUAL

*regularizing dual* *constraints* $\gamma > 0$

$$W_\gamma(\boldsymbol{\mu}, \boldsymbol{\nu}) = \sup_{\boldsymbol{\varphi}, \boldsymbol{\psi}} \int \boldsymbol{\varphi} d\boldsymbol{\mu} + \int \boldsymbol{\psi} d\boldsymbol{\nu} - \iota_C^\gamma(\boldsymbol{\varphi}, \boldsymbol{\psi})$$

$$\iota_C^\gamma(\boldsymbol{\varphi}, \boldsymbol{\psi}) = \gamma \iint e^{(\boldsymbol{\varphi} \oplus \boldsymbol{\psi} - \boldsymbol{D}^p)/\gamma} d\boldsymbol{\mu} d\boldsymbol{\nu}$$

REGULARIZED DUAL

11

# *W* is versatile

*Discrete* - *Discrete*

Network flow solver
Entropic regularization

# OT on Two Empirical Measures



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

$$(\Omega, D)$$

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

13

# OT on Two Empirical Measures

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



$(\Omega, D)$

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

# Dual regularization, *Discrete*

$$W_\gamma(\boldsymbol{\mu}, \boldsymbol{\nu}) = \sup_{\boldsymbol{\varphi}, \boldsymbol{\psi}} \int \boldsymbol{\varphi} d\boldsymbol{\mu} + \int \boldsymbol{\psi} d\boldsymbol{\nu} - \iota_C^\gamma(\boldsymbol{\varphi}, \boldsymbol{\psi})$$

$$\iota_C^\gamma(\boldsymbol{\varphi}, \boldsymbol{\psi}) = \gamma \iint e^{(\boldsymbol{\varphi} \oplus \boldsymbol{\psi} - \boldsymbol{D}^p)/\gamma} d\boldsymbol{\mu} d\boldsymbol{\nu}$$
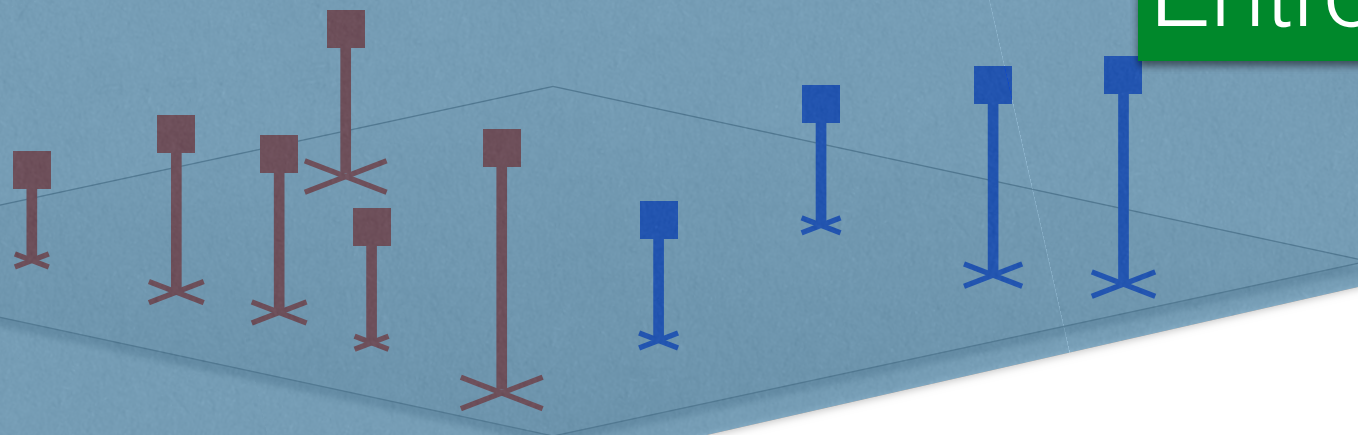
REGULARIZED DUAL

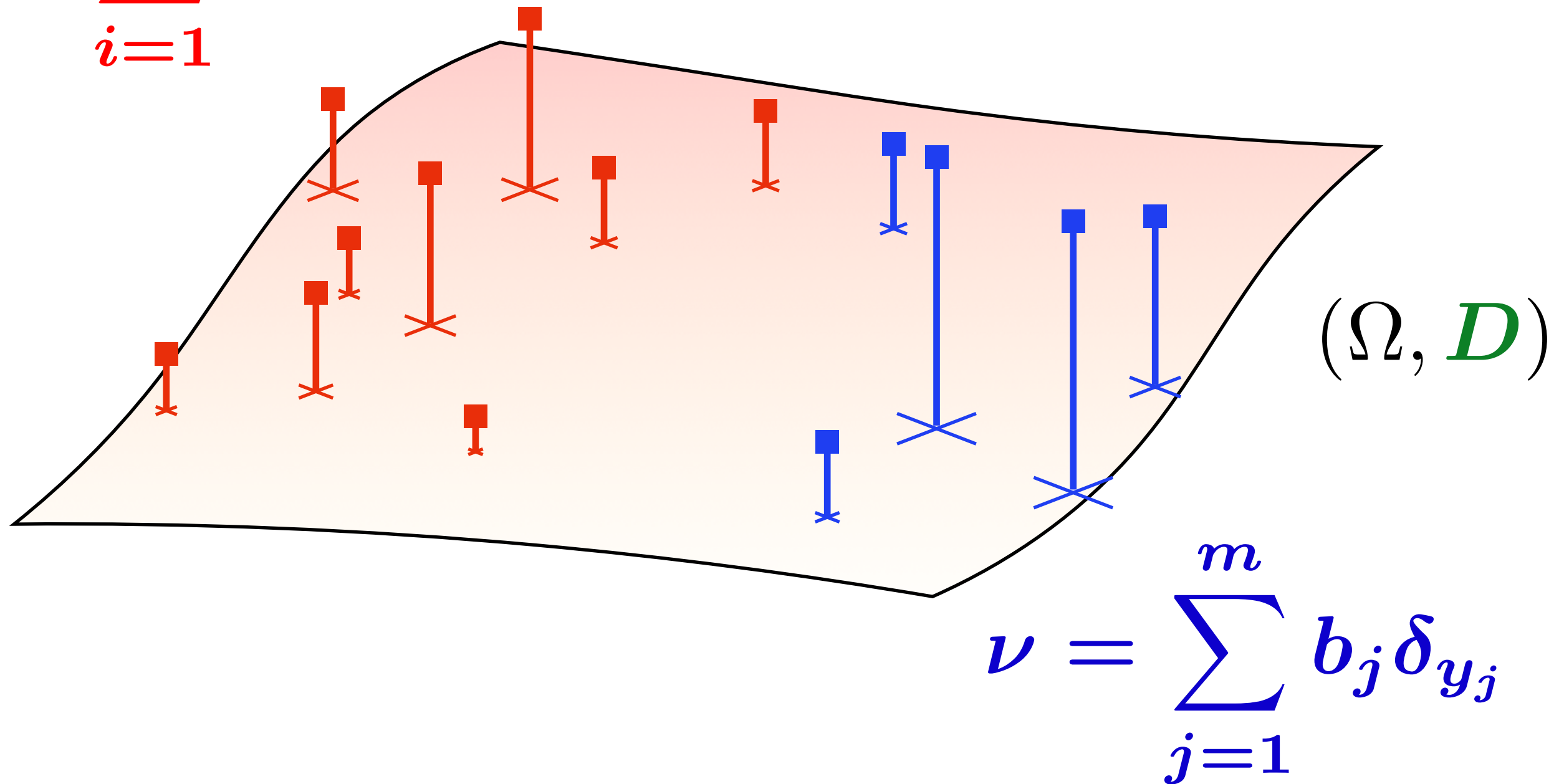$$\boldsymbol{\mu} = \sum_{i=1}^n \boldsymbol{a_i} \delta_{\boldsymbol{x_i}}$$

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

$$W_\gamma(\boldsymbol{\mu}, \boldsymbol{\nu}) = \max_{\boldsymbol{\alpha}, \boldsymbol{\beta}} \boldsymbol{\alpha}^T \boldsymbol{a} + \boldsymbol{\beta}^T \boldsymbol{b} - \gamma \sum_{\boldsymbol{ij}} \boldsymbol{a_i} \boldsymbol{b_j} \, e^{\frac{\boldsymbol{\alpha_i} + \boldsymbol{\beta_j} - \boldsymbol{D}^p(\boldsymbol{x_i}, \boldsymbol{y_j})}{\gamma}}$$

REGULARIZED DISCRETE DUAL

14

# Dual regularization, *Discrete*

$$W_\gamma(\boldsymbol{\mu}, \boldsymbol{\nu}) = \sup_{\boldsymbol{\varphi}, \boldsymbol{\psi}} \int \boldsymbol{\varphi} d\boldsymbol{\mu} + \int \boldsymbol{\psi} d\boldsymbol{\nu} - \iota_C^\gamma(\boldsymbol{\varphi}, \boldsymbol{\psi})$$

$$\iota_C^\gamma(\boldsymbol{\varphi}, \boldsymbol{\psi}) = \gamma \iint e^{(\boldsymbol{\varphi} \oplus \boldsymbol{\psi} - D^p)/\gamma} d\boldsymbol{\mu} d\boldsymbol{\nu}$$

REGULARIZED DUAL

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

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

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

$$\text{where } K = \left[ e^{-\frac{D^p(\boldsymbol{x_i}, \boldsymbol{y_j})}{\gamma}} \right]_{ij}$$

REGULARIZED DISCRETE DUAL

15

# Algorithm: Block Coordinate Ascent

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

REGULARIZED DISCRETE DUAL

$$\mathcal{E}(\boldsymbol{\alpha}, \boldsymbol{\beta}) = \boldsymbol{\alpha}^T \boldsymbol{a} + \boldsymbol{\beta}^T \boldsymbol{b} - \gamma(\boldsymbol{a} \odot e^{\boldsymbol{\alpha}/\gamma})^T \; K(\boldsymbol{b} \odot e^{\boldsymbol{\beta}/\gamma})$$

$$\nabla_{\boldsymbol{\alpha}} \mathcal{E} = \boldsymbol{a} - \boldsymbol{a} \odot e^{\boldsymbol{\alpha}/\gamma} \odot \; K(\boldsymbol{b} \odot e^{\boldsymbol{\beta}/\gamma})$$

$$\nabla_{\boldsymbol{\beta}} \mathcal{E} = \boldsymbol{b} - \boldsymbol{b} \odot e^{\boldsymbol{\beta}/\gamma} \odot \; K^T(\boldsymbol{a} \odot e^{\boldsymbol{\alpha}/\gamma})$$

# Algorithm: Block Coordinate Ascent

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

REGULARIZED DISCRETE DUAL

$$\mathcal{E}(\boldsymbol{\alpha}, \boldsymbol{\beta}) = \boldsymbol{\alpha}^T \boldsymbol{a} + \boldsymbol{\beta}^T \boldsymbol{b} - \gamma(\boldsymbol{a} \odot e^{\boldsymbol{\alpha}/\gamma})^T \ K(\boldsymbol{b} \odot e^{\boldsymbol{\beta}/\gamma})$$

$$\nabla_{\boldsymbol{\alpha}} \mathcal{E} = \boldsymbol{a} - \boldsymbol{a} \odot e^{\boldsymbol{\alpha}/\gamma} \odot \ K(\boldsymbol{b} \odot e^{\boldsymbol{\beta}/\gamma})$$

$$\boldsymbol{\alpha} \leftarrow -\gamma \log \ K(\boldsymbol{b} \odot e^{\boldsymbol{\beta}/\gamma})$$

$$\nabla_{\boldsymbol{\beta}} \mathcal{E} = \boldsymbol{b} - \boldsymbol{b} \odot e^{\boldsymbol{\beta}/\gamma} \odot \ K^T(\boldsymbol{a} \odot e^{\boldsymbol{\alpha}/\gamma})$$

$$\boldsymbol{\beta} \leftarrow -\gamma \log \ K^T(\boldsymbol{a} \odot e^{\boldsymbol{\alpha}/\gamma})$$

16

# *Algorithm: Block Coordinate Ascent*

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

REGULARIZED DISCRETE DUAL

17

# *Algorithm: Block Coordinate Ascent*

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

REGULARIZED DISCRETE DUAL

$$\boldsymbol{u} \leftarrow \frac{\boldsymbol{a}}{K\boldsymbol{v}}$$

$$\boldsymbol{v} \leftarrow \frac{\boldsymbol{b}}{K^T \boldsymbol{u}}$$

17

# Algorithm: Block Coordinate Ascent

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

REGULARIZED DISCRETE DUAL

$$(\boldsymbol{u}, \boldsymbol{v}) \overset{\text{def}}{=} (\boldsymbol{a} \odot e^{\boldsymbol{\alpha}/\gamma}, \boldsymbol{b} \odot e^{\boldsymbol{\beta}/\gamma})$$

$$\boldsymbol{\alpha} \leftarrow -\gamma \log K(\boldsymbol{b} \odot e^{\boldsymbol{\beta}/\gamma})$$

$$\boldsymbol{u} \leftarrow \frac{\boldsymbol{a}}{K\boldsymbol{v}}$$

$$\boldsymbol{\beta} \leftarrow -\gamma \log K^T(\boldsymbol{a} \odot e^{\boldsymbol{\alpha}/\gamma})$$

$$\boldsymbol{v} \leftarrow \frac{\boldsymbol{b}}{K^T \boldsymbol{u}}$$

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

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

$$W_\gamma(\textcolor{red}{\boldsymbol{\mu}}, \textcolor{blue}{\boldsymbol{\nu}}) \overset{\text{def}}{=} \min_{\textcolor{brown}{\boldsymbol{P}} \in U(\textcolor{red}{\boldsymbol{a}}, \textcolor{blue}{\boldsymbol{b}})} \langle \textcolor{brown}{\boldsymbol{P}}, M_{\textcolor{red}{\boldsymbol{X}}\textcolor{blue}{\boldsymbol{Y}}} \rangle + \gamma \mathrm{KL}\left(\textcolor{brown}{\boldsymbol{P}} \| \textcolor{red}{\boldsymbol{a}}\textcolor{blue}{\boldsymbol{b}}^T\right)$$

$$\mathrm{KL}\left(\textcolor{brown}{\boldsymbol{P}} \| \textcolor{red}{\boldsymbol{a}}\textcolor{blue}{\boldsymbol{b}}^T\right) = E(\textcolor{red}{\boldsymbol{a}}) + E(\textcolor{blue}{\boldsymbol{b}}) - E(\textcolor{brown}{\boldsymbol{P}})$$

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

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

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

$$W_\gamma(\textcolor{red}{\boldsymbol{\mu}}, \textcolor{blue}{\boldsymbol{\nu}}) \overset{\text{def}}{=} \min_{\boldsymbol{P} \in U(\textcolor{red}{\boldsymbol{a}}, \textcolor{blue}{\boldsymbol{b}})} \langle \textcolor{darkred}{\boldsymbol{P}}, M_{\textcolor{red}{\boldsymbol{X}}\textcolor{blue}{\boldsymbol{Y}}} \rangle + \gamma \mathrm{KL}\left(\textcolor{darkred}{\boldsymbol{P}} \| \textcolor{red}{\boldsymbol{a}}\textcolor{blue}{\boldsymbol{b}}^T\right)$$



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

18

# 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

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

# Fast & Scalable Algorithm

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

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

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

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

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

$$(\boldsymbol{\partial L / \partial P_{ij}} = 0) \Rightarrow P_{ij} = a_i e^{\frac{\alpha_i}{\gamma}} \; e^{-\frac{M_{ij}}{\gamma}} \; b_j e^{\frac{\beta_j}{\gamma}} = \boldsymbol{u_i} \; K_{ij} \boldsymbol{v_j}$$

# Fast & Scalable Algorithm

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

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

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

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

$$u \leftarrow a/Kv, \quad v \leftarrow b/K^T u$$

- $O(nm)$ complexity, GPGPU parallel [C'13].
- $O(n^{d+1})$ if $\Omega = \{1, \dots, n\}^d$ and $D^p$ separable.
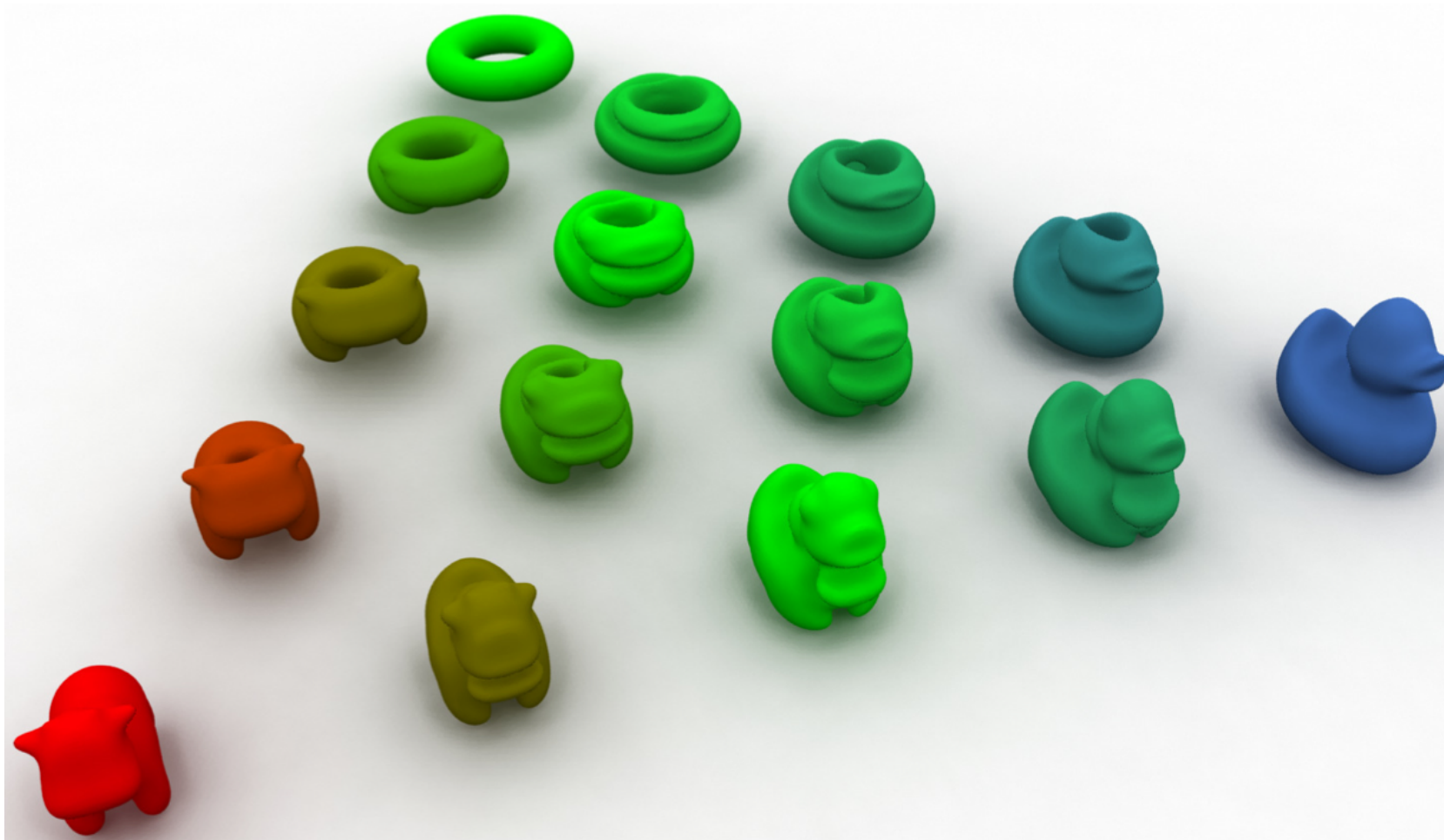
[S..C..'15]

19

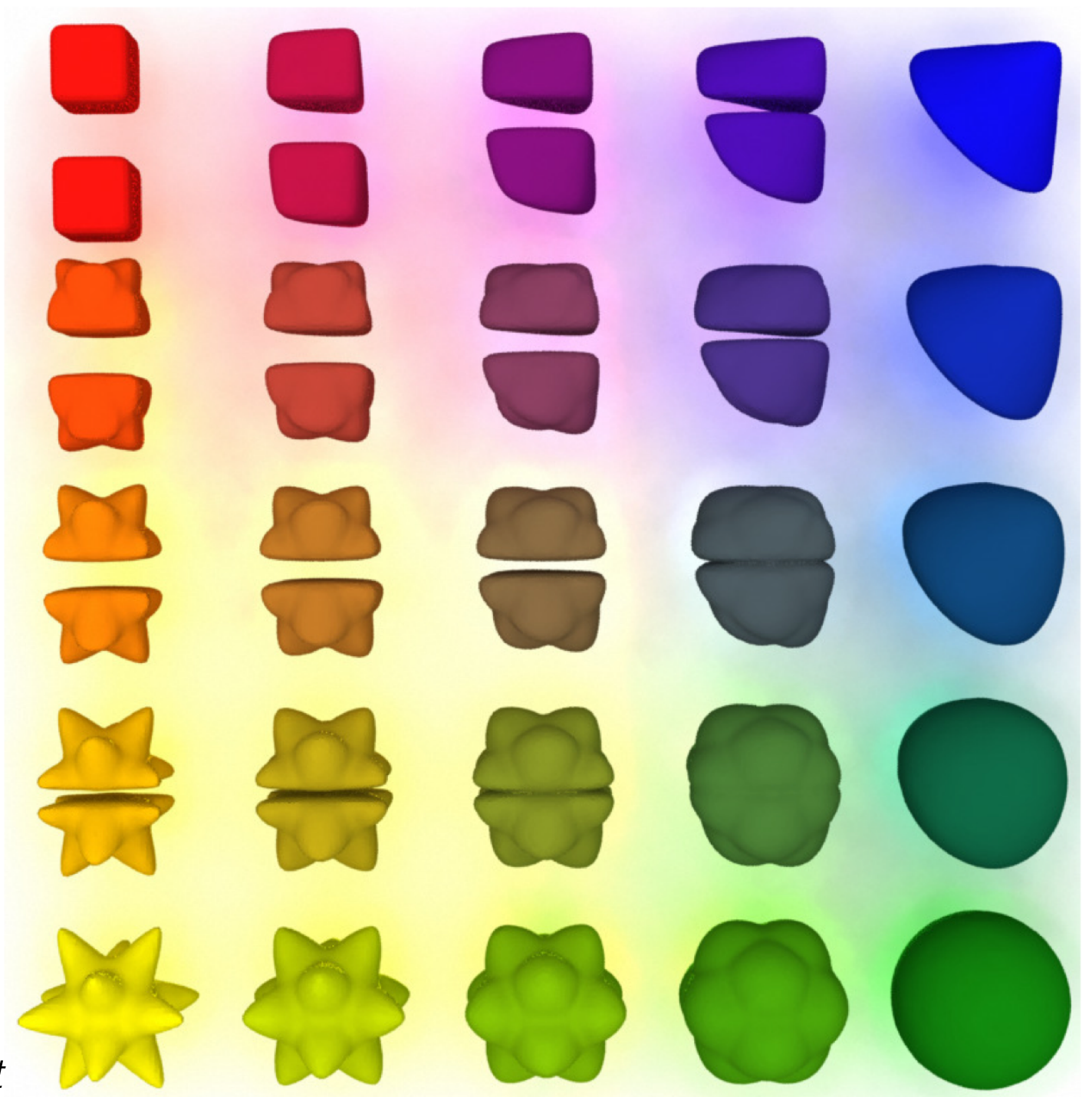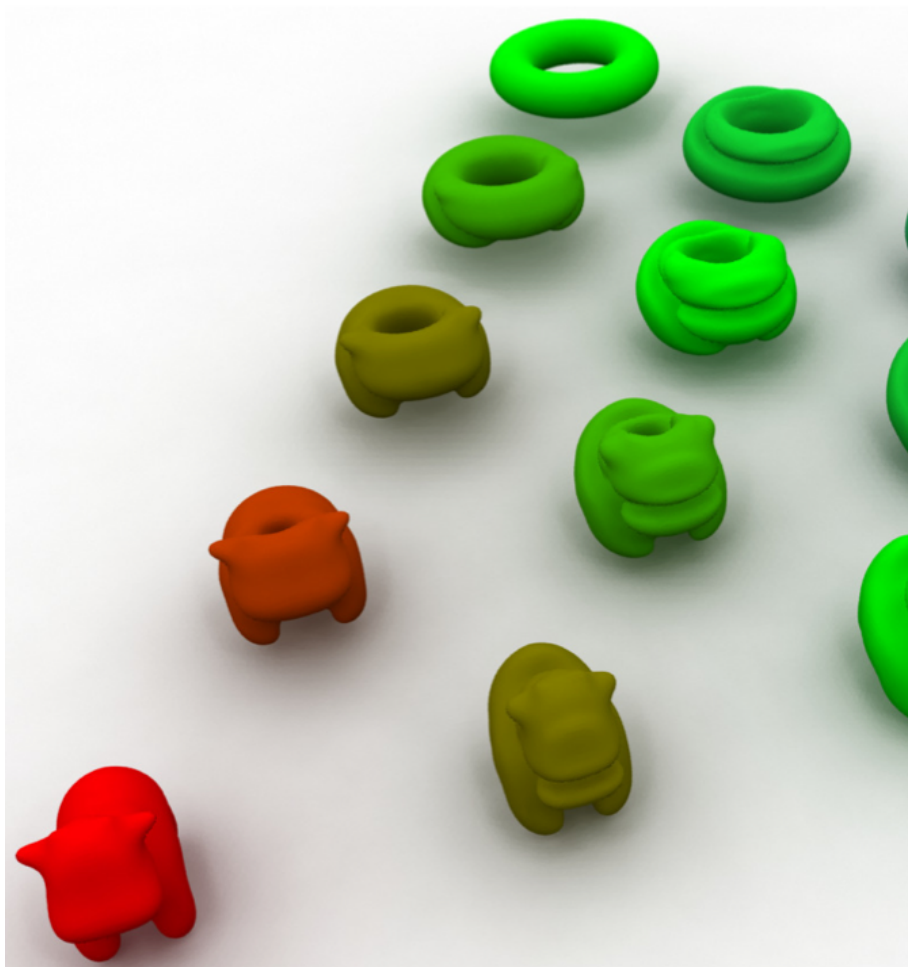# (Application: Barycenters)

# (Application: Barycenters)



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

20

# (Application: Barycenters)

[S..C..'15]

20

# (Application: Barycenters)

[S..C..'15]

20

# (Application: Barycenters)

[S..C..'15]

# (Application: Wasserstein Regression)



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

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

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

# (Application: Brain Regression)



Original       Euclidean       Wasserstein
projection      projection

# (Application: Brain Regression)

# Algorithmic Formulation

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

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

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

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

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

# Algorithmic Formulation

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

$$\underline{\mathrm{W}}_L(\textcolor{red}{\boldsymbol{\mu}}, \textcolor{blue}{\boldsymbol{\nu}}) \stackrel{\text{def}}{=} \gamma \boldsymbol{a}^T \log \boldsymbol{u}_L + \gamma \boldsymbol{b}^T \log \boldsymbol{v}_L,$$

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

**Prop.** $\frac{\partial \mathrm{W}_L}{\partial \textcolor{red}{\boldsymbol{X}}}, \frac{\partial \mathrm{W}_L}{\partial \textcolor{red}{\boldsymbol{a}}}$ can be computed recursively, in $O(L)$ kernel $\textcolor{blue}{K} \times$ vector products.

# Algorithmic Formulation

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

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

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

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

# 4. Algorithmic Formulation

**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
```

# *W* is versatile



*Discrete* - *Continuous*

*Continuous* - *Continuous*

# *W* is versatile



*Discrete* - *Continuous*

low dim.

**[M'116][KMB'16] [L'15]**

Stochastic
Optimization

**[GCPB'16]**

*Continuous* - *Continuous*

31

# $D$ transforms

$$W_p^p(\boldsymbol{\mu}, \boldsymbol{\nu}) = \sup_{\substack{\boldsymbol{\varphi} \in L_1(\boldsymbol{\mu}), \boldsymbol{\psi} \in L_1(\boldsymbol{\nu}) \\ \boldsymbol{\varphi}(x) + \boldsymbol{\psi}(y) \leq \boldsymbol{D}^p(x,y)}} \int \boldsymbol{\varphi} \, d\boldsymbol{\mu} + \int \boldsymbol{\psi} \, d\boldsymbol{\nu}.$$

DUAL

# $D$ transforms

$$W_p^p(\boldsymbol{\mu}, \boldsymbol{\nu}) = \sup_{\substack{\boldsymbol{\varphi} \in L_1(\boldsymbol{\mu}), \boldsymbol{\psi} \in L_1(\boldsymbol{\nu}) \\ \boldsymbol{\varphi}(x) + \boldsymbol{\psi}(y) \leq \boldsymbol{D}^p(x,y)}} \int \boldsymbol{\varphi} d\boldsymbol{\mu} + \int \boldsymbol{\psi} d\boldsymbol{\nu}.$$

For given $\boldsymbol{\varphi}$, cannot get a better $\psi$ than

$$\boldsymbol{\varphi}^{\boldsymbol{D}}(y) \overset{\text{def}}{=} \inf_x \boldsymbol{D}^p(x, y) - \boldsymbol{\varphi}(x).$$

$$W_p^p(\boldsymbol{\mu}, \boldsymbol{\nu}) = \sup_{\boldsymbol{\varphi}} \int \boldsymbol{\varphi} d\boldsymbol{\mu} + \int \boldsymbol{\varphi}^{\boldsymbol{D}} d\boldsymbol{\nu}.$$

# $D$ transforms

$$W_p^p(\boldsymbol{\mu}, \boldsymbol{\nu}) = \sup_{\boldsymbol{\varphi}} \int \boldsymbol{\varphi} \, d\boldsymbol{\mu} + \int \boldsymbol{\varphi}^{\boldsymbol{D}} \, d\boldsymbol{\nu}.$$

SEMI-DUAL

$$\boldsymbol{\varphi}^{\boldsymbol{D}}(y) \overset{\mathrm{def}}{=} \inf_{x} \boldsymbol{D}^p(x, y) - \boldsymbol{\varphi}(x).$$

$$\boldsymbol{\varphi}^{\boldsymbol{DD}}(x) = \inf_{y} \boldsymbol{D}^p(x, y) - \boldsymbol{\varphi}^{\boldsymbol{D}}(y).$$

$\boldsymbol{\varphi}$ is $\boldsymbol{D}$ concave if $\exists \boldsymbol{\phi} : \boldsymbol{\varphi} = \boldsymbol{\phi}^{\boldsymbol{D}}$

33

# $D$ transforms

$$\varphi^{D}(y) \overset{\text{def}}{=} \inf_{x} D^{p}(x, y) - \varphi(x).$$

$$\varphi^{DD}(x) = \inf_{y} D^{p}(x, y) - \varphi^{D}(y).$$

$$\varphi \text{ is } D \text{ concave if } \exists \phi : \varphi = \phi^{D}$$

$$W_{p}^{p}(\boldsymbol{\mu}, \boldsymbol{\nu}) = \sup_{\varphi \text{ is } D\text{-concave}} \int \varphi d\boldsymbol{\mu} + \int \varphi^{D} d\boldsymbol{\nu}.$$

SEMI-DUAL

# Reminder: dual regularization

$$W_p^p(\boldsymbol{\mu}, \boldsymbol{\nu}) = \sup_{\boldsymbol{\varphi}, \boldsymbol{\psi}} \int \boldsymbol{\varphi} \, d\boldsymbol{\mu} + \int \boldsymbol{\psi} \, d\boldsymbol{\nu} - \iota_C(\boldsymbol{\varphi}, \boldsymbol{\psi})$$

$$C = \{(\boldsymbol{\varphi}, \boldsymbol{\psi}) | \boldsymbol{\varphi} \oplus \boldsymbol{\psi} \le \boldsymbol{D}^p\}$$

*regularizing dual*    *constraints*    $\gamma > 0$

$$W_\gamma(\boldsymbol{\mu}, \boldsymbol{\nu}) = \sup_{\boldsymbol{\varphi}, \boldsymbol{\psi}} \int \boldsymbol{\varphi} \, d\boldsymbol{\mu} + \int \boldsymbol{\psi} \, d\boldsymbol{\nu} - \iota_C^\gamma(\boldsymbol{\varphi}, \boldsymbol{\psi})$$

$$\iota_C^\gamma(\boldsymbol{\varphi}, \boldsymbol{\psi}) = \gamma \iint e^{(\boldsymbol{\varphi} \oplus \boldsymbol{\psi} - \boldsymbol{D}^p)/\gamma} \, d\boldsymbol{\mu} \, d\boldsymbol{\nu}$$

# Smoothed $D$ transforms

$$W_\gamma(\boldsymbol{\mu}, \boldsymbol{\nu}) = \sup_{\boldsymbol{\varphi}, \boldsymbol{\psi}} \int \boldsymbol{\varphi} \, d\boldsymbol{\mu} + \int \boldsymbol{\psi} \, d\boldsymbol{\nu} - \iota_C^\gamma(\boldsymbol{\varphi}, \boldsymbol{\psi})$$

$$\iota_C^\gamma(\boldsymbol{\varphi}, \boldsymbol{\psi}) = \gamma \iint e^{(\boldsymbol{\varphi} \oplus \boldsymbol{\psi} - \boldsymbol{D}^p)/\gamma} \, d\boldsymbol{\mu} \, d\boldsymbol{\nu}$$

REGULARIZED DUAL

$$\nabla_{\boldsymbol{\psi}} = 0 \qquad \gamma > 0$$

$$W_\gamma(\boldsymbol{\mu}, \boldsymbol{\nu}) = \sup_{\boldsymbol{\varphi}} \int \boldsymbol{\varphi} \, d\boldsymbol{\mu} + \int \boldsymbol{\varphi}^{\boldsymbol{D}, \gamma} \, d\boldsymbol{\nu}.$$

$$\boldsymbol{\varphi}^{\boldsymbol{D}, \gamma}(\boldsymbol{y}) = -\gamma \log \int e^{\frac{\boldsymbol{\varphi}(\boldsymbol{x}) - \boldsymbol{D}(\boldsymbol{x}, \boldsymbol{y})^p}{\gamma}} \, d\boldsymbol{\mu}(\boldsymbol{x})$$

REGULARIZED SEMI-DUAL

# Regularized Semidual Wasserstein

$$W_\gamma(\textcolor{red}{\boldsymbol{\mu}}, \textcolor{blue}{\boldsymbol{\nu}}) = \sup_{\textcolor{red}{\boldsymbol{\varphi}}} \int \textcolor{red}{\boldsymbol{\varphi}} d\textcolor{red}{\boldsymbol{\mu}} + \int \textcolor{blue}{\boldsymbol{\varphi}}^{\textcolor{green}{\boldsymbol{D}},\gamma} d\textcolor{blue}{\boldsymbol{\nu}}.$$

$$\textcolor{green}{\boldsymbol{\varphi}}^{\textcolor{green}{\boldsymbol{D}},\gamma}(\textcolor{blue}{\boldsymbol{y}}) = -\gamma \log \int e^{\frac{\textcolor{red}{\boldsymbol{\varphi}(\boldsymbol{x})} - \textcolor{green}{\boldsymbol{D}(\boldsymbol{x},\boldsymbol{y})}^p}{\gamma}} d\textcolor{red}{\boldsymbol{\mu}(\boldsymbol{x})}$$

REGULARIZED SEMI-DUAL

substituting

$$\sup_{\textcolor{red}{\boldsymbol{\varphi}}} \int_{\textcolor{blue}{\boldsymbol{y}}} \left[ \int_{\textcolor{red}{\boldsymbol{x}}} \textcolor{red}{\boldsymbol{\varphi}(\boldsymbol{x})} d\textcolor{red}{\boldsymbol{\mu}(\boldsymbol{x})} - \gamma \log \int_{\textcolor{red}{\boldsymbol{x}}} e^{\frac{\textcolor{red}{\boldsymbol{\varphi}(\boldsymbol{x})} - \textcolor{green}{\boldsymbol{D}(\boldsymbol{x},\boldsymbol{y})}^p}{\gamma}} d\textcolor{red}{\boldsymbol{\mu}(\boldsymbol{x})} \right] d\textcolor{blue}{\boldsymbol{\nu}(\boldsymbol{y})}.$$

REGULARIZED SEMI-DUAL

# Semi-discrete case: Stochastic Opt.

$$\sup_{\boldsymbol{\varphi}} \int_{\boldsymbol{y}} \left[ \int_{\boldsymbol{x}} \boldsymbol{\varphi}(\boldsymbol{x}) d\boldsymbol{\mu}(\boldsymbol{x}) - \gamma \log \int_{\boldsymbol{x}} e^{\frac{\boldsymbol{\varphi}(\boldsymbol{x}) - D(\boldsymbol{x}, \boldsymbol{y})^p}{\gamma}} d\boldsymbol{\mu}(\boldsymbol{x}) \right] d\boldsymbol{\nu}(\boldsymbol{y}).$$

REGULARIZED SEMI-DUAL

# Semi-discrete case: Stochastic Opt.

$$\sup_{\boldsymbol{\varphi}} \int_{\boldsymbol{y}} \left[ \int_{\boldsymbol{x}} \boldsymbol{\varphi}(\boldsymbol{x}) d\boldsymbol{\mu}(\boldsymbol{x}) - \gamma \log \int_{\boldsymbol{x}} e^{\frac{\boldsymbol{\varphi}(\boldsymbol{x}) - D(\boldsymbol{x},\boldsymbol{y})^p}{\gamma}} d\boldsymbol{\mu}(\boldsymbol{x}) \right] d\boldsymbol{\nu}(\boldsymbol{y}).$$

REGULARIZED SEMI-DUAL

What if $\boldsymbol{\mu}$ is a discrete measure?   $\boldsymbol{\mu} = \sum_{i=1}^{n} \boldsymbol{a_i} \delta_{\boldsymbol{x_i}}$

$\boldsymbol{\varphi} \in L_1(\boldsymbol{\mu})$ is now just a vector $\boldsymbol{\alpha} \in \mathbb{R}^n$!

# Semi-discrete case: Stochastic Opt.

$$\sup_{\varphi} \int_{y} \left[ \int_{x} \varphi(x)d\mu(x) - \gamma \log \int_{x} e^{\frac{\varphi(x) - D(x,y)^p}{\gamma}} d\mu(x) \right] d\nu(y).$$

REGULARIZED SEMI-DUAL

What if $\mu$ is a discrete measure?   $\mu = \sum_{i=1}^{n} a_i \delta_{x_i}$

$\varphi \in L_1(\mu)$ is now just a vector $\alpha \in \mathbb{R}^n$!

$$\sup_{\alpha \in \mathbb{R}^n} \int_{y} \left[ \sum_{i=1}^{n} \alpha_i a_i - \gamma \log \sum_{i=1}^{n} e^{\frac{\alpha_i - D(x_i,y)^p}{\gamma}} a_i \right] d\nu(y)$$

$$= \sup_{\alpha \in \mathbb{R}^n} \mathbb{E}_{\nu}[f(\alpha, y)]$$

STOCHASTIC REGULARIZED SEMI-DUAL

# (in Discrete Setting)

$$\sup_{\boldsymbol{\varphi}} \int_{\boldsymbol{y}} \left[ \int_{\boldsymbol{x}} \boldsymbol{\varphi}(\boldsymbol{x}) d\boldsymbol{\mu}(\boldsymbol{x}) - \gamma \log \int_{\boldsymbol{x}} e^{\frac{\boldsymbol{\varphi}(\boldsymbol{x}) - D(\boldsymbol{x},\boldsymbol{y})^p}{\gamma}} d\boldsymbol{\mu}(\boldsymbol{x}) \right] d\boldsymbol{\nu}(\boldsymbol{y}).$$

REGULARIZED SEMI-DUAL

# (in Discrete Setting)

$$\sup_{\boldsymbol{\varphi}} \int_{\boldsymbol{y}} \left[ \int_{\boldsymbol{x}} \boldsymbol{\varphi}(\boldsymbol{x}) d\boldsymbol{\mu}(\boldsymbol{x}) - \gamma \log \int_{\boldsymbol{x}} e^{\frac{\boldsymbol{\varphi}(\boldsymbol{x}) - D(\boldsymbol{x}, \boldsymbol{y})^p}{\gamma}} d\boldsymbol{\mu}(\boldsymbol{x}) \right] d\boldsymbol{\nu}(\boldsymbol{y}).$$

REGULARIZED SEMI-DUAL

## What if $\boldsymbol{\nu}$ is **also** a discrete measure?

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

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

# (in Discrete Setting)

$$\sup_{\boldsymbol{\varphi}} \int_{\boldsymbol{y}} \left[ \int_{\boldsymbol{x}} \boldsymbol{\varphi}(\boldsymbol{x}) d\boldsymbol{\mu}(\boldsymbol{x}) - \gamma \log \int_{\boldsymbol{x}} e^{\frac{\boldsymbol{\varphi}(\boldsymbol{x}) - D(\boldsymbol{x}, \boldsymbol{y})^p}{\gamma}} d\boldsymbol{\mu}(\boldsymbol{x}) \right] d\boldsymbol{\nu}(\boldsymbol{y}).$$

REGULARIZED SEMI-DUAL

What if $\boldsymbol{\nu}$ is **also** a discrete measure?

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

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

$$\sup_{\boldsymbol{\alpha} \in \mathbb{R}^n} \int_y \left[ \sum_{i=1}^{n} \boldsymbol{\alpha_i} \boldsymbol{a_i} - \gamma \log \sum_{i=1}^{n} e^{\frac{\boldsymbol{\alpha_i} - D(\boldsymbol{x_i}, \boldsymbol{y})^p}{\gamma}} \boldsymbol{a_i} \right] d\boldsymbol{\nu}(\boldsymbol{y})$$

$$\sup_{\boldsymbol{\alpha} \in \mathbb{R}^n} \boldsymbol{\alpha}^T \boldsymbol{a} - \gamma \boldsymbol{b}^T \log K^T (\boldsymbol{a} \odot e^{\frac{\boldsymbol{\alpha}}{\gamma}})$$

REGULARIZED SEMI-DUAL

# Minimum Kantorovich Estimators

# Minimum Kantorovich Estimators

$$\mathcal{P}(\Omega) \qquad \boldsymbol{\nu}_{\mathrm{data}} \qquad \{p_\theta, \theta \in \Theta\}$$

$$p_{\boldsymbol{\theta}^\star}$$

$$\min_{\boldsymbol{\theta} \in \Theta} \mathrm{KL}(\boldsymbol{\nu}_{\mathrm{data}} \| \boldsymbol{p_\theta} \|) \qquad \text{MLE}$$

$$\min_{\boldsymbol{\theta} \in \Theta} W(\boldsymbol{\nu}_{\mathrm{data}}, \boldsymbol{p_\theta}) \qquad \text{MKE}$$

**[Bassetti'06]**

40

# In a discrete setting

- Suppose $\Omega$ is a discrete, finite space.

$$W_\gamma(p_{\boldsymbol{\theta}}, \boldsymbol{\nu}_{\text{data}}) = \max_{\boldsymbol{\alpha}, \boldsymbol{\beta}} \langle \boldsymbol{\alpha}, p_{\boldsymbol{\theta}} \rangle + \langle \boldsymbol{\beta}, \boldsymbol{\nu}_{\text{data}} \rangle - \gamma \langle e^{\boldsymbol{\alpha}/\gamma}, \ K e^{\boldsymbol{\beta}/\gamma} \rangle$$

$$\nabla_\theta W_\gamma = \left( \frac{\partial p_\theta}{\partial \theta} \right)^T \boldsymbol{\alpha}^\star$$

- Used for discrete models with very large state spaces in **[MMC'16]**.

- Considered for restricted Boltzmann machines, using stochastic approximation & regularization.

# In a continuous observation setting

# In a continuous observation setting

# In a continuous observation setting

$$W_\gamma(p_{\boldsymbol{\theta}}, \boldsymbol{\nu}_{\text{data}}) = \max_{\boldsymbol{f}, \boldsymbol{b}} \int_\Omega \boldsymbol{f} \, dp_{\boldsymbol{\theta}} + \boldsymbol{b}^T \mathbf{1}_m - \gamma \langle e^{\boldsymbol{f}/\gamma}, \ K e^{\boldsymbol{b}/\gamma} \rangle_{p_{\boldsymbol{\theta}}}$$



$\boldsymbol{\nu}_{\text{data}}$

# In a continuous observation setting

$$W_\gamma(p_{\boldsymbol{\theta}}, \boldsymbol{\nu}_{\text{data}}) = \max_{\boldsymbol{f}, \boldsymbol{b}} \int_\Omega \boldsymbol{f} \, dp_{\boldsymbol{\theta}} + \boldsymbol{b}^T \mathbf{1}_m - \gamma \langle e^{\boldsymbol{f}/\gamma}, \ K e^{\boldsymbol{b}/\gamma} \rangle_{p_{\boldsymbol{\theta}}}$$

$$\sup_{\boldsymbol{\beta} \in \mathbb{R}^m} \int_{\boldsymbol{x}} \left[ \sum_{i=1}^m \boldsymbol{\beta_j}/m - \gamma \log \frac{1}{m} \sum_{j=1}^m e^{\frac{\boldsymbol{\beta_j} - D^p(\boldsymbol{x}, \boldsymbol{y_j})}{\gamma}} \right] p_{\boldsymbol{\theta}}(\boldsymbol{x})$$

$$\sup_{\boldsymbol{\beta} \in \mathbb{R}^m} \mathbb{E}_{p_{\boldsymbol{\theta}}} [h(\boldsymbol{\beta}, \boldsymbol{x})]$$

$\boldsymbol{\nu}_{\text{data}}$

# In a continuous observation setting

$$W_\gamma(p_{\boldsymbol{\theta}}, \boldsymbol{\nu}_{\text{data}}) = \max_{\boldsymbol{f}, \boldsymbol{b}} \int_\Omega \boldsymbol{f} dp_{\boldsymbol{\theta}} + \boldsymbol{b}^T \mathbf{1}_m - \gamma \langle e^{\boldsymbol{f}/\gamma}, \ K e^{\boldsymbol{b}/\gamma} \rangle_{p_{\boldsymbol{\theta}}}$$

$$\sup_{\boldsymbol{\beta} \in \mathbb{R}^m} \int_{\boldsymbol{x}} \left[ \sum_{i=1}^m \boldsymbol{\beta_j}/m - \gamma \log \frac{1}{m} \sum_{j=1}^m e^{\frac{\boldsymbol{\beta_j} - D^p(\boldsymbol{x}, \boldsymbol{y_j})}{\gamma}} \right] p_{\boldsymbol{\theta}}(\boldsymbol{x})$$

$$\sup_{\boldsymbol{\beta} \in \mathbb{R}^m} \mathbb{E}_{p_{\boldsymbol{\theta}}}[h(\boldsymbol{\beta}, \boldsymbol{x})]$$

$$\boldsymbol{f}^\star = (\boldsymbol{b}^\star)^{D, \gamma} = \boldsymbol{x} \mapsto -\gamma \log \frac{1}{m} \sum_{i=1}^m e^{\frac{\boldsymbol{b}_j^\star - D^p(\boldsymbol{y_j}, \boldsymbol{x})}{\gamma}}$$
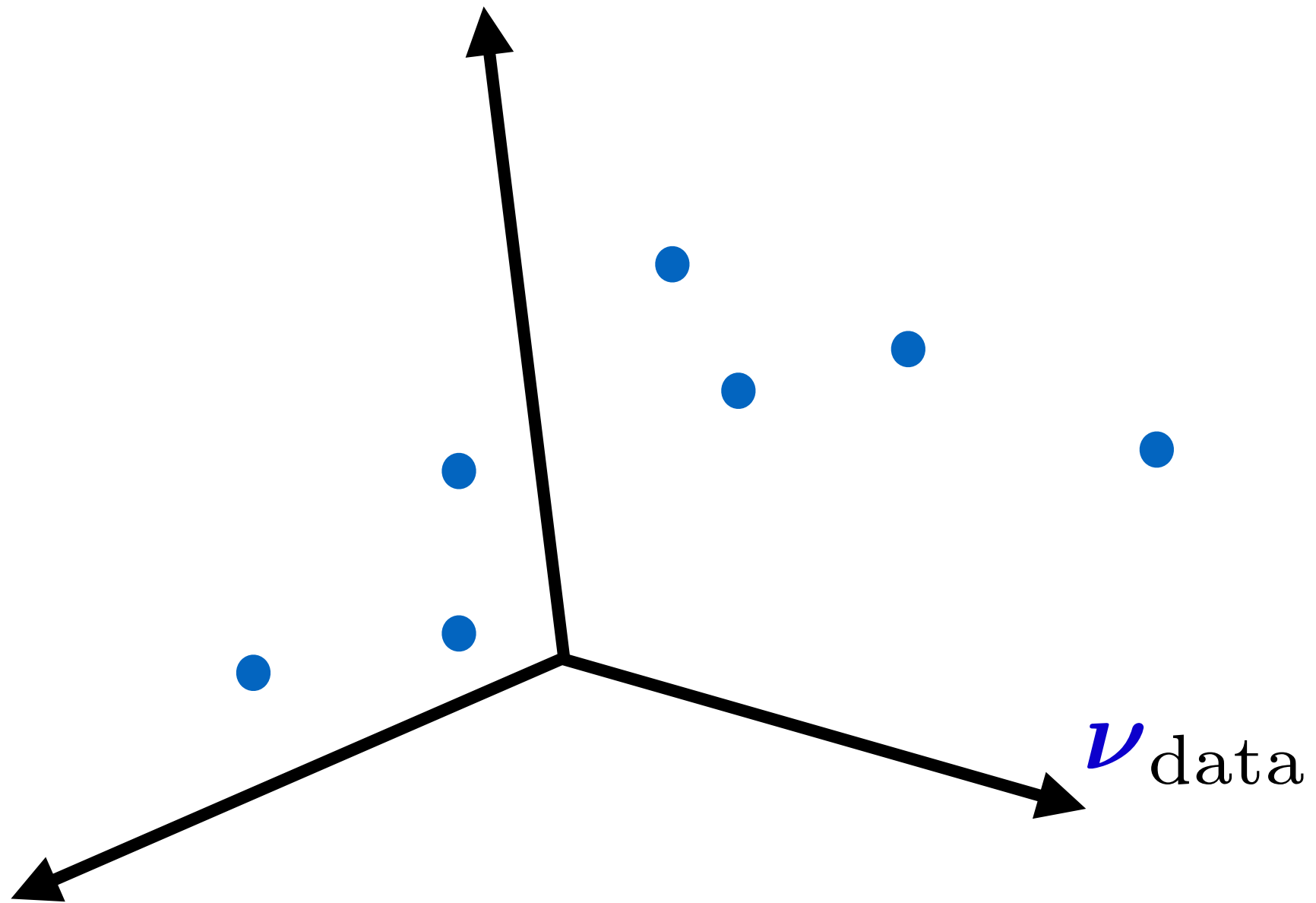
42

# In a continuous observation setting

$$W_\gamma(p_{\boldsymbol{\theta}}, \boldsymbol{\nu}_{\text{data}}) = \max_{\boldsymbol{f}, \boldsymbol{b}} \int_\Omega \boldsymbol{f} \, dp_{\boldsymbol{\theta}} + \boldsymbol{b}^T \mathbf{1}_m - \gamma \langle e^{\boldsymbol{f}/\gamma}, \; K e^{\boldsymbol{b}/\gamma} \rangle_{p_{\boldsymbol{\theta}}}$$

$$\sup_{\boldsymbol{\beta} \in \mathbb{R}^m} \int_{\boldsymbol{x}} \left[ \sum_{i=1}^m \boldsymbol{\beta_j}/m - \gamma \log \frac{1}{m} \sum_{j=1}^m e^{\frac{\boldsymbol{\beta_j} - D^p(\boldsymbol{x}, \boldsymbol{y_j})}{\gamma}} \right] p_{\boldsymbol{\theta}}(\boldsymbol{x})$$

$$\sup_{\boldsymbol{\beta} \in \mathbb{R}^m} \mathbb{E}_{p_{\boldsymbol{\theta}}}[h(\boldsymbol{\beta}, \boldsymbol{x})]$$

$$\boldsymbol{f}^\star = (\boldsymbol{b}^\star)^{D, \gamma} = \boldsymbol{x} \mapsto -\gamma \log \frac{1}{m} \sum_{i=1}^m e^{\frac{\boldsymbol{b}_j^\star - D^p(\boldsymbol{y_j}, \boldsymbol{x})}{\gamma}}$$

$$\nabla_\theta W_\gamma = \left( \frac{\partial p_\theta}{\partial \theta} \right)^* \boldsymbol{f}^\star$$

**[GCPB'16]**

42

# In a generative model setting



$\nu_{\text{data}}$
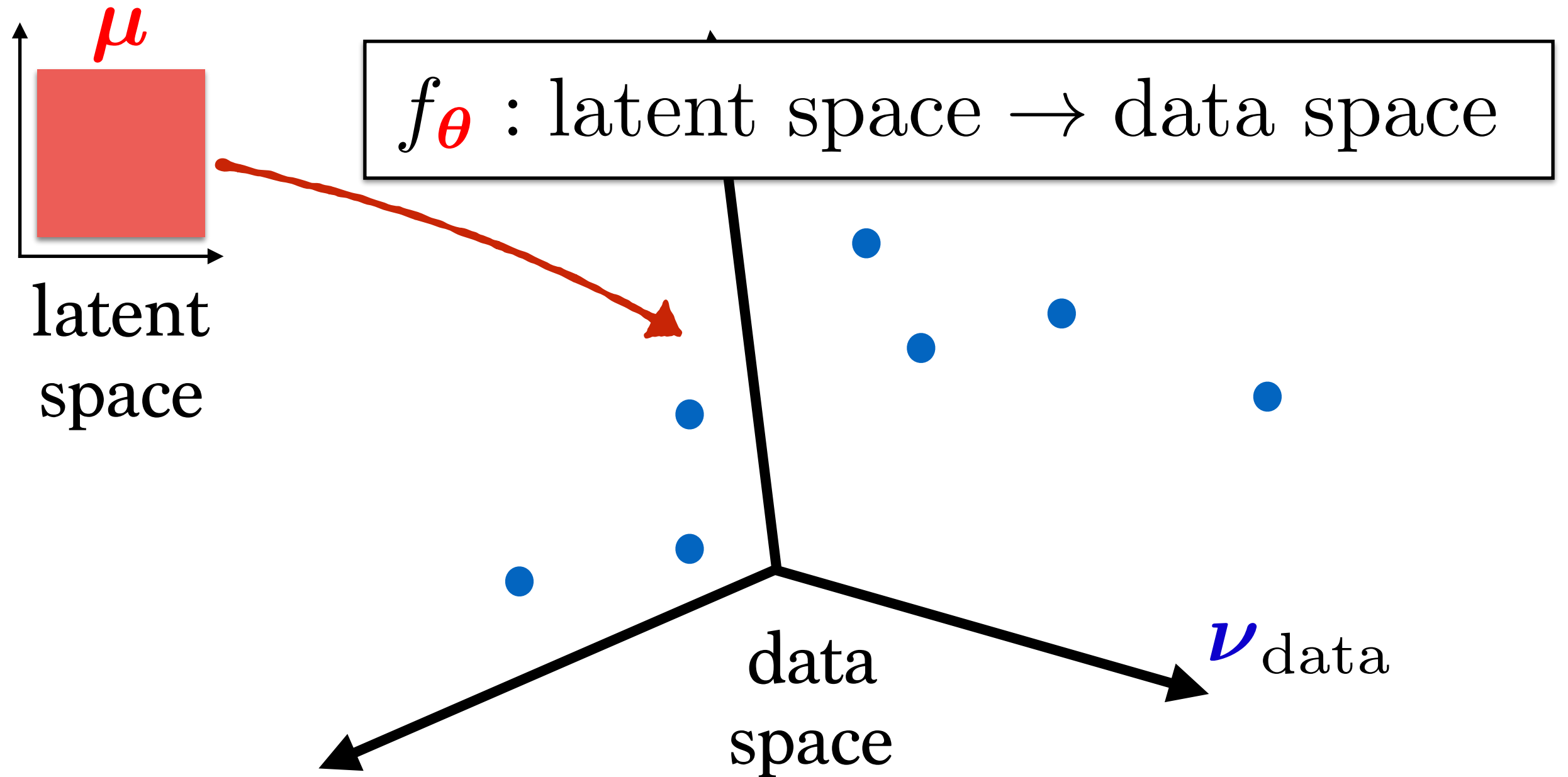
# In a generative model setting

# In a generative model setting



$f_{\boldsymbol{\theta}} :$ latent space $\rightarrow$ data space

$\boldsymbol{\mu}$

latent space

data space

$\boldsymbol{\nu}_{\text{data}}$

# In a generative model setting



$f_{\boldsymbol{\theta}}$ : latent space $\rightarrow$ data space

$\boldsymbol{\mu}$

latent
space

$f_{\boldsymbol{\theta}\sharp}\boldsymbol{\mu}$

$\boldsymbol{\nu}_{\mathrm{data}}$

data
space

# In a generative model setting



latent space

$f_{\boldsymbol{\theta}}$ : latent space $\rightarrow$ data space

$f_{\boldsymbol{\theta}\sharp}\boldsymbol{\mu}$

$\boldsymbol{\nu}_{\mathrm{data}}$

data space

MLE

$$\min_{\boldsymbol{\theta}\in\Theta} W(\boldsymbol{\nu}_{\mathrm{data}}, f_{\boldsymbol{\theta}\sharp}\boldsymbol{\mu})$$

GM-MKE
W-GAN

**[ACB'17]**

**[BGTSS'17]**

43

# Our algorithmic proposal
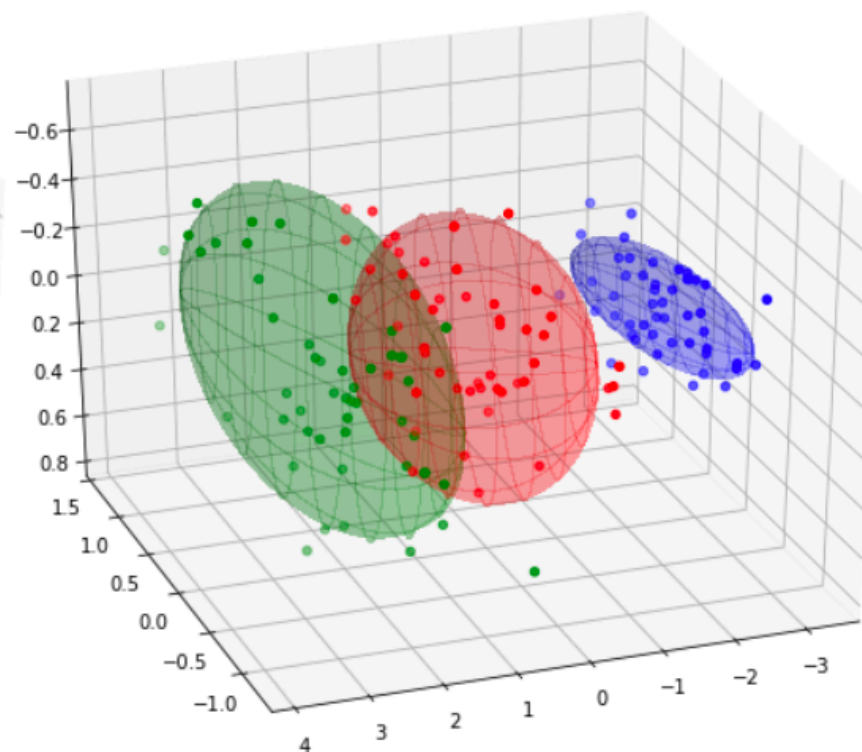
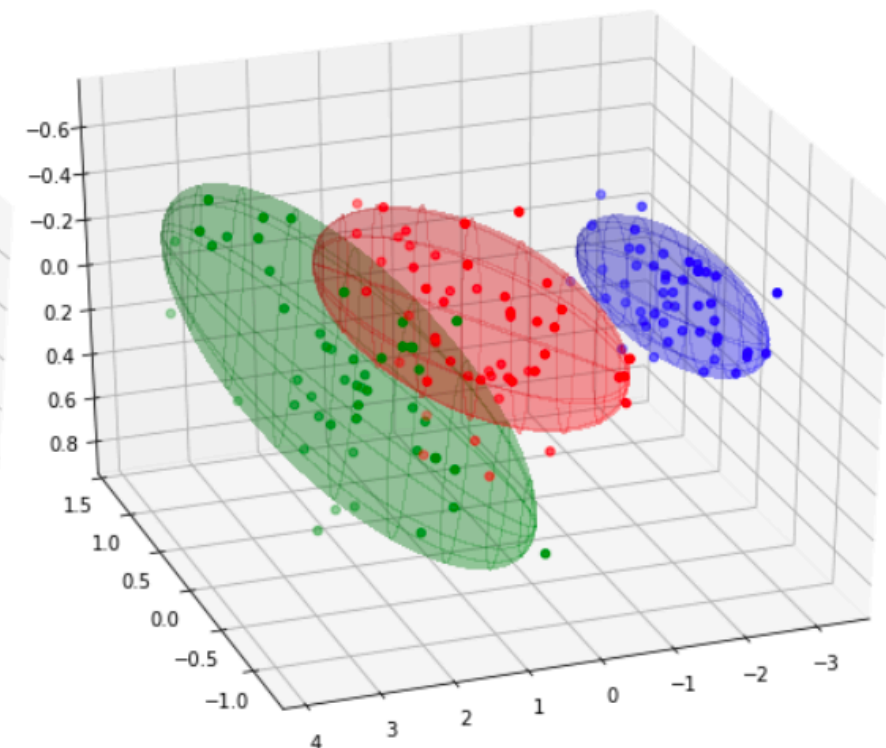Approximate regularized $W$ loss by $W_L$.

# Example: Fitting Ellipses

- $k$-means problem can be seen as a MKE when the model = atomic measures with $k$ atoms.

- We generalize by estimating uniform ellipsoid measures that approximate clouds of points.



(a) Initialization (unit balls, kmeans centers)

(b) After 3 gradient steps

(c) At convergence (15 steps)