# Empirical (Regularized) Optimal Transport: Statistical Theory and Applications

Optimal Transport, Topological Data Analysis and Applications to Shape and Machine Learning

Marcel Klatt<sup>1</sup>

Carla Tameling<sup>1</sup>

Yoav Zemel<sup>2</sup>

Axel Munk<sup>1</sup>

July 29, 2020

Institute for Mathematical Stochastics<sup>1</sup> University of Göttingen

Centre for Mathematical Sciences<sup>2</sup> University of Cambridge

#### Based on joint work with...



**Yoav Zemel** 



**Carla Tameling** 



Axel Munk



**M. Klatt, C. Tameling and A. Munk** *Empirical Regularized Optimal Transport: Statistical Theory and Applications*, SIAM Journal on Mathematics of Data Science (2020)

M. Klatt, A. Munk and Y. Zemel Limit Laws for Empirical Optimal Solutions in Stochastic Linear Programs, arXiv preprint (2020)

#### **Optimal Transport & Entropy Regularization**



Two probability measures

$$\mathbf{r} = \sum_{i} \mathbf{r}_{i} \delta_{x_{i}} \qquad \mathbf{s} = \sum_{j} \mathbf{s}_{j} \delta_{x_{j}}$$

Transport costs

$$d_{i,j}^p := d^p(x_i, x_j), \quad p \ge 1$$

Task: Find the most efficient way to transport one measure into the other.

$$\min_{\pi \in \Pi(\mathbf{r},\mathbf{s})} \sum_{i,j} d^p_{ij} \pi_{ij}$$

#### **Optimal Transport & Entropy Regularization**



**M. Cuturi** Sinkhorn Distances: Lightspeed Computation of Optimal Transportation Distances, Advances in neural information processing systems (2013)

**G. Peyré and M. Cuturi** *Computational Optimal Transport* Foundations and Trends in Machine Learning (2019)

#### **Optimal Transport & Entropy Regularization**



Finite discrete metric space  $(\mathcal{X}, d)$ 

Wasserstein distance

$$W_p(\mathbf{r}, \mathbf{s}) = (\min_{\pi \in \Pi(\mathbf{r}, \mathbf{s})} \sum_{i,j} d_{ij}^p \pi_{ij})^{\frac{1}{p}}$$

OT plan

$$\pi(\mathbf{r}, \mathbf{s}) \in \arg\min_{\pi \in \Pi(\mathbf{r}, \mathbf{s})} \sum_{i,j} d_{ij}^p \pi_{ij}$$

Sinkhorn divergence

$$W_{p,\lambda}(\mathbf{r},\mathbf{s}) = (\min_{\pi \in \Pi(\mathbf{r},\mathbf{s})} \sum_{i,j} d_{ij}^p \pi_{ij} - \lambda E(\pi))^{\frac{1}{p}}$$

ROT plan

**M. Cuturi** Sinkhorn Distances: Lightspeed Computation of Optimal Transportation Distances, Advances in neural information processing systems (2013)

**G. Peyré and M. Cuturi** *Computational Optimal Transport* Foundations and Trends in Machine Learning (2019)

$$\pi_{\lambda}(\mathbf{r}, \mathbf{s}) = \arg \min_{\pi \in \Pi(\mathbf{r}, \mathbf{s})} \sum_{i, j} d_{ij}^{p} \pi_{ij} - \lambda E(\pi)$$

#### In this talk, we focus on statistical properties of these quantities!

Suppose, we observe independent and identically distributed random variables

$$X_1, \ldots, X_n \stackrel{i.i.d.}{\sim} r$$
.

This naturally leads to estimate  $\Gamma$  by the empirical measure

$$\hat{\boldsymbol{r}}_{\boldsymbol{n}} = \frac{1}{n} \sum_{i=1}^{n} \delta_{X_i}$$

#### **Central Question:**

How do these random quantities relate to their population counterpart, respectively?

Wasserstein distance

$$W_{pp}((\mathbf{r},\mathbf{s})) = ((\min_{\pi \in \Pi(\widehat{r}_{n},s)}\sum_{i,j}d_{ij}^{p}\pi_{ij}))^{\frac{1}{p}}$$

OT plan

$$\pi\pi(\mathbf{r},\mathbf{s}) \in \arg\min_{\pi \in \Pi(\hat{\mathbf{r}},\mathbf{s})} \sum_{i \neq j \neq j} d_{i \neq j}^{p} \pi_{i \neq j}$$

Sinkhorn divergence

$$WV_{p,p,k}((r,,s)) = (\min_{\pi \in \Pi(\hat{r}_{n},s)} \sum_{i \in jj} d^{p}_{ijj} \pi_{ijj} - \lambda H((\pi)))^{\frac{1}{p}}$$

ROT plan

$$\pi_{\mathcal{A}}(\mathbf{r},s) = \arg\min_{\pi \in \Pi(\hat{\mathbf{r}},s)} \sum_{i \neq j} \mathcal{A}_{ijj}^{p} \pi_{ijj} - \mathcal{A}_{H}(\pi)$$

#### In this talk, we focus on statistical properties of these quantities!

Goal:

Quantify random fluctuation by asymptotic limit laws!

Wasserstein distance

$$a_n\left\{W_p^p(\hat{r}_n, s) - W_p^p(r, s)\right\} \xrightarrow{\mathscr{D}} ??$$

OT plan

$$a_n\left\{\pi(\hat{r}_n, s) - \pi(r, s)\right\} \xrightarrow{\mathscr{D}} ??$$

Sinkhorn divergence

$$a_n\left\{W_{p,\lambda}^p(\hat{r}_n,s) - W_{p,\lambda}^p(r,s)\right\} \xrightarrow{\mathscr{D}} ??$$

ROT plan

$$a_n\left\{\pi_{\lambda}(\hat{\boldsymbol{r}}_n,\boldsymbol{s})-\pi_{\lambda}(\boldsymbol{r},\boldsymbol{s})\right\} \xrightarrow{\mathscr{D}} ??$$

7

#### In this talk, we focus on statistical properties of these quantities!

# Distributional limit laws for the Wasserstein distance are well understood:



**M. Sommerfeld and A. Munk** Inference for empirical Wasserstein distances on finite spaces., Journal of the Royal Statistical Society: Series B (Methodological) (2018)

**C. Tameling, M. Sommerfeld and A. Munk** *Empirical optimal transport on countable metric spaces: Distributional limits and statistical applications*, Annals of Applied Probability (2019)

#### • $\Gamma_{\star}$ set of dual solutions

$$\max_{\alpha,\beta\in\mathbb{R}^{N}}\sum_{i}\alpha_{i}r_{i}+\sum_{j}\beta_{j}s_{j}$$
  
s.t.  $\alpha_{i}+\beta_{j}\leq d_{ij}^{p}$ 

• G the Gaussian limit of  $\sqrt{n} (\hat{r}_n - r)$ 

#### Wasserstein distance

$$\sqrt{n} \left\{ W_p^p(\hat{\boldsymbol{r}}_n, \boldsymbol{s}) - W_p^p(\boldsymbol{r}, \boldsymbol{s}) \right\} \stackrel{\mathcal{D}}{\to} \max_{\alpha_\star \in \Gamma_\star} \langle G, \alpha_\star \rangle$$

OT plan

$$a_n\left\{\pi(\hat{r}_n, s) - \pi(r, s)\right\} \xrightarrow{\mathscr{D}} ??$$

Sinkhorn divergence

$$a_n \left\{ W_{p,\lambda}^p(\hat{\boldsymbol{r}}_n, \boldsymbol{s}) - W_{p,\lambda}^p(\boldsymbol{r}, \boldsymbol{s}) \right\} \xrightarrow{\mathscr{D}} ??$$

ROT plan

$$a_n\left\{\pi_{\lambda}(\hat{\boldsymbol{r}}_n,\boldsymbol{s})-\pi_{\lambda}(\boldsymbol{r},\boldsymbol{s})\right\} \xrightarrow{\mathscr{D}} ??$$

**ROT** plan

$$\pi_{\lambda}(\hat{r}_n, s) = \arg\min_{\pi \in \Pi(\hat{r}_n, s)} \sum_{i,j} d_{ij}^p \pi_{ij} - \lambda E(\pi), \quad \lambda > 0$$

Theorem (K., Tameling & Munk (2020)):

With the sample size *n* approaching infinity, it holds that

$$\sqrt{n} \left\{ \pi_{\lambda}(\hat{\boldsymbol{r}}_{n}, \boldsymbol{s}) - \pi_{\lambda}(\boldsymbol{r}, \boldsymbol{s}) \right\} \xrightarrow{\mathcal{D}} \mathcal{N}_{N^{2}} \left( 0, \Sigma_{\lambda}(\boldsymbol{r} \mid \boldsymbol{s}) \right)$$

# Why Gaussian fluctuation?

$$\sqrt{n} \left\{ \pi_{\lambda}(\hat{\boldsymbol{r}}_{n}, \boldsymbol{s}) - \pi_{\lambda}(\boldsymbol{r}, \boldsymbol{s}) \right\} \xrightarrow{\mathcal{D}} \mathcal{N}_{N^{2}} \left( 0, \Sigma_{\lambda}(\boldsymbol{r} \mid \boldsymbol{s}) \right)$$



# **Further consequences**

$$\sqrt{n} \left\{ \pi_{\lambda}(\hat{\boldsymbol{r}}_{n}, \boldsymbol{s}) - \pi_{\lambda}(\boldsymbol{r}, \boldsymbol{s}) \right\} \xrightarrow{\mathscr{D}} \mathcal{N}_{N^{2}} \left( 0, \Sigma_{\lambda}(\boldsymbol{r} \mid \boldsymbol{s}) \right)$$

Limit distributions for...

• Sinkhorn divergence

$$\sqrt{n}\left\{W_{p,\lambda}^{p}(\hat{\boldsymbol{r}}_{n},\boldsymbol{s})-W_{p,\lambda}^{p}(\boldsymbol{r},\boldsymbol{s})\right\}\overset{\mathcal{D}}{\rightarrow}\langle G,\alpha_{\lambda}\rangle$$

• 
$$\alpha_{\lambda} = \arg \max_{\alpha, \beta \in \mathbb{R}^{\mathcal{X}}} \alpha^{T} r + \beta^{T} s - \lambda \sum_{i,j} \exp\left(\frac{\alpha_{i} + \beta_{j} - d_{ij}^{p}}{\lambda}\right) - 1.$$

• *G* the Gaussian limit of  $\sqrt{n} (\hat{r}_n - r)$ 

**J. Bigot, E. Cazelles and N. Papadakis** Central limit theorems for entropy-regularized optimal transport on finite spaces and statistical applications, Electronic Journal of Statistics (2019)



#### **Further consequences**

 $\sqrt{n} \left\{ \pi_{\lambda}(\hat{\boldsymbol{r}}_{n}, \boldsymbol{s}) - \pi_{\lambda}(\boldsymbol{r}, \boldsymbol{s}) \right\} \xrightarrow{\mathscr{D}} \mathcal{N}_{N^{2}} \left( 0, \Sigma_{\lambda}(\boldsymbol{r} \mid \boldsymbol{s}) \right)$ 

Limit distributions for...

Sinkhorn divergence

$$\sqrt{n}\left\{W_{p,\lambda}^{p}(\hat{\boldsymbol{r}}_{n},\boldsymbol{s})-W_{p,\lambda}^{p}(\boldsymbol{r},\boldsymbol{s})\right\}\overset{\mathcal{D}}{\to}\langle G,\alpha_{\lambda}\rangle$$

... and many more related quantities:

- Sinkhorn divergence II
- Sinkhorn loss

$$\begin{split} \tilde{W}_{p,\lambda}^{p}(\boldsymbol{r},\boldsymbol{s}) &:= \sum_{i,j} d_{ij}^{p} \pi_{\lambda}(\boldsymbol{r},\boldsymbol{s}) \\ S_{p,\lambda}^{p}(\boldsymbol{r},\boldsymbol{s}) &:= W_{p,\lambda}^{p}(\boldsymbol{r},\boldsymbol{s}) - \frac{1}{2} \left( W_{p,\lambda}^{p}(\boldsymbol{r},\boldsymbol{s}) + W_{p,\lambda}^{p}(\boldsymbol{r},\boldsymbol{s}) \right) \end{split}$$

$$R^{p}_{p,\lambda}(\boldsymbol{r},\boldsymbol{s}) := 2S^{p}_{p,\lambda}(\boldsymbol{r},\boldsymbol{s}) - S^{p}_{p,\sqrt{2}\lambda}(\boldsymbol{r},\boldsymbol{s})$$

**J. Bigot, E. Cazelles and N. Papadakis** Central limit theorems for entropy-regularized optimal transport on finite spaces and statistical applications, Electronic Journal of Statistics (2019)



**A. Genevay, G. Peyré and M. Cuturi** *Learning generative models with Sinkhorn divergences*, Proceedings of the 21st International Conference on Artificial Intelligence and Statistics (2018)

L. Chizat, P. Roussillon, F. Léger, F.-X. Vialard and G. Peyré Faster Wasserstein distance estimation with Sinkhorn divergence, arXiv preprint (2020)

# **Further consequences**

 $\sqrt{n} \left\{ \pi_{\lambda}(\hat{\boldsymbol{r}}_{n}, \boldsymbol{s}) - \pi_{\lambda}(\boldsymbol{r}, \boldsymbol{s}) \right\} \xrightarrow{\mathcal{D}} \mathcal{N}_{N^{2}} \left( 0, \Sigma_{\lambda}(\boldsymbol{r} \mid \boldsymbol{s}) \right)$ 

Limit distributions for...

Sinkhorn divergence

$$\sqrt{n}\left\{W_{p,\lambda}^{p}(\hat{\boldsymbol{r}}_{n},\boldsymbol{s})-W_{p,\lambda}^{p}(\boldsymbol{r},\boldsymbol{s})\right\}\overset{\mathcal{D}}{\rightarrow}\langle G,\alpha_{\lambda}\rangle$$

There are extensions to countable metric spaces...



**S. Hundrieser, M. Klatt and A. Munk** *Limit Distributions for Entropic Optimal Transport on Countable Discrete Spaces*, In preparation

**J. Bigot, E. Cazelles and N. Papadakis** Central limit theorems for entropy-regularized optimal transport on finite spaces and statistical applications, Electronic Journal of Statistics (2019)



**A. Genevay, G. Peyré and M. Cuturi** *Learning generative models with Sinkhorn divergences*, Proceedings of the 21st International Conference on Artificial Intelligence and Statistics (2018)

L. Chizat, P. Roussillon, F. Léger, F.-X. Vialard and G. Peyré Faster Wasserstein distance estimation with Sinkhorn divergence, arXiv preprint (2020)

### **Application: Colocalization Analysis**

$$\sqrt{n} \left\{ \pi_{\lambda}(\hat{\boldsymbol{r}}_{n}, \boldsymbol{s}) - \pi_{\lambda}(\boldsymbol{r}, \boldsymbol{s}) \right\} \stackrel{\mathcal{D}}{\to} \mathcal{N}_{N^{2}} \left( 0, \Sigma_{\lambda}(\boldsymbol{r} \mid \boldsymbol{s}) \right)$$



(a) ATP Synthase



(b) MIC60

Multiscale Cluster of Excellence From Molecular Machines to Networks of Excitable Cells

**Figure 1:** Staining of two different proteins (Jakobs lab, Department of NanoBiophotonics, Max-Planck Institute for Biophysical Chemistry, Göttingen)

Aim: Analyse the interaction between fluorescently-labeled molecules by quantifying the co-occurrence and correlation between them!

$$\sqrt{n} \left\{ \pi_{\lambda}(\hat{\boldsymbol{r}}_{n}, \boldsymbol{s}) - \pi_{\lambda}(\boldsymbol{r}, \boldsymbol{s}) \right\} \stackrel{\mathcal{D}}{\to} \mathcal{N}_{N^{2}} \left( 0, \Sigma_{\lambda}(\boldsymbol{r} \mid \boldsymbol{s}) \right)$$

#### **Conventional Methods:**





(a) ATP Synthase

(b) MIC60

**Pixel based** intensity correlation analysis (Pearsons's correlation coefficient) or co-occurence (Manders' split coefficients)



These methods are very sensitive to the resolution of the images to be compared!

Figure 2: Illustration of Confocal and STED images of two proteins which are located at a distance of 45nm. The resolution of the confocal image is 244nm and for the STED image it is 40nm.



$$\sqrt{n} \left\{ \pi_{\lambda}(\hat{\boldsymbol{r}}_{n}, \boldsymbol{s}) - \pi_{\lambda}(\boldsymbol{r}, \boldsymbol{s}) \right\} \stackrel{\mathcal{D}}{\to} \mathcal{N}_{N^{2}} \left( 0, \Sigma_{\lambda}(\boldsymbol{r} \mid \boldsymbol{s}) \right)$$

#### Our approach:

- Compute the entropy regularized transport plan
- Colocalization measure **RCol** based on (regularized) optimal transport for  $t \in [0, diam(\mathcal{X})]$





Theorem (K., Tameling & Munk (2020)):

As sample size *n* approaches infinity

$$\sqrt{n} \left\{ \hat{\mathbf{RCol}}_n - \hat{\mathbf{RCol}} \right\} \xrightarrow{\mathcal{D}} \hat{\mathbf{RCol}}(G)$$

with *G* the random variable with distribution given by the limit law for the empirical regularized OT plan.

17

# **Application: Colocalization Analysis**

 $\sqrt{n} \left\{ \pi_{\lambda}(\hat{\boldsymbol{r}}_{n}, \boldsymbol{s}) - \pi_{\lambda}(\boldsymbol{r}, \boldsymbol{s}) \right\} \xrightarrow{\mathcal{D}} \mathcal{N}_{N^{2}} \left( 0, \Sigma_{\lambda}(\boldsymbol{r} \mid \boldsymbol{s}) \right)$ 

$$\sqrt{n} \left\{ \hat{\mathbf{RCol}}_n - \hat{\mathbf{RCol}} \right\} \xrightarrow{\mathscr{D}} \hat{\mathbf{RCol}}(G)$$

This yields  $1 - \alpha$  approximate uniform confidence bands, i.e.,

$$\lim_{n\to\infty}\mathbb{P}\left(\mathbf{RCol}\in I_n\right)=1-\alpha,$$

where

$$I_n := \left[ -\frac{\mathfrak{u}_{1-\alpha}}{\sqrt{n}} + \mathbf{R}\hat{\mathbf{C}}\mathbf{ol}_n, \frac{\mathfrak{u}_{1-\alpha}}{\sqrt{n}} + \mathbf{R}\hat{\mathbf{C}}\mathbf{ol}_n \right]$$

and  $\mathfrak{u}_{1-\alpha}$  is the  $1-\alpha$  quantile from the distribution  $\|\mathbf{RCol}(G)\|_{\infty}$ .

The quantile  $\mathfrak{u}_{1-\alpha}$  can be consistently approximated by its *n* out of *n* bootstrap analogue.

(a) ATP Synthase  
(b) MIC60  

$$x$$
  
 $r$   
 $r$   
 $\pi_{\lambda}(r, s)$   
 $RCol(\pi_{\lambda}(r, s))(t) = \sum_{i,j} \pi_{\lambda}(r, s)_{ij} 1_{\{d_{ij} \leq t\}}$   
 $\int_{i,j} \int_{i,j} \int_{i,$ 

#### **Application: Colocalization Analysis**













empirical OT plan? ( $\lambda = 0$ )





 $\lambda = 0$ Limit Theory for the OT plan





 $\min_{\pi \in \Pi(\mathbf{r}, \mathbf{s})} \sum_{i, i} d_{ij}^p \pi_{ij} \qquad \max_{\alpha, \beta \in \mathbb{R}^N} \sum_i \alpha_i \mathbf{r}_i + \sum_j \beta_j \mathbf{s}_j$  $s.t. \quad \alpha_i + \beta_j \le d_{ij}^p$ 

Assumption I : Unique OT plan

Assumption II : Optimal dual solutions nondegenerate



$$\sqrt{n} \left\{ \pi(\hat{\boldsymbol{r}}_n, \boldsymbol{s}) - \pi(\boldsymbol{r}, \boldsymbol{s}) \right\} \xrightarrow{\mathscr{D}} \sum_{k=1}^K \mathbb{1}_{G \in H_k \setminus \bigcup_{j < k} H_j} \pi(I_k, G)$$

OT plan degenerate



#### OT plan nondegenerate





Limit Theory for the OT plan  $\lambda = 0$ 



Dual

 $\min_{\pi \in \Pi(\mathbf{r}, \mathbf{s})} \sum_{i,j} d_{ij}^p \pi_{ij} \qquad \max_{\alpha, \beta \in \mathbb{R}^N} \sum_i \alpha_i \mathbf{r}_i + \sum_j \beta_j \mathbf{s}_j$  $s \cdot t \cdot \alpha_i + \beta_j \leq d_{ij}^p$ 

Assumption I : Unique OT plan

Assumption II : Optimal dual solutions nondegenerate

Theorem (K., Munk & Zemel (2020)): As sample size *n* approaches infinity

$$\sqrt{n} \left\{ \pi(\hat{\mathbf{r}}_{n}, \mathbf{s}) - \pi(\mathbf{r}, \mathbf{s}) \right\} \xrightarrow{\mathscr{D}}_{\mathscr{K}} \sum_{\mathscr{K}} \mathbb{1}_{G \in H_{\mathscr{K}} \setminus \bigcup_{k \notin \mathscr{K}} H_{j}} \alpha^{\mathscr{K}} \otimes \pi(I_{\mathscr{K}}, G)$$



Limit Theory for the OT plan  $\lambda = 0$ 



 $\min_{\pi \in \Pi(\mathbf{r},\mathbf{s})} \sum_{i,j} d^p_{ij} \pi_{ij}$ 

$$\max_{\alpha,\beta\in\mathbb{R}^{N}}\sum_{i}\alpha_{i}r_{i}+\sum_{j}\beta_{j}s_{j}$$
  
s.t.  $\alpha_{i}+\beta_{j}\leq d_{ij}^{p}$ 

Assumption I: Unique OT plan

Assumption II : Optimal dual solutions nondegenerate





# **Overview: (R)OT limit laws for finite metric spaces**

 $\lambda = 0$ 

OT plan

$$\sqrt{n}\left\{\pi(\hat{\boldsymbol{r}}_{n},\boldsymbol{s})-\pi(\boldsymbol{r},\boldsymbol{s})\right\} \xrightarrow{\mathcal{D}} \sum_{\mathcal{K}} \mathbb{1}_{G \in H_{\mathcal{K}} \setminus \bigcup_{k \notin \mathcal{K}} H_{j}} \alpha^{\mathcal{K}} \otimes \pi(I_{\mathcal{K}},G)$$

Wasserstein distance

$$\sqrt{n} \left\{ W_p^p(\hat{\boldsymbol{r}}_n, \boldsymbol{s}) - W_p^p(\boldsymbol{r}, \boldsymbol{s}) \right\} \stackrel{\mathcal{D}}{\to} \max_{\alpha_\star \in \Gamma_\star} \langle G, \alpha_\star \rangle$$



 $\begin{aligned} \lambda &> 0\\ \text{ROT plan} \\ \sqrt{n} \left\{ \pi_{\lambda}(\hat{\boldsymbol{r}}_{n}, \boldsymbol{s}) - \pi_{\lambda}(\boldsymbol{r}, \boldsymbol{s}) \right\} \xrightarrow{\mathcal{D}} \mathcal{N}_{N^{2}} \left( 0, \Sigma_{\lambda}(\boldsymbol{r} \mid \boldsymbol{s}) \right). \end{aligned}$ 

Sinkhorn divergence

$$\sqrt{n}\left\{W_{p,\lambda}^{p}(\hat{\boldsymbol{r}}_{n},\boldsymbol{s})-W_{p,\lambda}^{p}(\boldsymbol{r},\boldsymbol{s})\right\}\overset{\mathcal{D}}{\to}\langle G,\alpha_{\lambda}\rangle$$

