# Inverse optimal transport and related problems

Clarice Poon (University of Warwick)

Joint work with: Francisco Andrade and Gabriel Peyré (ENS Paris)
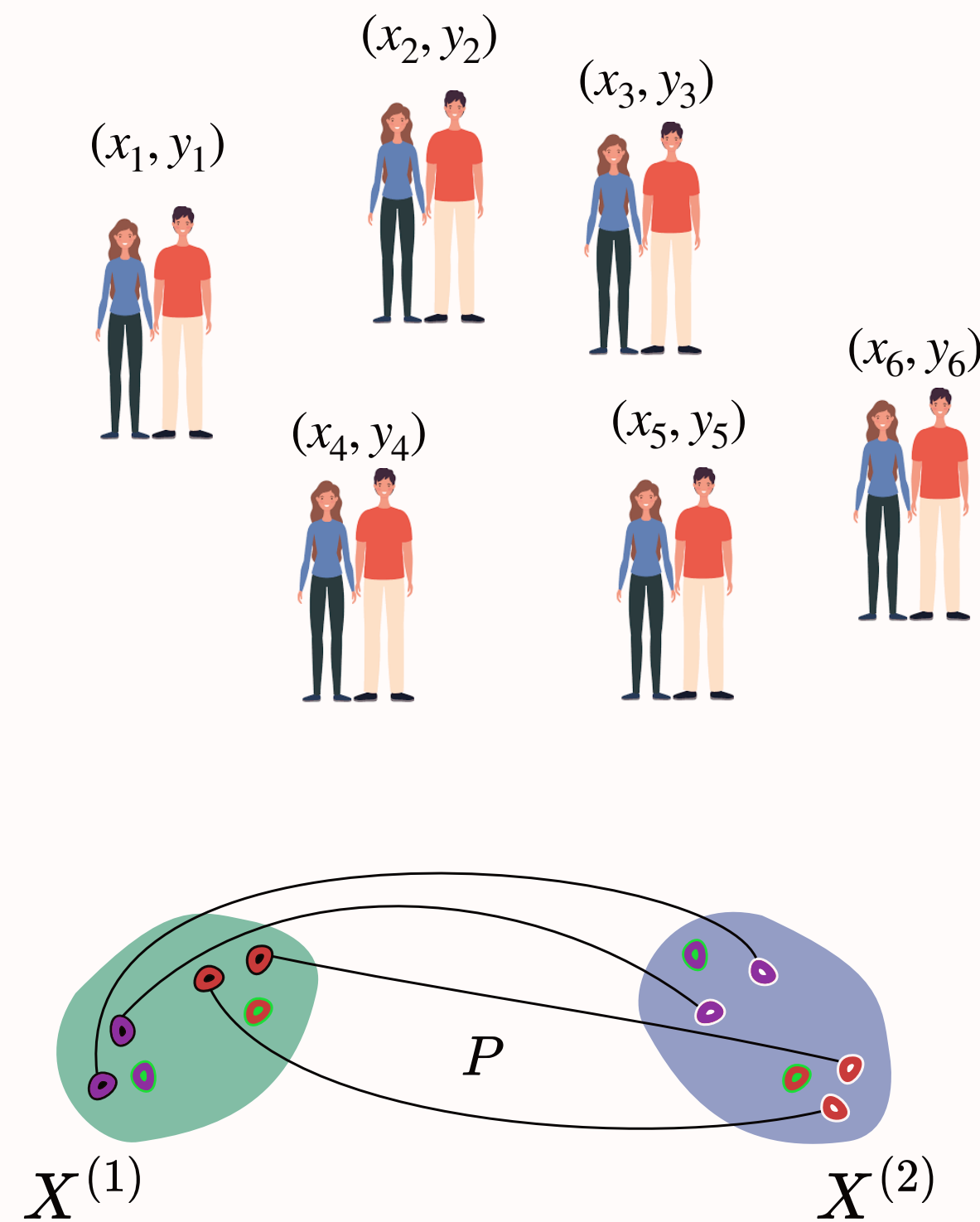
# Outline



Inverse problems in OT

$(x_1, y_1)$

$(x_2, y_2)$

$(x_3, y_3)$

$(x_4, y_4)$

$(x_5, y_5)$

$(x_6, y_6)$

$P$

$X^{(1)}$

$X^{(2)}$

Learning framework

$J(\theta)$

$\theta$

$\theta^\star$

$$\inf_\theta J(\theta) + \lambda R(\theta)$$

$c(x, y)$

$\hat{\pi}$

Recovery guarantees

$m_\alpha$

$m_\beta$

Number of wrongly estimated positions

$N=100$
$N=1400$
$N=2700$
$N=4000$

60
50
40
30
20
10
0

$\lambda$

2

# Optimal transport

$c(x, y)$ is the **cost** of moving from position $x$ to $y$
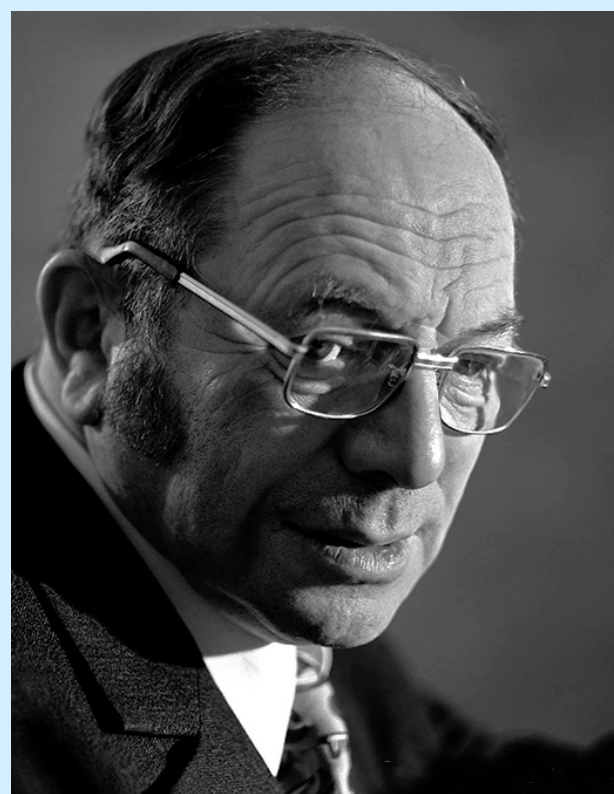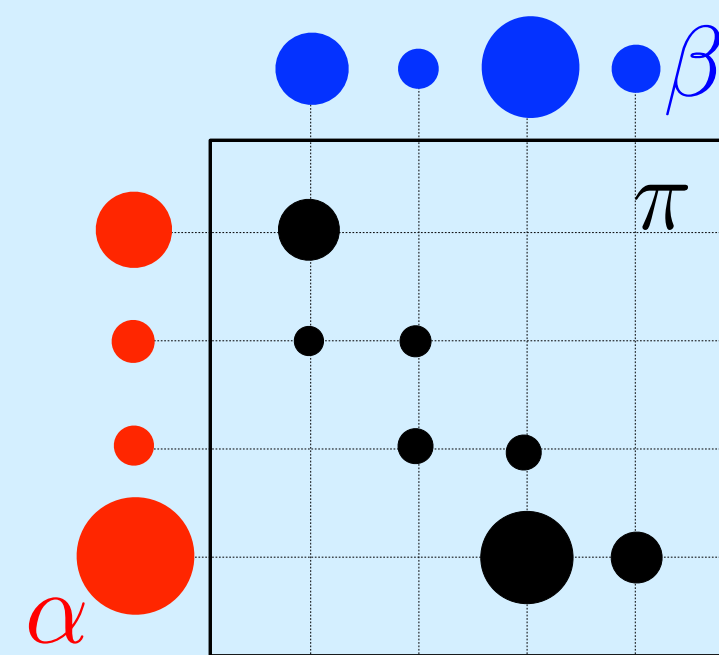
Monge 1781: Given probability measures $\alpha \in \mathscr{P}(\mathscr{X})$ and $\beta \in \mathscr{P}(\mathscr{Y})$, find the optimal way of transporting $\alpha$ to $\beta$.

$$\text{OT}(\alpha, \beta) = \inf_{T_\sharp \alpha = \beta} \int c(x, T(x)) d\alpha(x)$$

$$\text{OT}(\alpha, \beta) = \inf_{\pi_1 = \alpha, \pi_2 = \beta} \int c(x, y) d\pi(x, y)$$

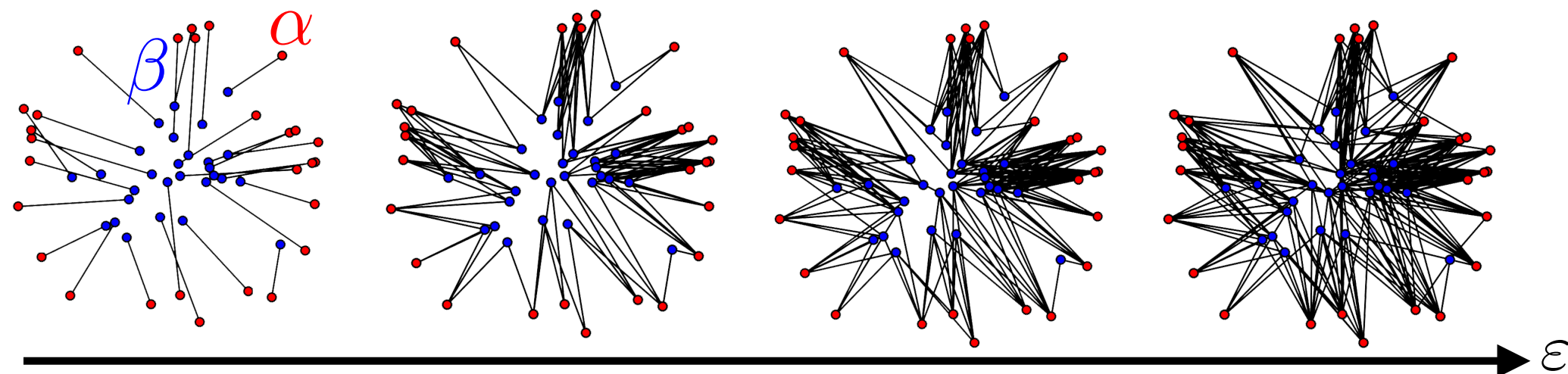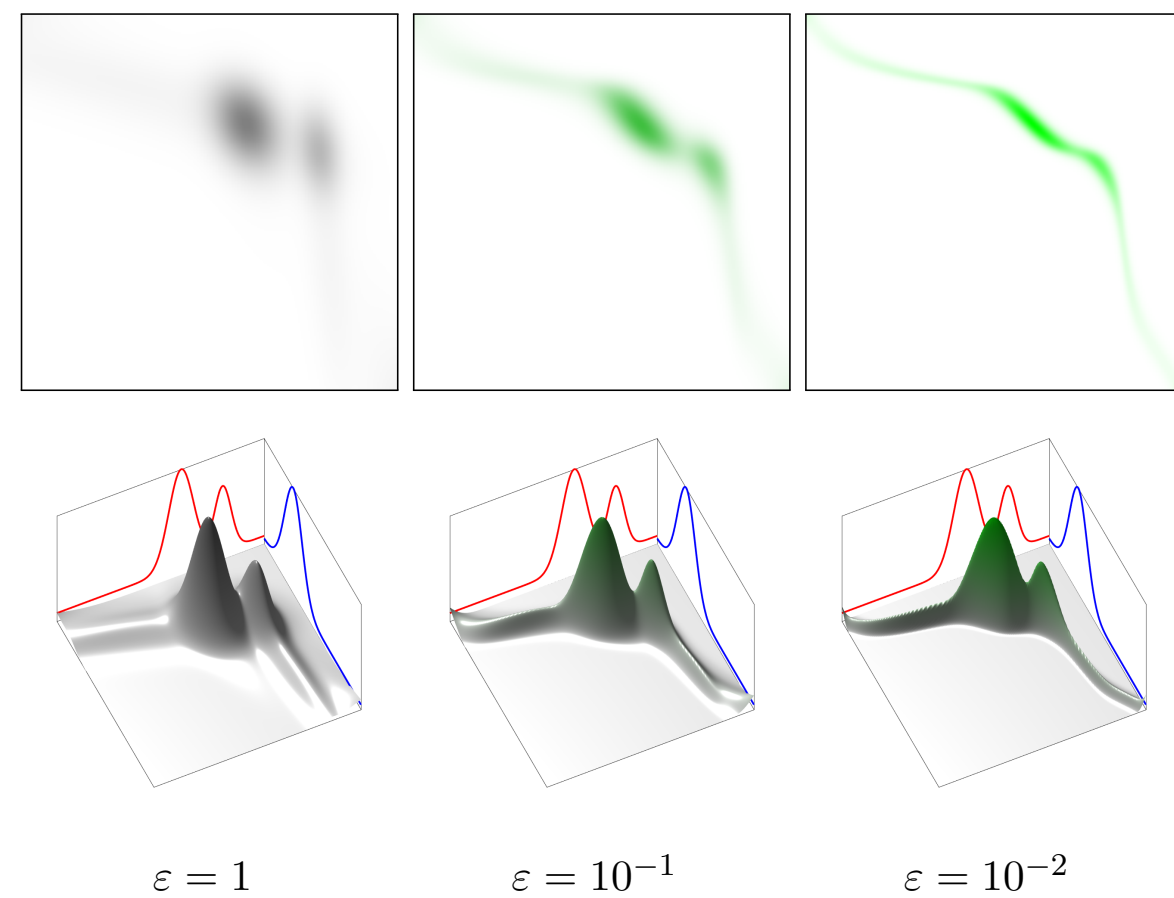$$\longleftrightarrow \quad \sup_{f \oplus g \leq c} \int f d\alpha + \int g d\beta$$

**Convex**

Kantorovich 1942

3

# Entropic optimal transport

**Regularize** with $\mathrm{KL}(\pi|\mu) := \int \log(d\pi/d\mu)d\pi$ and $\epsilon > 0$

$$\mathrm{eOT}(\alpha, \beta) = \inf_{\pi \in \mathcal{U}(\alpha,\beta)} \int c(x, y)d\pi(x, y) + \epsilon\mathrm{KL}(\pi\,|\,\alpha \otimes \beta)$$



$\varepsilon = 1 \qquad \varepsilon = 10^{-1} \qquad \varepsilon = 10^{-2}$



Images credit [Peyré & Cuturi '19]

# Entropic optimal transport

**Regularize** with $\mathrm{KL}(\pi|\mu) := \int \log(d\pi/d\mu)d\pi$ and $\epsilon > 0$

$$\mathrm{eOT}(\alpha, \beta) = \inf_{\pi \in \mathscr{U}(\alpha,\beta)} \int c(x, y)d\pi(x, y) + \epsilon\mathrm{KL}(\pi | \alpha \otimes \beta)$$

- **Natural** modelling assumption

- Alleviates the **curse of dimensionality** [Genevay et al '19, Mena & Weed '19]

$$\mathrm{eOT}(\alpha, \beta) - \mathrm{eOT}(\alpha_n, \beta_n) = \mathcal{O}(n^{-1/2})$$

- **Fast algorithms** available [Sinkhorn '64, Cuturi '13]

$$\sup_{f,g} \int f d\alpha + \int g d\beta - \epsilon \int \exp\left(\frac{f(x) + g(y) - c(x, y)}{\epsilon}\right) d\alpha(x)d\beta(y)$$

# Sinkhorn algorithm

- **Fast algorithms** available [Sinkhorn '64, Cuturi '13]

$$\sup_{f,g} \int f d\alpha + \int g d\beta - \epsilon \int \exp\left(\frac{f(x) + g(y) - c(x,y)}{\epsilon}\right) d\alpha(x) d\beta(y)$$

First order optimality:

$$\exp(-f/\epsilon) = \int \exp((g(y) - c(x,y))/\epsilon) d\beta(y)$$

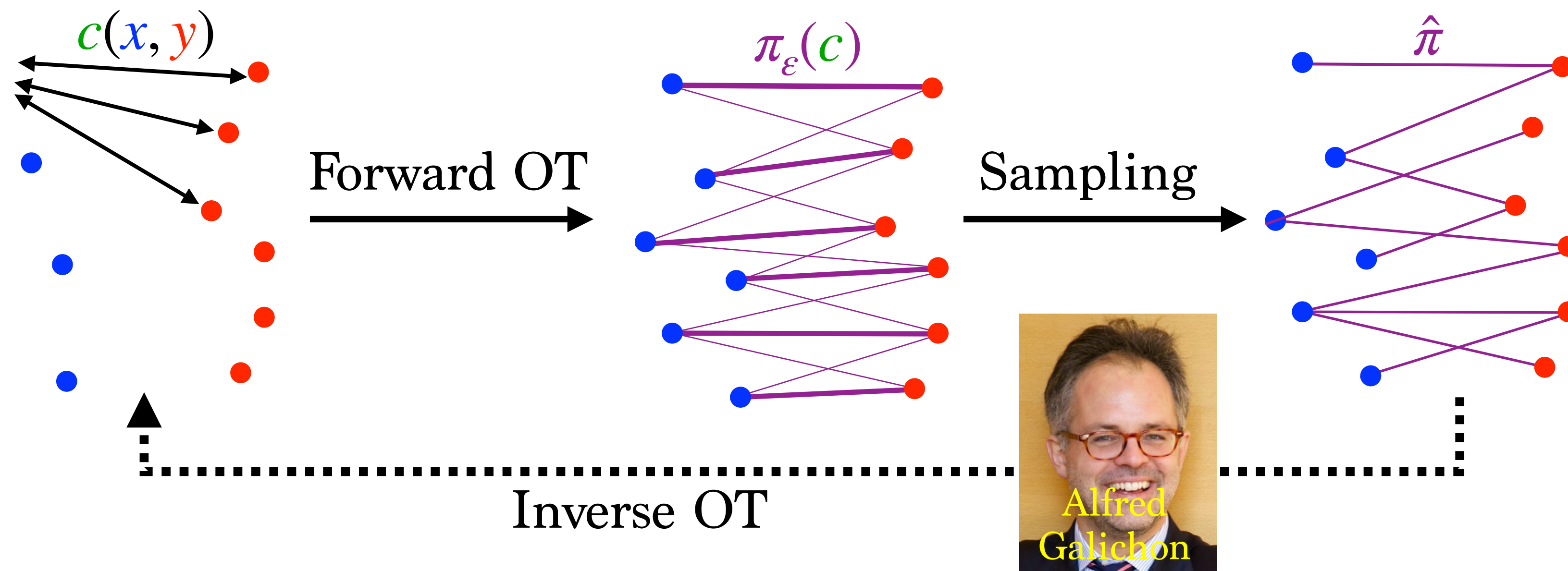$$\exp(-g/\epsilon) = \int \exp((f(x) - c(x,y))/\epsilon) d\alpha(y)$$

Sinkhorn is alternating minimisation:

$$f_{k+1} = -\epsilon \log\left(\int \exp((g_k(y) - c(x,y))/\epsilon) d\beta(y)\right)$$

$$g_{k+1} = -\epsilon \log\left(\int \exp((f_{k+1}(x) - c(x,y))/\epsilon) d\alpha(y)\right)$$

# Inverse optimal transport

*Given probability measures $\alpha, \beta$ and a ground cost $c(x, y)$, compute the optimal coupling.*



**Suppose you observe how two populations are coupled. How can we infer the 'cost' that led to this coupling?**
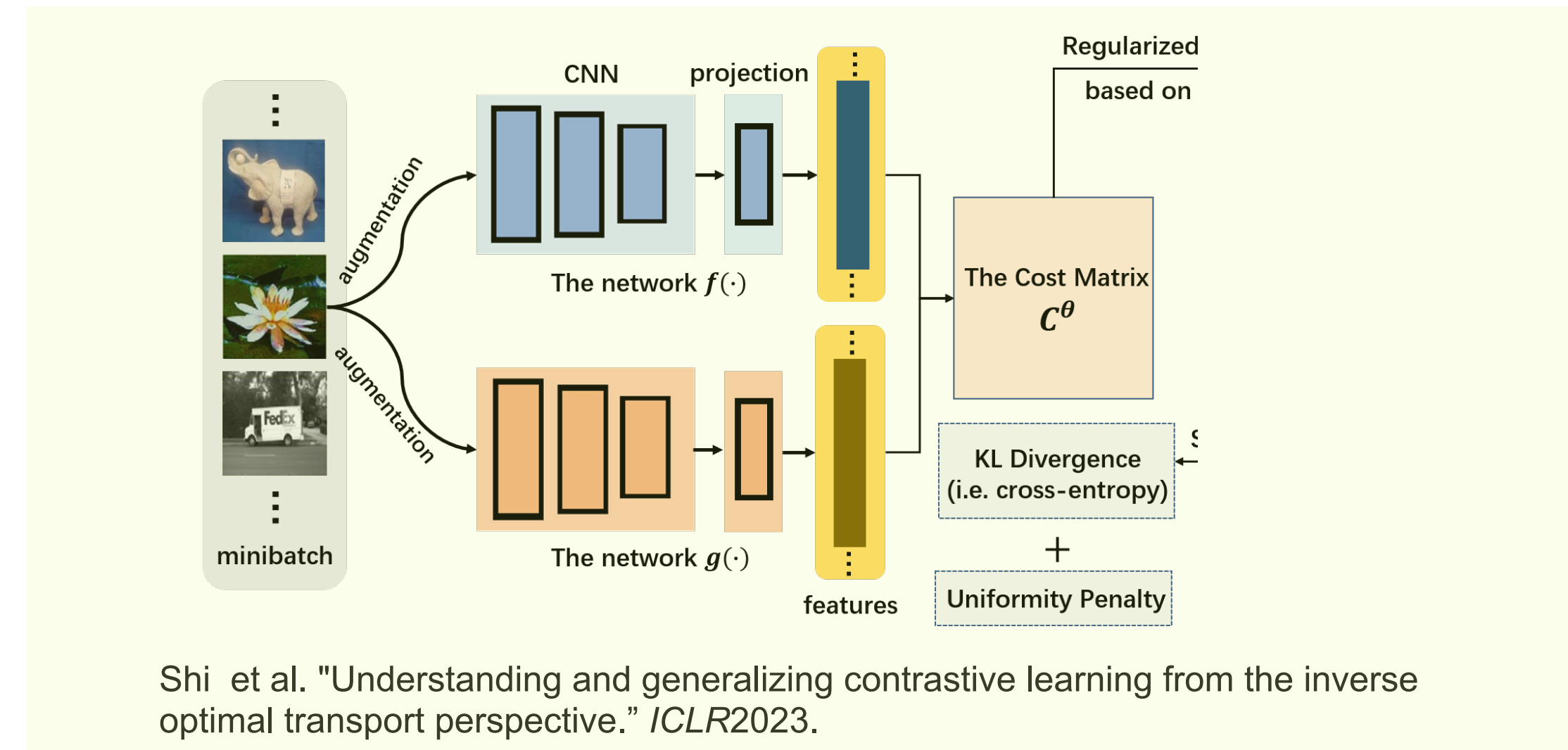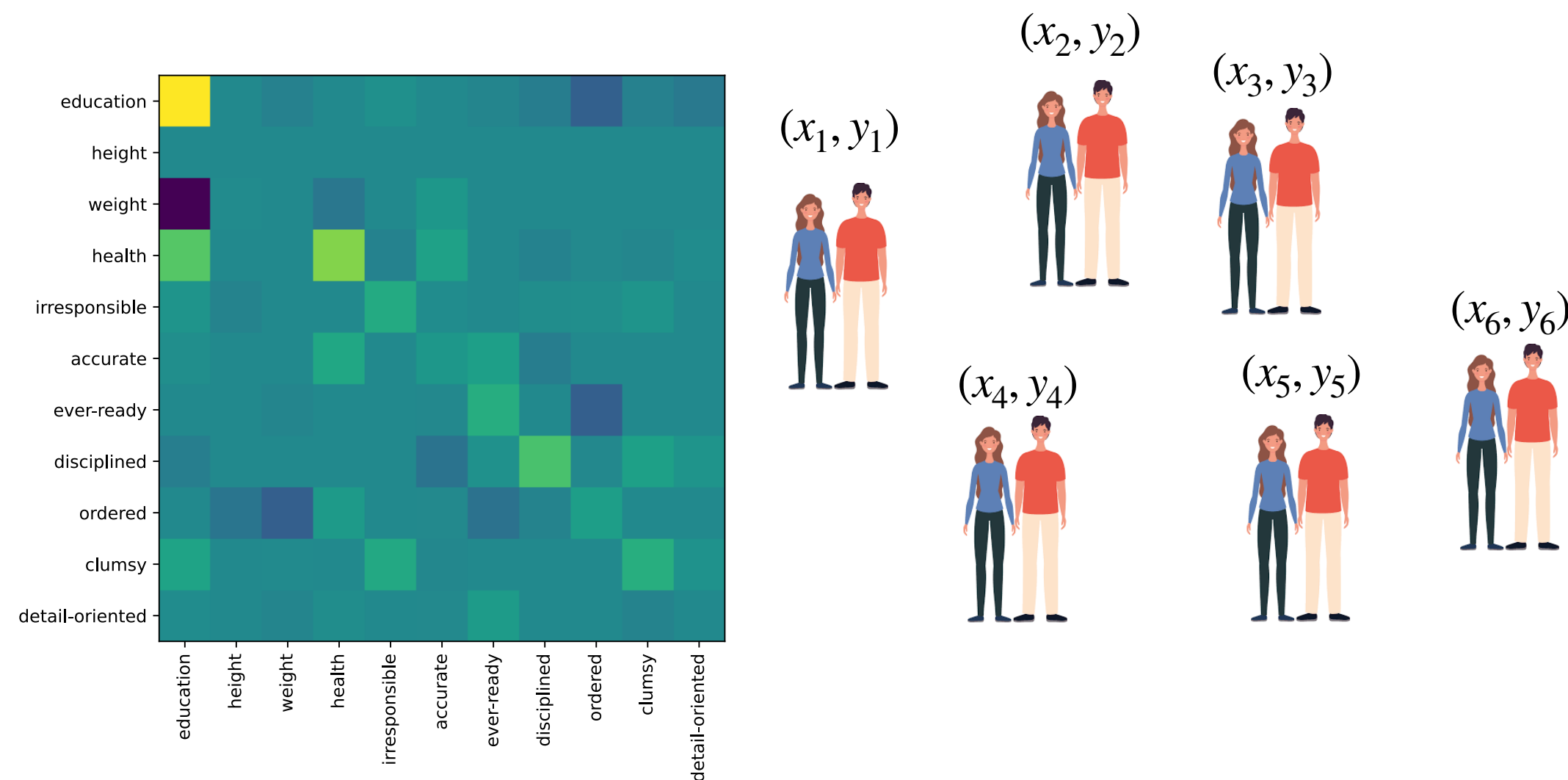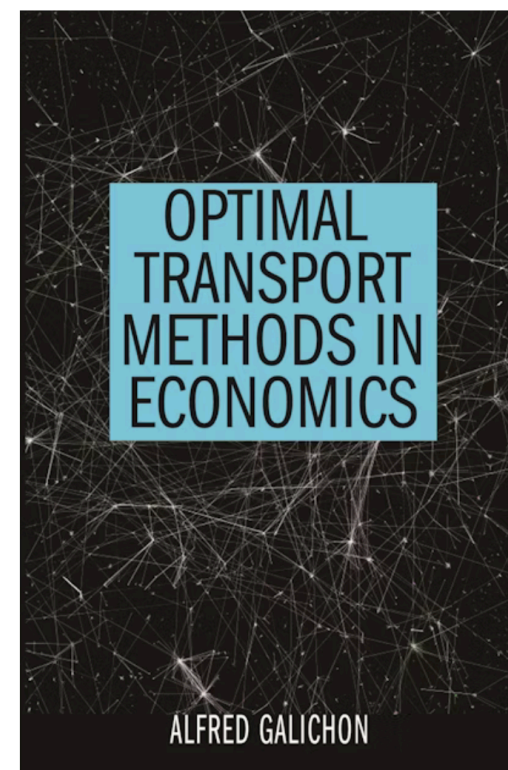
Galichon, Alfred, and Bernard Salanié. "Cupid's Invisible Hand: Social Surplus and Identification in Matching Models." (2015).
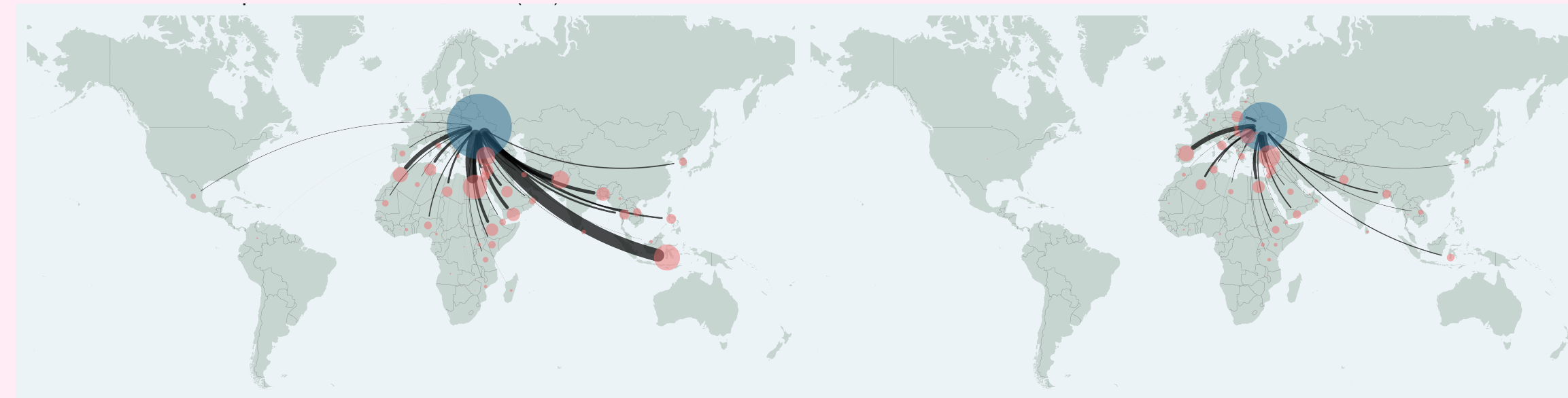Galichon, Alfred. *Optimal transport methods in economics*. Princeton University Press, 2016.
Dupuy, Arnaud, Alfred Galichon, and Yifei Sun. "Estimating matching affinity matrices under low-rank constraints." *Information and Inference: A Journal of the IMA* 8.4 (2019): 677-689.

# Understanding matching

Dupuy & Galichon 2014. "Personality traits and the marriage market."



$(x_1, y_1)$ $(x_2, y_2)$ $(x_3, y_3)$ $(x_4, y_4)$ $(x_5, y_5)$ $(x_6, y_6)$



Shi et al. "Understanding and generalizing contrastive learning from the inverse optimal transport perspective." *ICLR*2023.
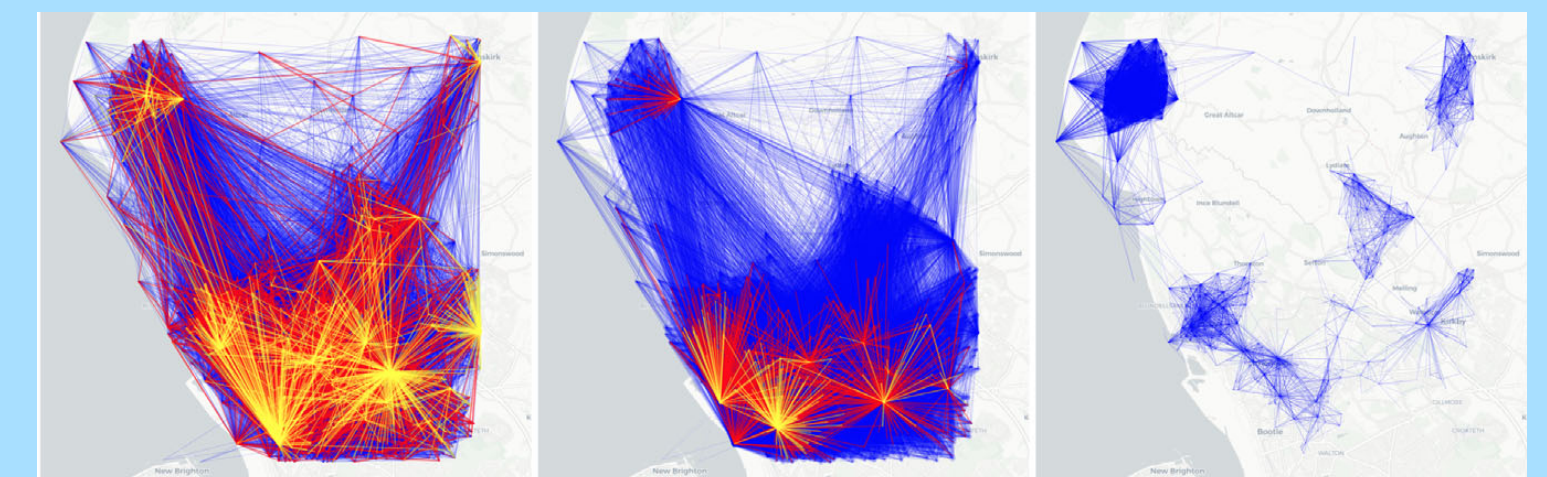


**Modelling Global Trade with Optimal Transport**

**Thomas Gaskin, Marie-Therese Wolfram, Andrew Duncan, Guven Demirel**

A Deep Gravity model for mobility flows generation

Filippo Simini [1,2,3], Gianni Barlacchi[4], Massimilano Luca [5,6] & Luca Pappalardo [7]

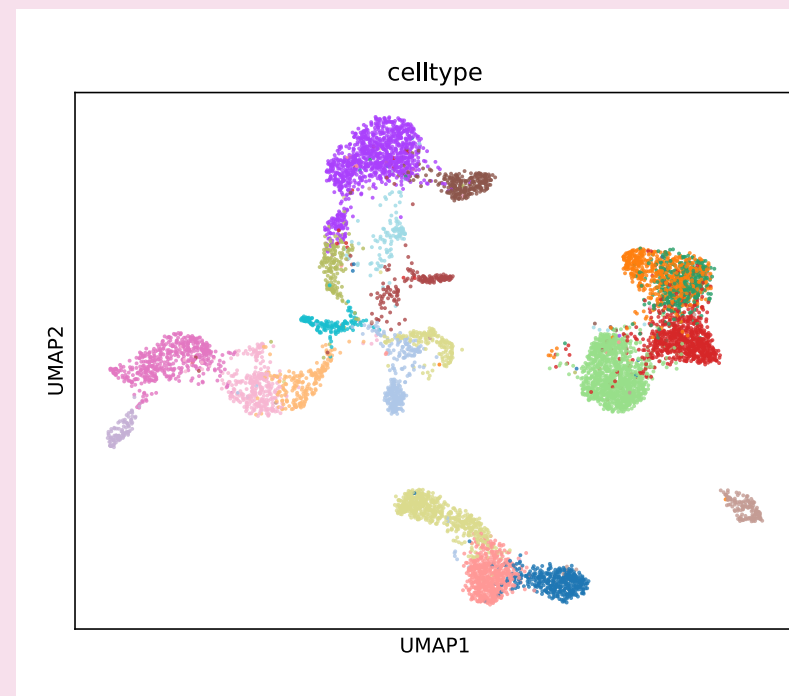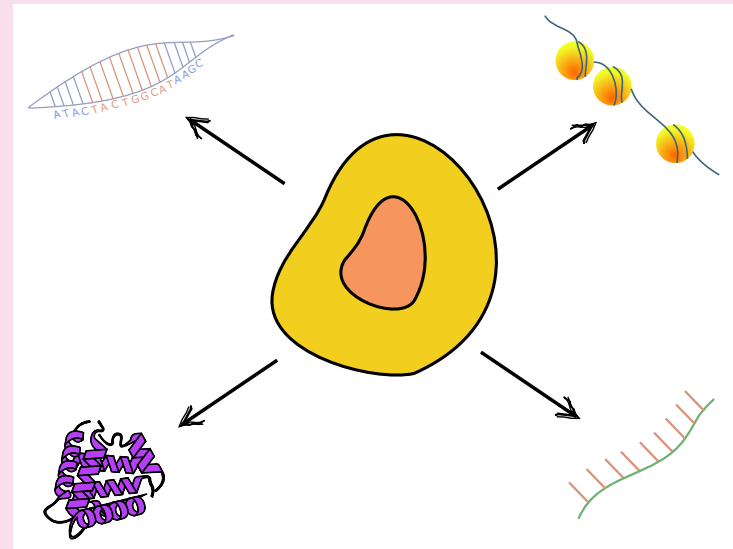a) Observed Flows   b) DG (CPC = 0.41)   c) G (CPC = 0.12)
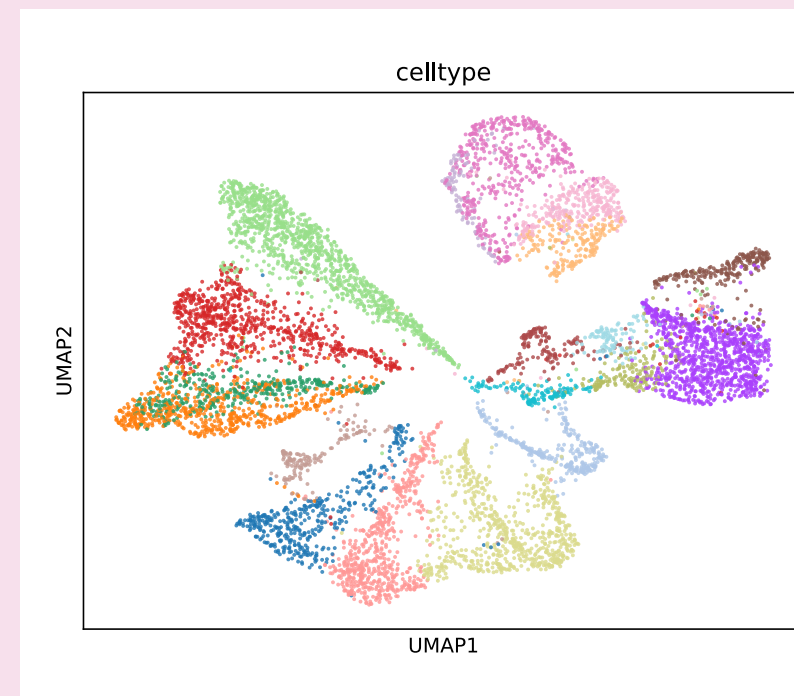
# Cell genomics

## Jules Samaran, Gabriel Peyré, Laura Cantini
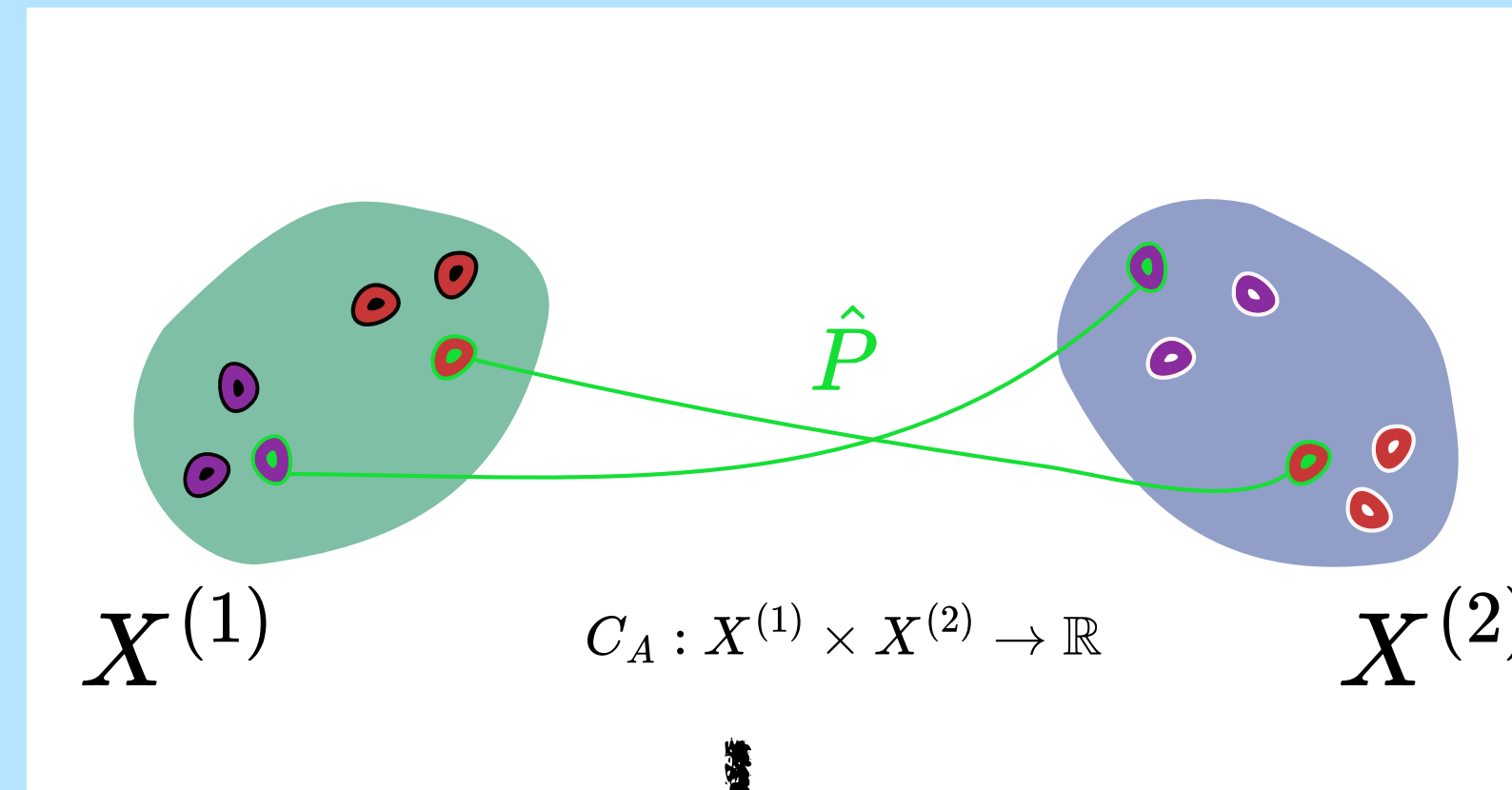


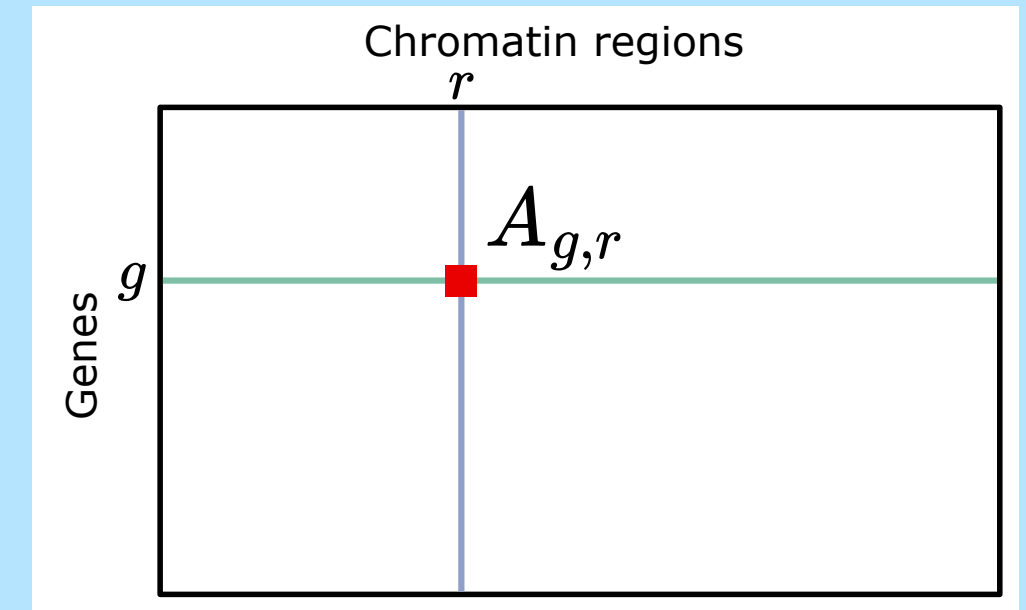**Cellular data is multimodal**

Gene expression

Chromatin accessibility

- B1 B
- CD4+ T activated
- CD4+ T naive
- CD8+ T
- CD14+ Mono
- CD16+ Mono
- Erythroblast
- G/M prog
- HSC
- Lymph prog
- MK/E prog
- NK
- Naive CD20+ B
- Normoblast
- Plasma cell
- Proerythroblast
- Transitional B
- cDC2
- pDC

$X^{(1)}$  $\quad C_A : X^{(1)} \times X^{(2)} \to \mathbb{R} \quad$  $X^{(2)}$

$\hat{P}$

Jules Samaran

$C_A(x_1, x_2) = x_1^T A x_2$

Chromatin regions

$A_{g,r}$

Genes

***Learn the missing links***

$X^{(1)}$  $\quad P \quad$  $X^{(2)}$

89% accuracy

9

# How do populations evolve?

$$\frac{dX_i}{dt} = \mathbf{v}_t(X_i) \qquad X_i \overset{iid}{\sim} \rho_t \qquad \partial_t \rho_t + \nabla \cdot (\rho_t \mathbf{v}_t) = 0$$

*Goal: Given iid samples of $\rho_{t_0}, \rho_{t_1}, \rho_{t_2}, \ldots \rho_{t_T}$, find $\mathbf{v}_t$.*

Examples:

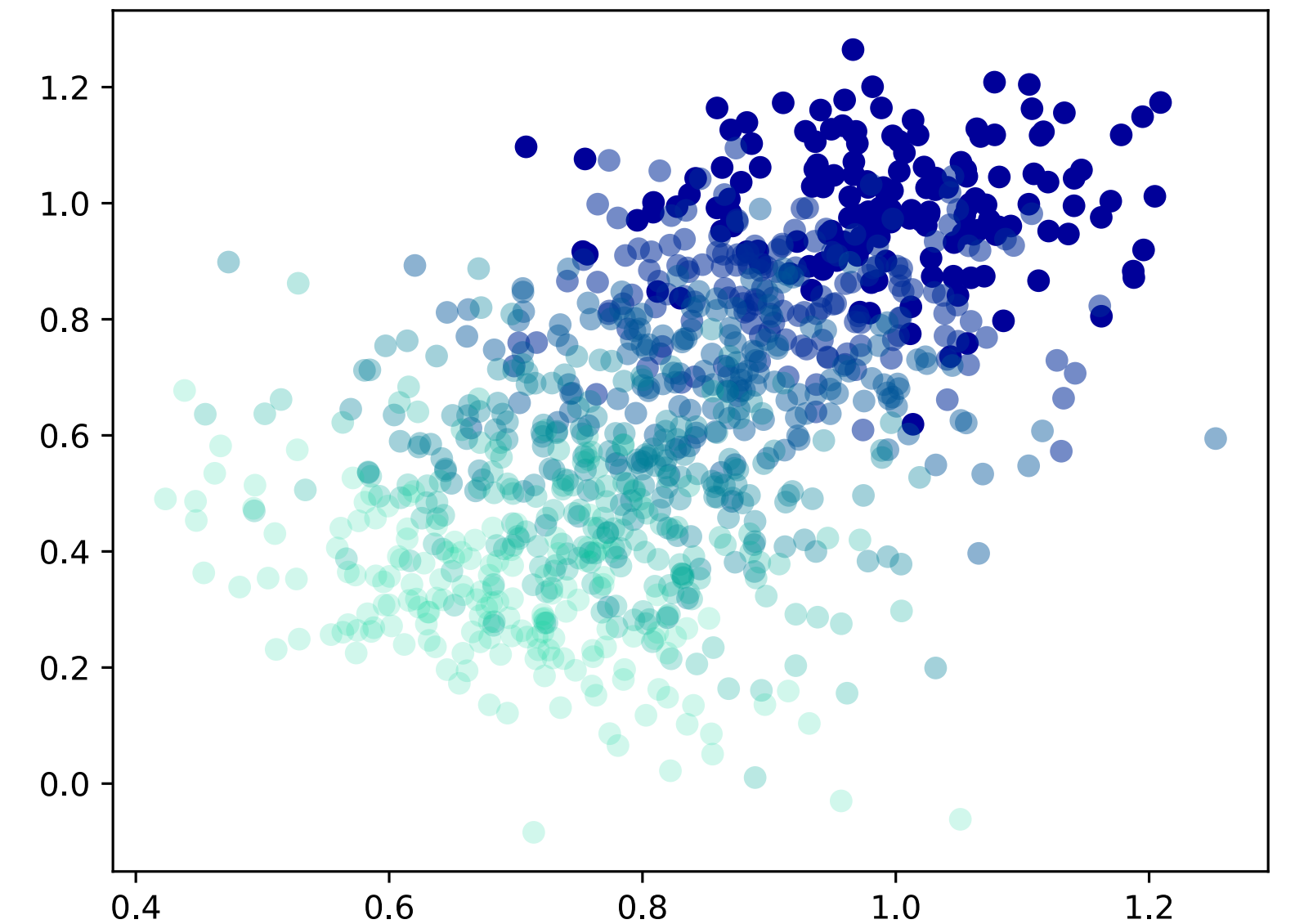$$\mathbf{v} = \nabla V \text{ or } \mathbf{v} = \int \nabla_1 W(\cdot, x) d\rho(x)$$



**Modelling using Wasserstein gradient flows:**

Suppose $\mathbf{v}_t = -\nabla \delta \mathscr{F}(\rho_t)$ for some $\mathscr{F} : \mathscr{P}(X) \to \mathbb{R}$

Model using discretised dynamics: $\alpha_{k+1} = \underset{\alpha \in \mathscr{P}(X)}{\text{argmin}} \, \mathscr{F}(\alpha) + \frac{1}{2\tau} W_2^2(\alpha, \alpha_k)$

Jordon Kinderlehrer Otto scheme (JKO)
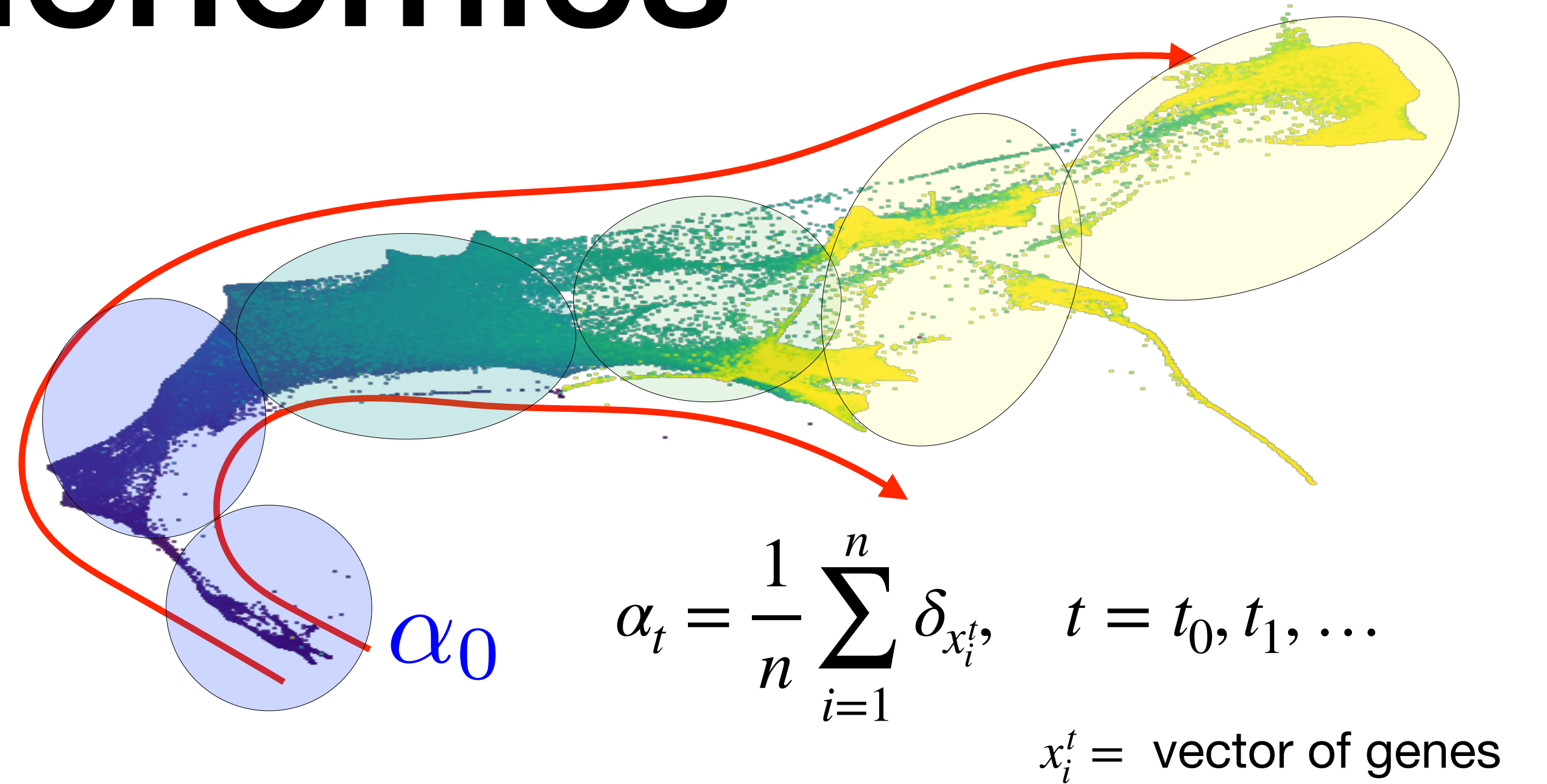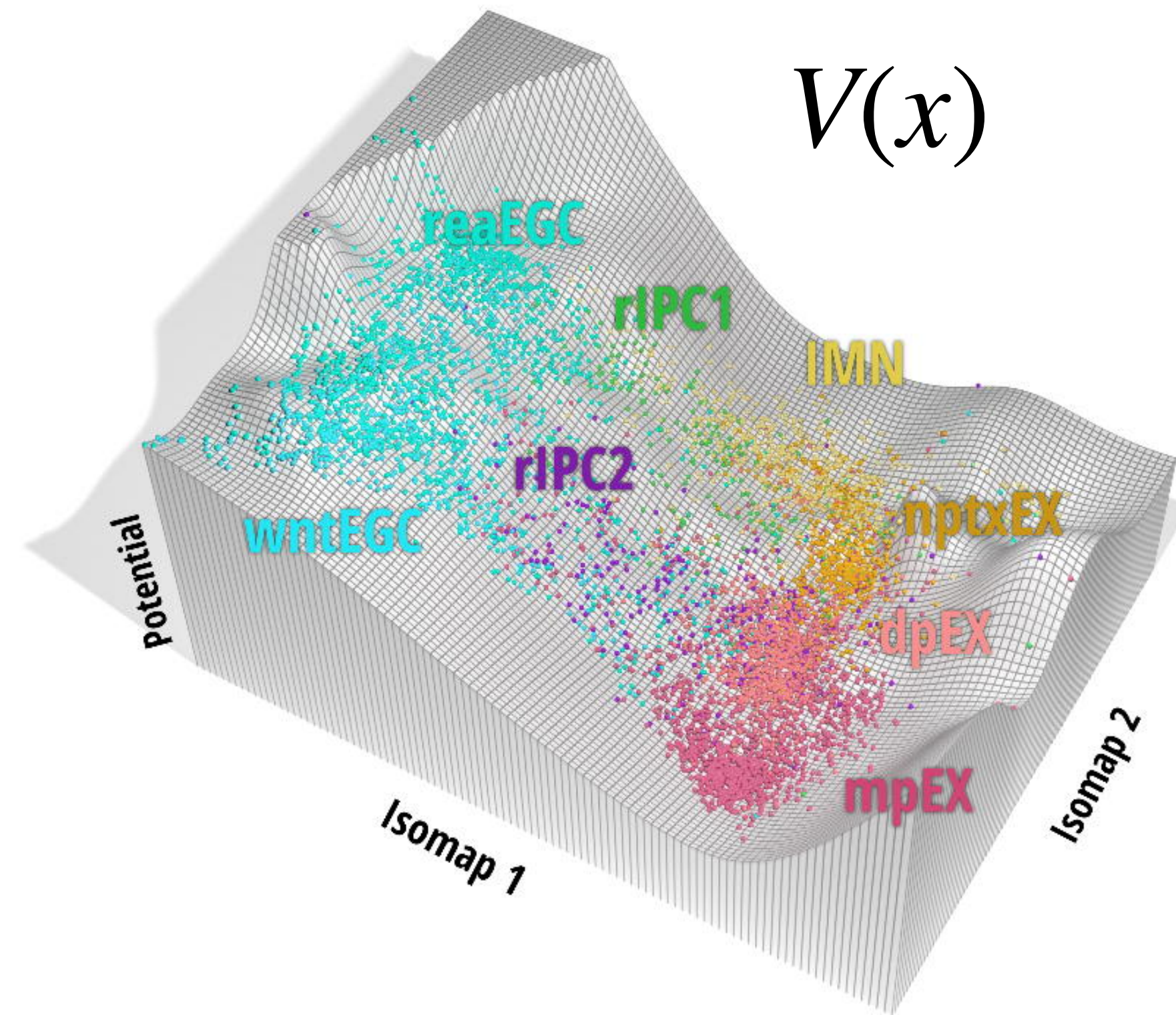
Discretisation of curve of probability measures $\alpha_k \approx \rho_{k\tau}$

Example: $\mathscr{F}(\rho) = \int V(x) d\rho(x)$ so $\delta \mathscr{F} = V$ and $\mathbf{v} = \nabla V$

Example: $\mathscr{F}(\rho) = \frac{1}{2} \int W(x, y) d\rho(x) d\rho(y), \, \delta \mathscr{F}(\rho) = \int W(\cdot, x) d\rho(x)$ and $\mathbf{v} = \int \nabla_1 W(\cdot, x) d\rho(x)$

# Cell Genomics

$V(x)$



$$\alpha_t = \frac{1}{n} \sum_{i=1}^{n} \delta_{x_i^t}, \quad t = t_0, t_1, \dots$$
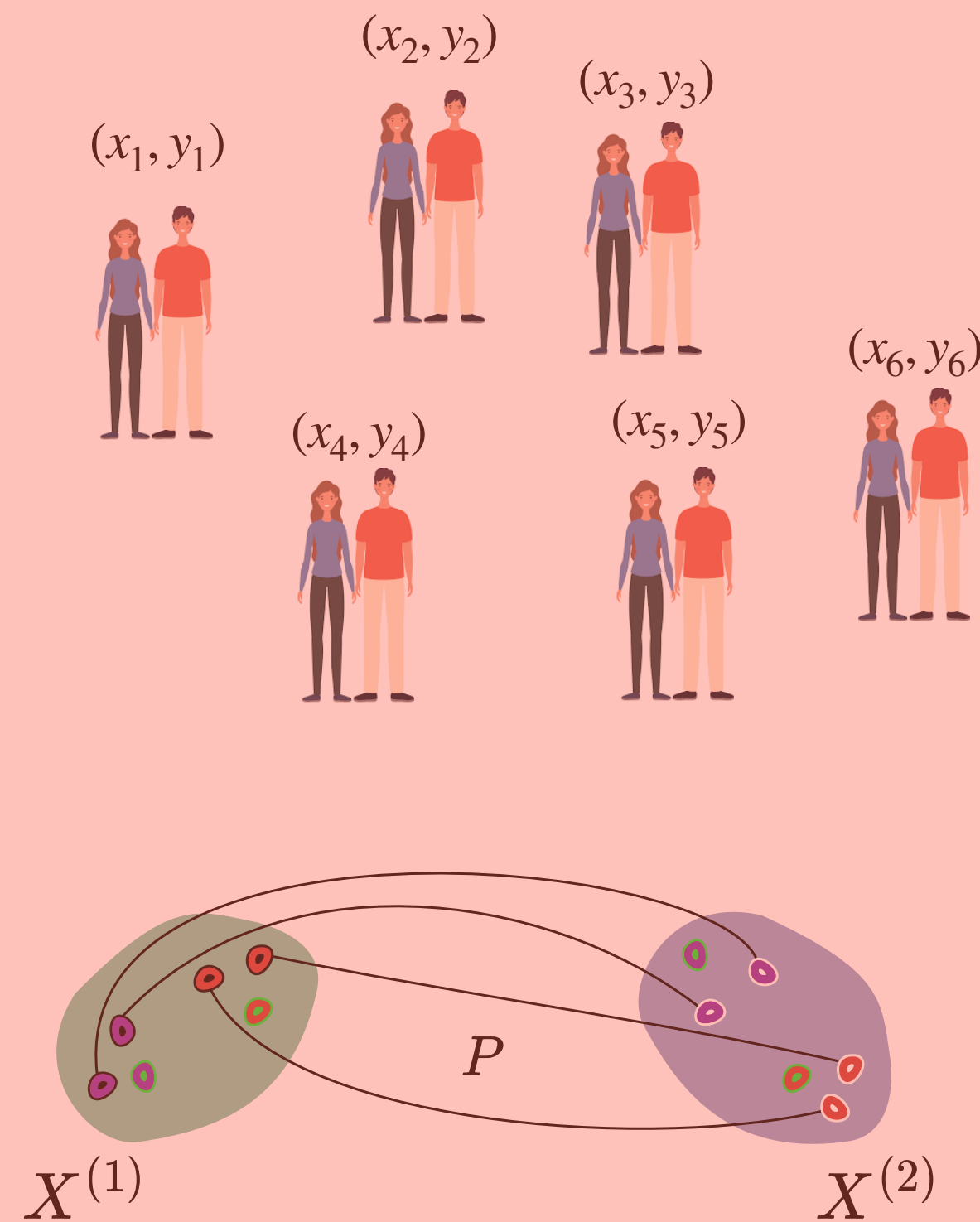
$x_i^t =$ vector of genes

*Model as JKO flow*

$$\alpha_{t+1} = \text{argmin}_\rho \int V(x) d\rho(x) + W_2^2(\rho, \alpha_t)$$

*Reference: Learning cell fate landscapes from spatial transcriptomics using Fused Gromov-Wasserstein.*
***Geert-Jan Huizing, Gabriel Peyré, Laura Cantini***

# Outline



Inverse problems in OT

$(x_1, y_1)$ $(x_2, y_2)$ $(x_3, y_3)$ $(x_4, y_4)$ $(x_5, y_5)$ $(x_6, y_6)$

$P$

$X^{(1)}$ $X^{(2)}$

Learning framework

$J(\theta)$

$\theta^\star$ $\theta$

$\inf_{\theta} J(\theta) + \lambda R(\theta)$

$c(x, y)$ $\hat{\pi}$

Recovery guarantees

$m_\alpha$ $m_\beta$

Number of wrongly estimated positions

$N=100$
$N=1400$
$N=2700$
$N=4000$

60
50
40
30
20
10
0

$\lambda$

# Outline



Inverse problems in OT

$(x_1, y_1)$   $(x_2, y_2)$   $(x_3, y_3)$

$(x_4, y_4)$   $(x_5, y_5)$   $(x_6, y_6)$

$P$

$X^{(1)}$     $X^{(2)}$

Learning framework

$J(\theta)$

$\theta$

$\theta^\star$

$\inf_\theta J(\theta) + \lambda R(\theta)$

$c(x, y)$     $\hat{\pi}$

Recovery guarantees

$m_\alpha$     $m_\beta$

*Number of wrongly estimated positions*

$N=100$
$N=1400$
$N=2700$
$N=4000$

$\lambda$

13

# Stability guarantee

Suppose that $\pi^\star = P_\Omega(c_\theta^\star) := \mathrm{argmin}_c \langle c_{\theta^\star}, \pi \rangle + \Omega(\pi)$. Assume that $\hat{\pi}$ is noisy version of $\pi^\star$. Solve:

$$\inf_{\theta \in \mathbb{R}^p} \mathscr{L}(c_\theta, \textcolor{red}{\hat{\pi}, \hat{\Omega}}) + \lambda R(\theta)$$

***Theorem:*** Let $\gamma > 0$ be the 'noise' level. If there is:

◆ Measurement stability, $|\langle \psi, \hat{\pi} - \pi^\star \rangle| \leq \gamma$ for all basis elements $\psi$

◆ Forward stability: $|\langle \psi, P_\Omega(c_\theta^\star) - P_{\hat{\Omega}}(c_\theta^\star) \rangle| \leq \gamma$ for all basis elements $\psi$

◆ Local curvature: $J(\theta) := \mathscr{L}(c_\theta | \hat{\pi}, \hat{\Omega})$ is locally strongly convex and Lipschitz smooth.

Then, the minimizer $\theta$ satisfies $\|\theta - \theta^\star\| = \mathcal{O}(\lambda + \gamma)$.

# Assumptions

Loss to minimize in the case of **iOT**:

$$\inf_{f,g,c} \langle c - (f \oplus g), \hat{\pi} \rangle + \epsilon \int \exp\left( \frac{f(x) + g(y) - c(x, y)}{\epsilon} \right) d\hat{\alpha}(x) d\hat{\beta}(y)$$

<u>Non-uniqueness for iOT</u>
1. If $f, g$ minimizers $\iff f + a, g - a$ are minimisers for all constants $a$.
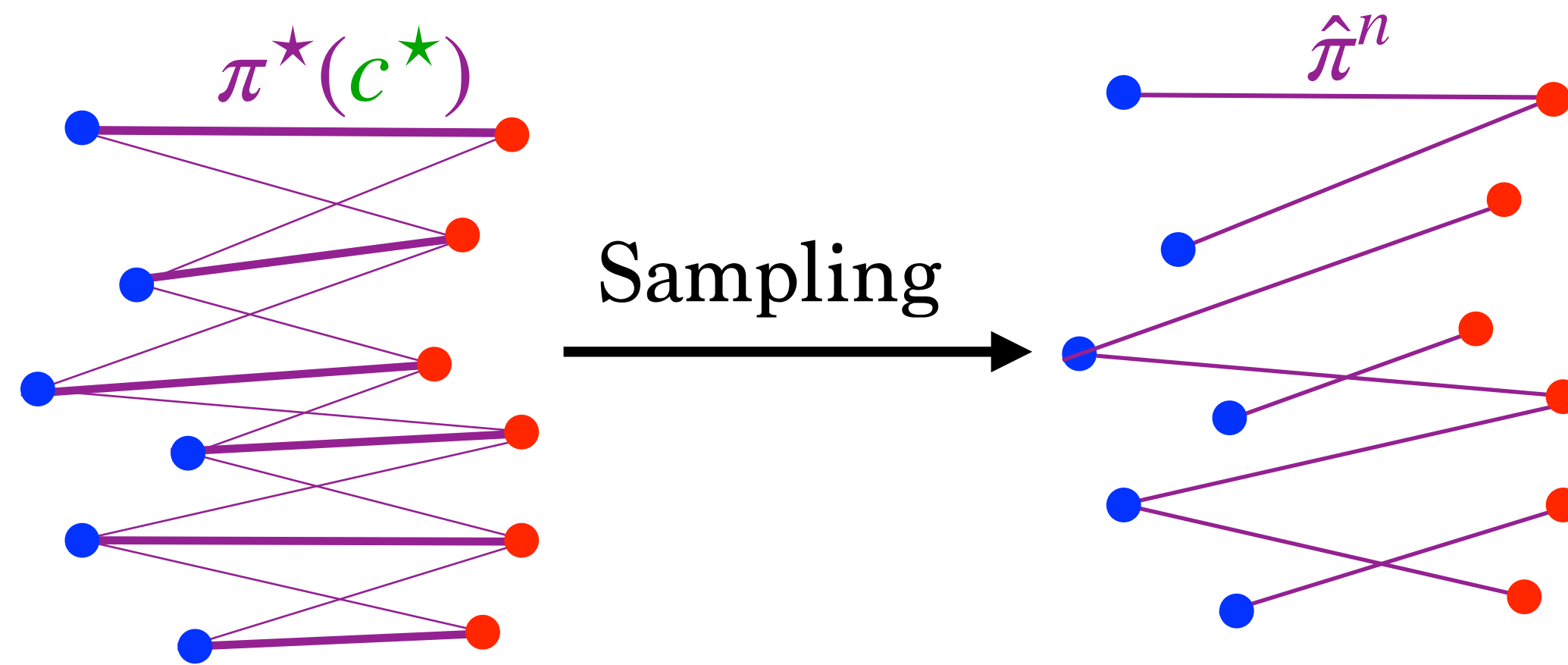2. Replace $c$ with $c - u \oplus v \iff f + u, g + v$ are minimisers.

Assume $c_\theta(x, y) = \sum_k \theta_k c_k(x, y)$ are <u>centred</u>:

$$\int c_k(x, y) d\alpha(x) = 0, \quad \int c_k(x, y) d\beta(y) = 0$$

and linearly independent.

Assume that $\alpha, \beta$ have compact supports.

15

# Sample complexity for iUOT



$\pi^\star(c^\star)$

$\hat{\pi}^n$

Sampling

*How accurate is $c_{\theta^{n,\lambda}}$ constructed from $\hat{\pi}^n$ ?*

**Theorem**:

Let $c^\star = \Phi\theta^\star$ be the cost that gave rise to $\pi^\star, \alpha^\star, \beta^\star$. Let $R$ be a convex, l.s.c. regularizer. Given $(x_i, y_i) \sim \pi^\star$ iid with $i = 1,\ldots,n$, the solution to $\theta = \underset{\theta}{\arg\min}\, L(\theta; \hat{\pi}_n, \hat{\alpha}_n, \hat{\beta}_n) + \lambda R(\theta)$ is <span style="color:red">unique</span> and with probability at least $1 - e^{-t}$,
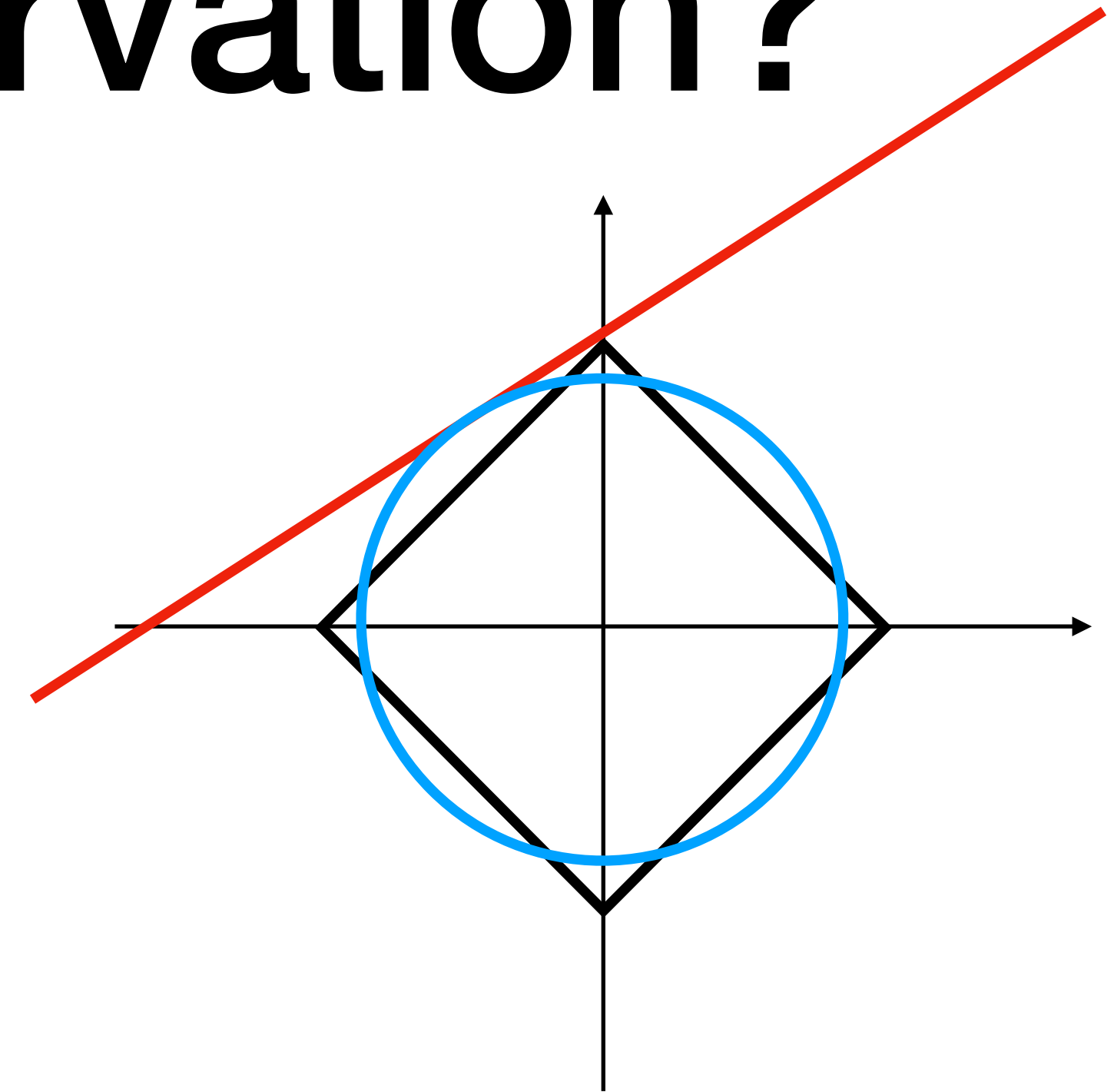
$$\|\theta - \theta^\star\| = \mathcal{O}\left(\lambda + \frac{\sqrt{\log(n) + t}}{\sqrt{n}}\right)$$

# Structure preservation?

**Typical regularisers:**

Sparsity.  $R(\theta) = \sum_{i=1}^{n} |\theta_i|$

Low rank  $R(\theta) = \sum_{i=1}^{n} \sigma_i(\theta)$ where $\sigma_i(\theta)$ = singular values of $\theta$.

$$\hat{\theta} = \underset{\theta}{\text{argmin}}\, L(\theta; \hat{\pi}_n, \hat{\alpha}_n, \hat{\beta}_n) + \lambda R(\theta)$$
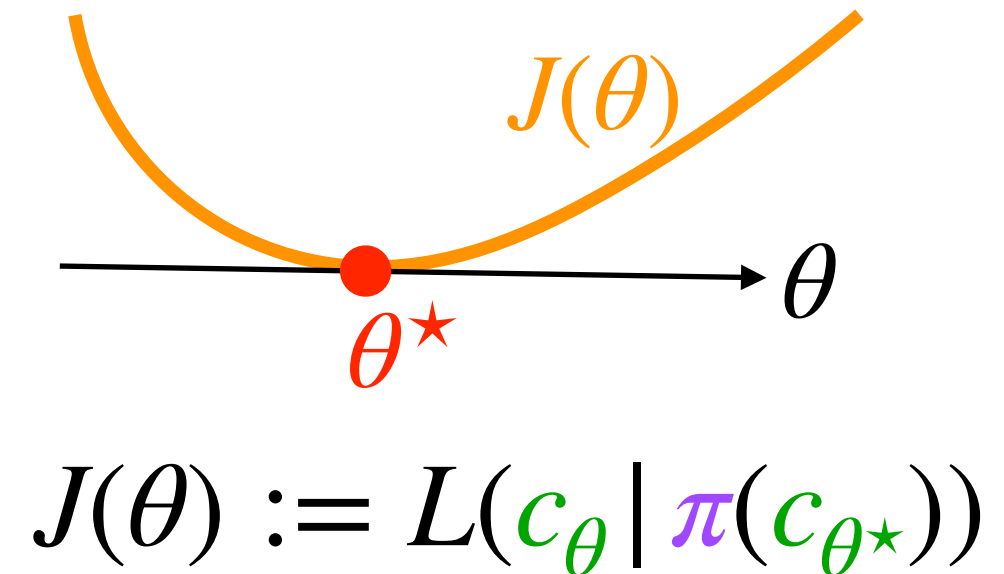
**Question:**

*If $\theta^\star$ that generated $\pi^\star, \alpha^\star, \beta^\star$ is sparse/low rank, is the solution $\hat{\theta}$ also of the same sparsity/rank when $n$ is large enough and $\lambda$ is small enough?*

# Structure preservation

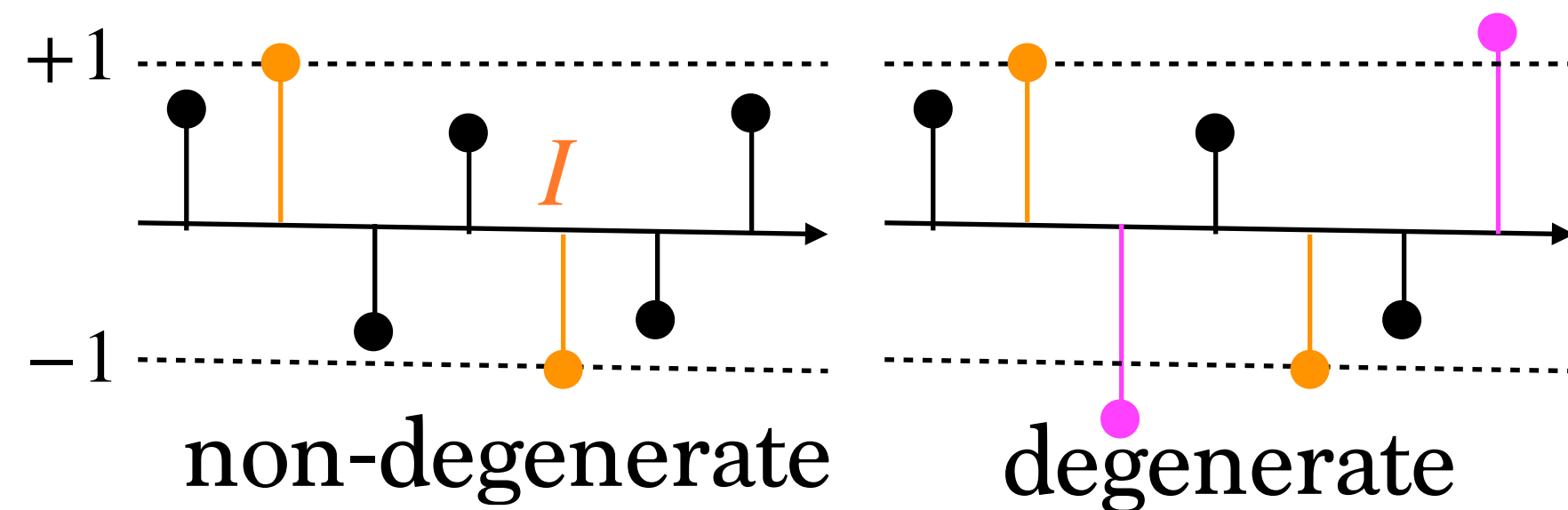**Certificate:** Let $M := \nabla^2 J(\theta^\star)$, and define

○ $\hat{z}_1 = M_{(:,I)} M_{(I,I)}^{-1} \text{Sign}(\theta_I^\star)$ in the case of l1 where $I$ is the support of $\theta^\star$

○ $\hat{z}_* = M P_T (P_T M P_T)^{-1} \text{Sign}(\theta_I^\star)$ where $P_T$ is the projection onto the row/column space spanned by $\theta^\star$.



$$J(\theta) := L(c_\theta \mid \pi(c_{\theta^\star}))$$

We say that

$\hat{z}_1$ is *non degenerate* if $\|\hat{z}_{I^c}^1\|_\infty < 1$

$\hat{z}*$ is non degenerate if $\|P_T^\perp \hat{z}*\|_2 < 1$.
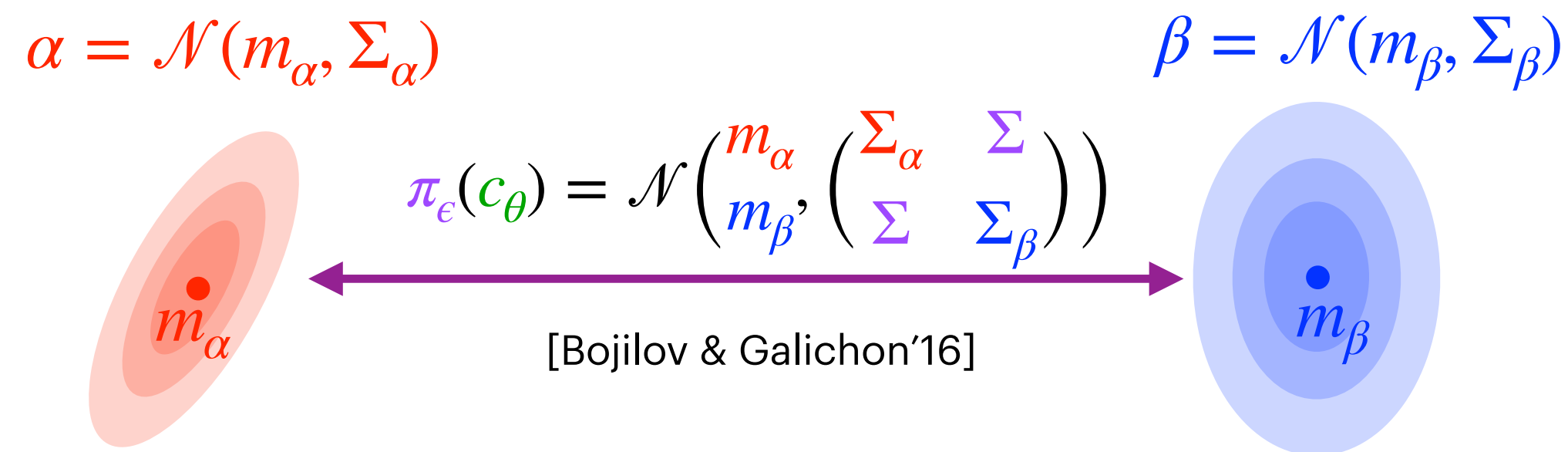


non-degenerate        degenerate

**Theorem**: Suppose that $\hat{z}$ is nondegenerate.

Then, provided that $1 \gtrsim \lambda \gtrsim \sqrt{\dfrac{t + \log(n)}{n}}$, with

probability at least $1 - e^{-t}$,

○ $\text{Supp}(\theta_{n,\lambda}) = \text{Supp}(\theta^\star)$ in the case of $R(\theta) = \|\theta\|_1$

○ $\text{Rank}(\theta_{n,\lambda}) = \text{Rank}(\theta^\star)$ in the case of $R(\theta) = \|\theta\|_*$.

# The iOT loss in the Gaussian setting

$\alpha = \mathcal{N}(m_\alpha, \Sigma_\alpha)$

$\beta = \mathcal{N}(m_\beta, \Sigma_\beta)$

$$\pi_\varepsilon(c_\theta) = \mathcal{N}\left(\begin{matrix} m_\alpha \\ m_\beta \end{matrix}, \begin{pmatrix} \Sigma_\alpha & \Sigma \\ \Sigma & \Sigma_\beta \end{pmatrix}\right)$$

$m_\alpha$

$m_\beta$

[Bojilov & Galichon'16]

*Proposition:*

$$\partial^2 J(\theta^\star) = 2\varepsilon\left[4\varepsilon^2(\Sigma_\beta - \Sigma^T\Sigma_\alpha\Sigma)^{-1} \otimes (\Sigma_\alpha - \Sigma\Sigma_\beta^{-1}\Sigma^\top)^{-1} + (\theta^{\star\top} \otimes \theta^\star)\right]^{-1}$$

$\rightarrow$ Numerically check when $\eta^\star$ is non-degenerated.

$\lambda = \lambda_0\varepsilon$
$\varepsilon \to 0$

$$\min_\theta L(c_\theta \mid \hat{\pi}) + \lambda\|\theta\|_1$$

$\lambda = \lambda_0/\varepsilon$
$\varepsilon \to +\infty$

$$\min_{\theta > 0} \frac{1}{2}\log\det(\theta) + \frac{1}{2}\langle\theta, \hat{\theta}^{-1}\rangle + \lambda_0\|\theta\|_1$$

Graphical-Lasso

$$\min_\theta \frac{1}{2}\|(\Sigma_\beta^{\frac{1}{2}} \otimes \Sigma_\alpha^{\frac{1}{2}})(\theta - \hat{\theta})\|_F^2 + \lambda_0\|\theta\|_1$$
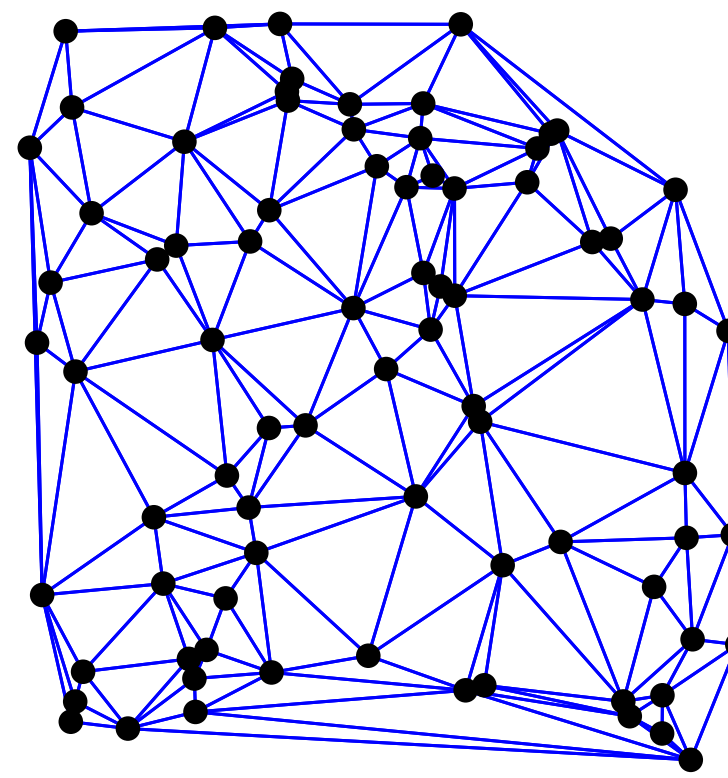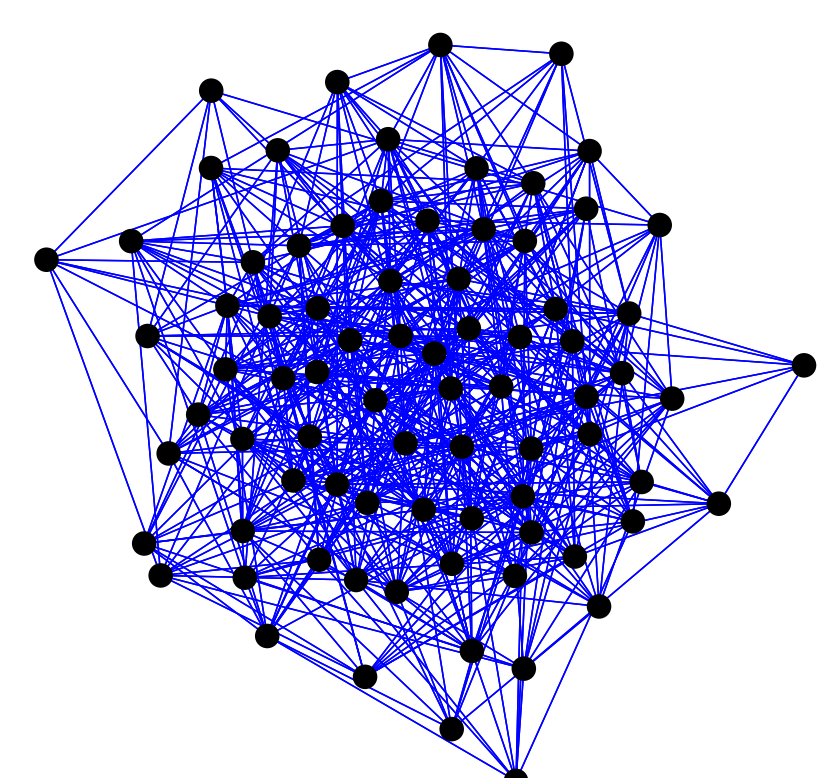
Lasso

# Numerical illustrations of $\ell_1$ certificates
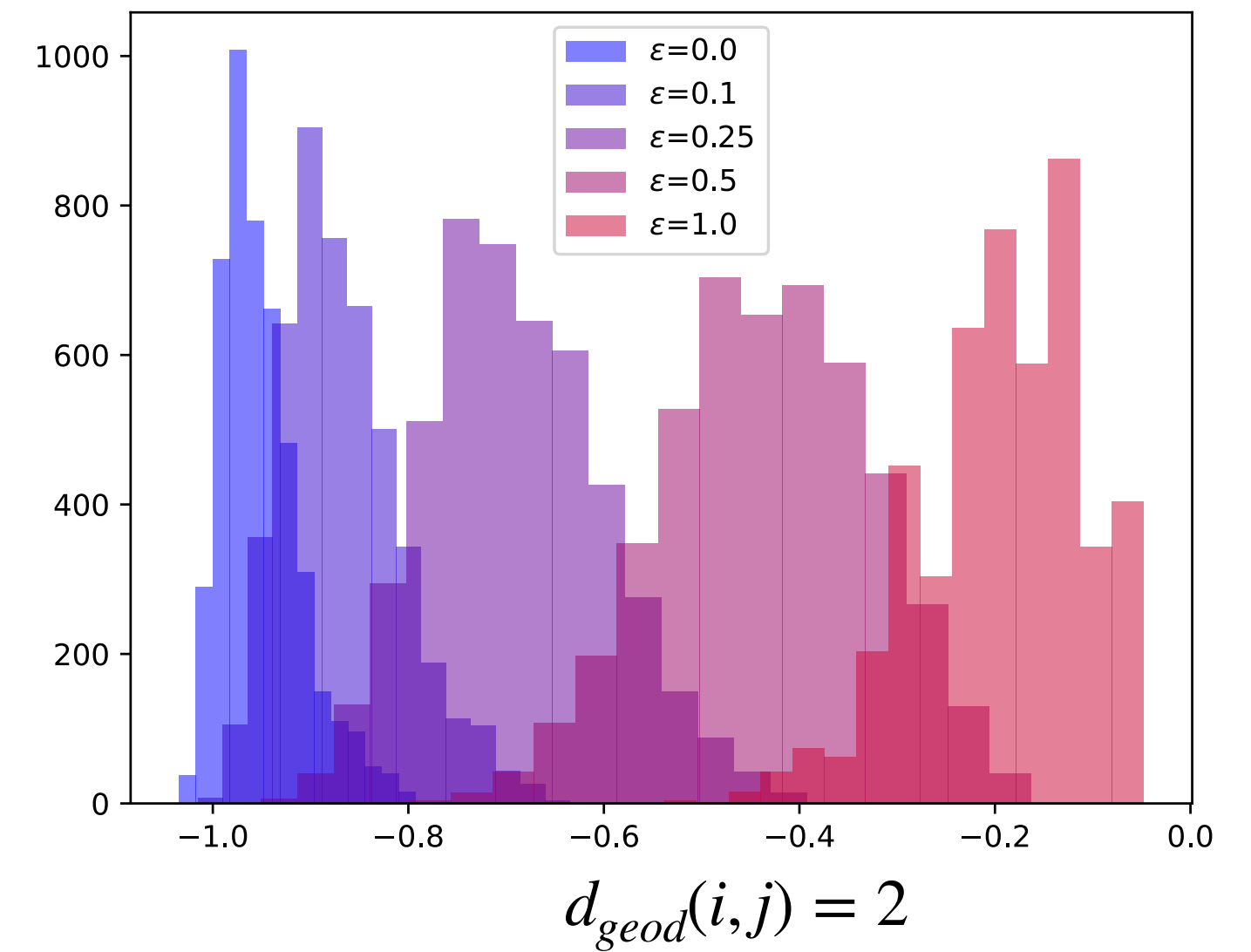
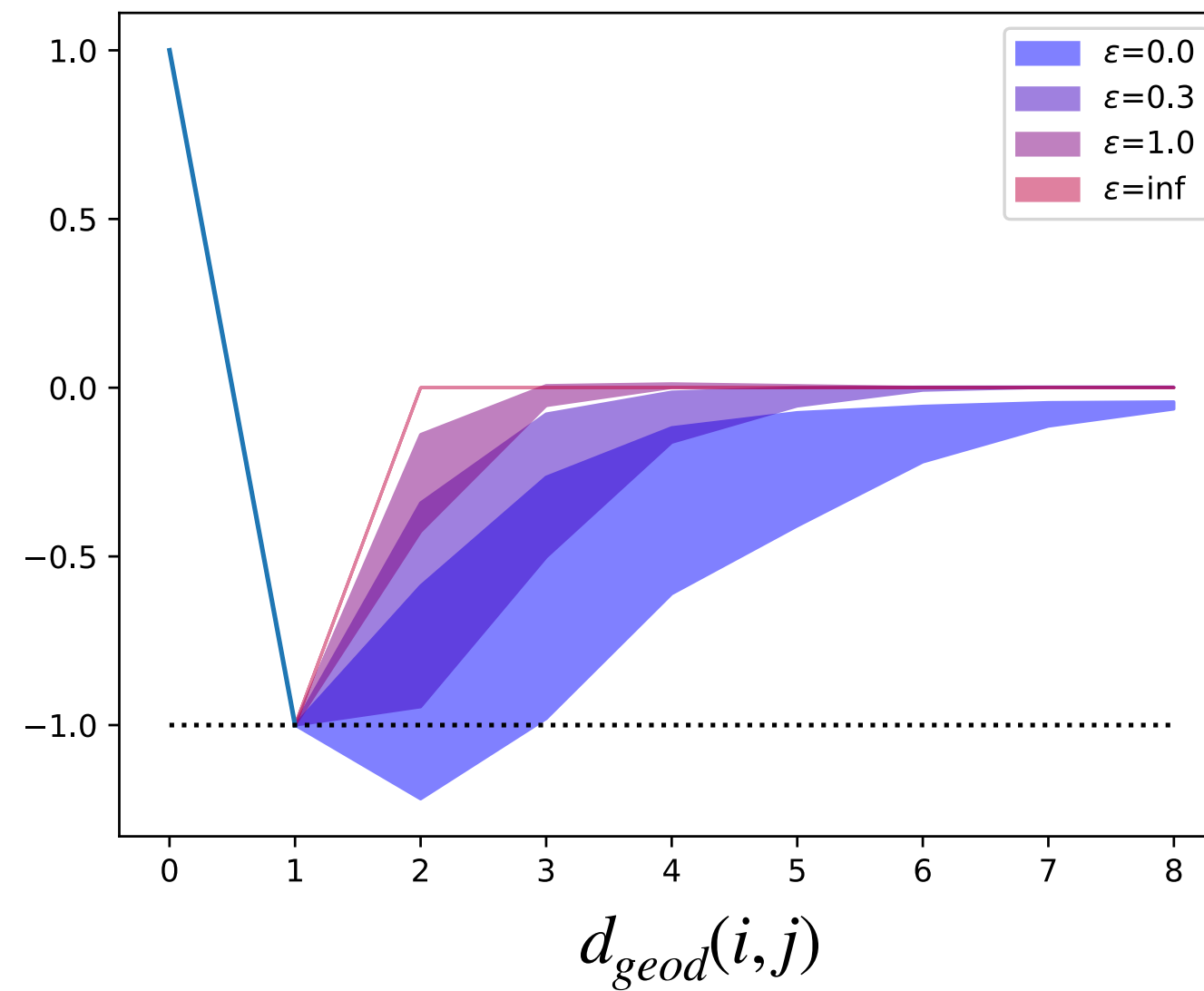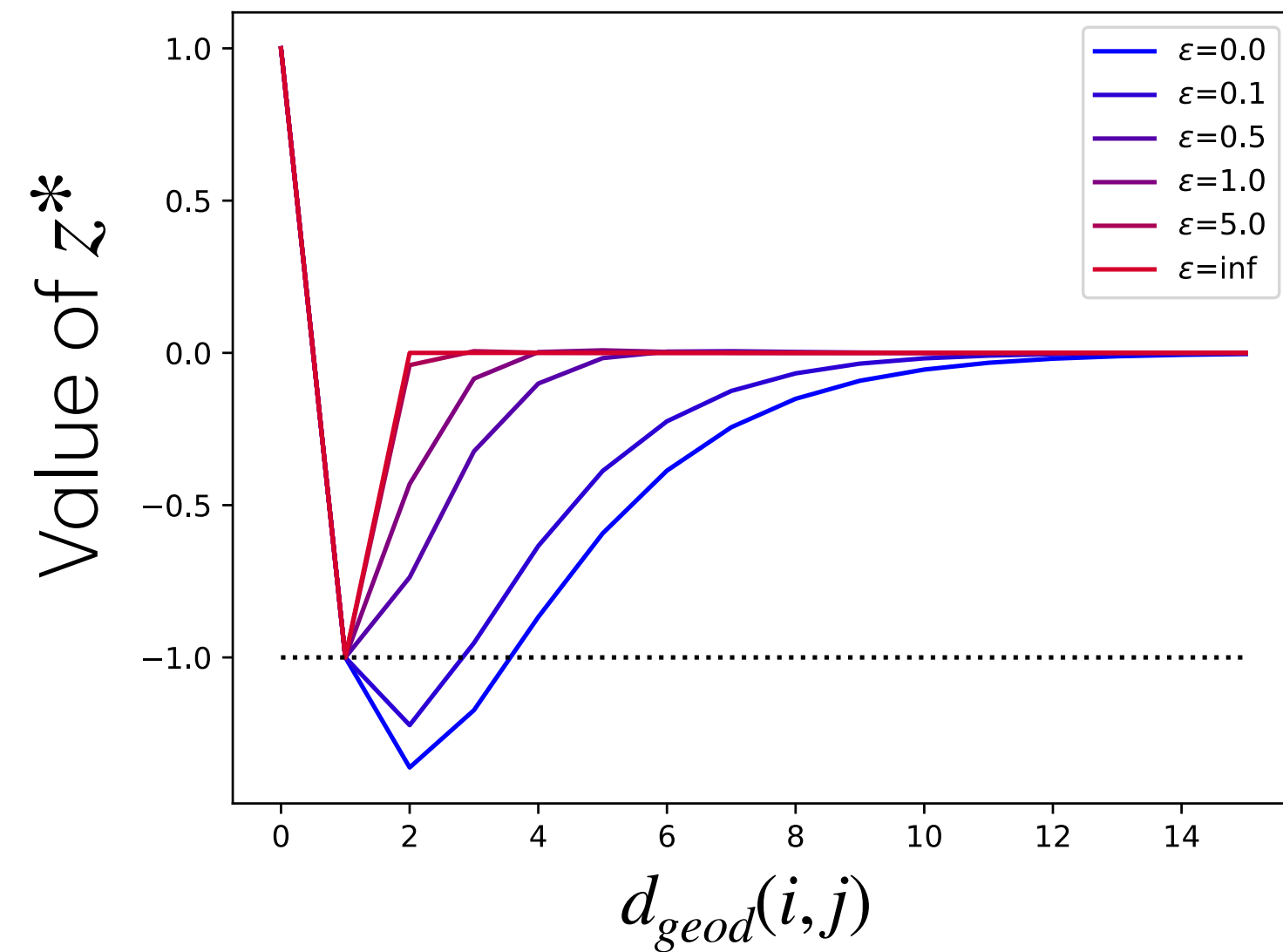$\theta^\star = \delta I + \mathrm{diag}(G1) - G$



Circular

Planar

Erdös-Rényi

# Numerical illustrations

$\epsilon = 0.1$

$\epsilon = 1$

$\epsilon = 10$



Recovery performance of $\ell_1$-iOT for a circular graph.

# The iJKO loss

$$L_r(\theta) := \langle V_\theta, \alpha^{k+1} \rangle - \inf_{\alpha \in \mathscr{P}(\mathscr{X})} \langle V_\theta, \alpha \rangle + \frac{1}{\tau} W_{2,\epsilon}^2(\alpha, \alpha^k \,|\, \alpha \otimes \alpha^k) + r\mathrm{KL}(\alpha \,|\, \alpha^{k+1}).$$

$$\theta^s = \operatorname*{argmin}_{\theta} \frac{\lambda}{r} R(\theta) + L_r(\theta) \quad \xrightarrow{\;r \to \infty\;} \quad \operatorname*{argmin}_{\theta} \lambda R(\theta) + \mathrm{Var}_{\alpha^{k+1}} \left[ V_\theta + \tau^{-1} f^*(\alpha^{k+1}, \alpha^k) \right]$$

**Gaussian experiment**

Consider $\dfrac{d}{dt} X_t = -\nabla V(X_t)$ where $X_0 \sim \mathscr{N}(\mathbf{m}^\star, \Sigma^\star)$ and $V(x) = x^\top \theta^\star x$.

Then, $\alpha_t := \mathrm{law}(X_t) = \mathscr{N}(\mathbf{m}_t, \Sigma_t)$ with $\mathbf{m}_t = e^{-2t\theta^\star} \mathbf{m}^\star$ and $\Sigma_t = e^{-2t\theta^\star} \Sigma^\star e^{-2t\theta^\star}$.
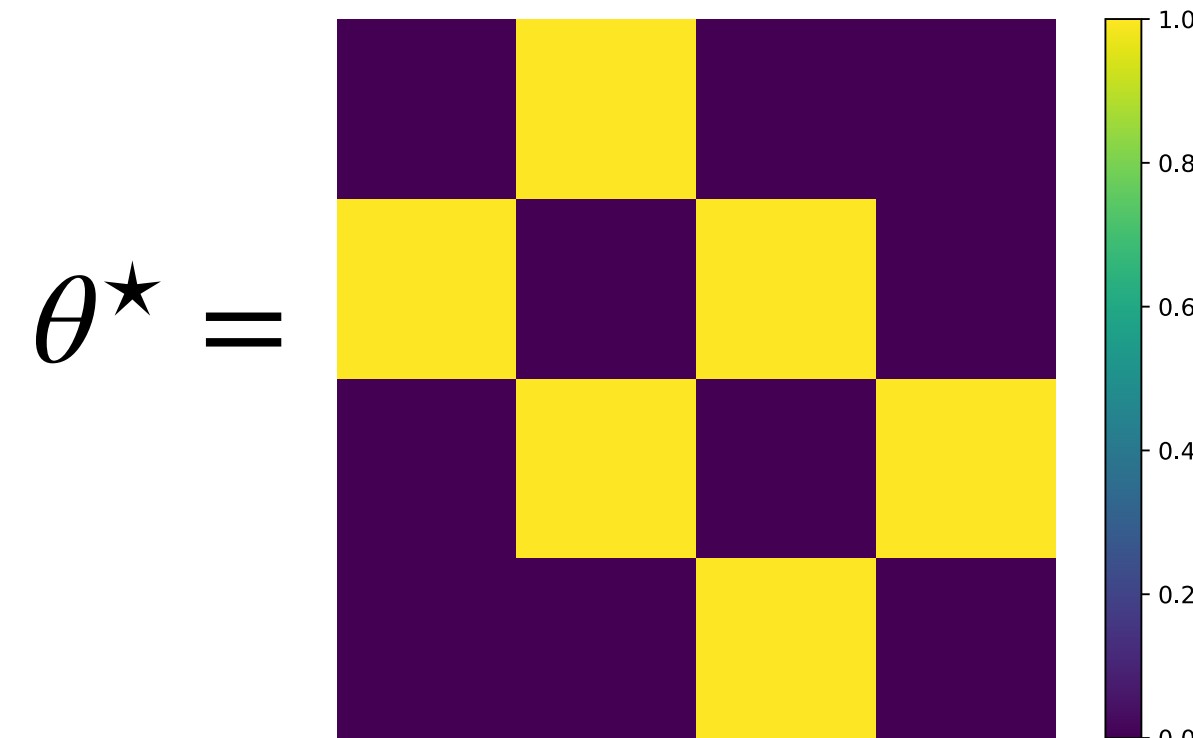
Suppose we observe samples of $\alpha_t$ at discrete time point.

*Question: what kind of $\theta^\star$ are easy to recover?*

c.f. loss $\| \nabla V_\theta + \tau^{-1} \nabla f^*(\alpha^{k+1}, \alpha^k) \|_{L^2(\alpha^{k+1})}^2$ of Terpin, Lanzetti, Gadea, and Dörfler. *Learning diffusion at lightspeed*
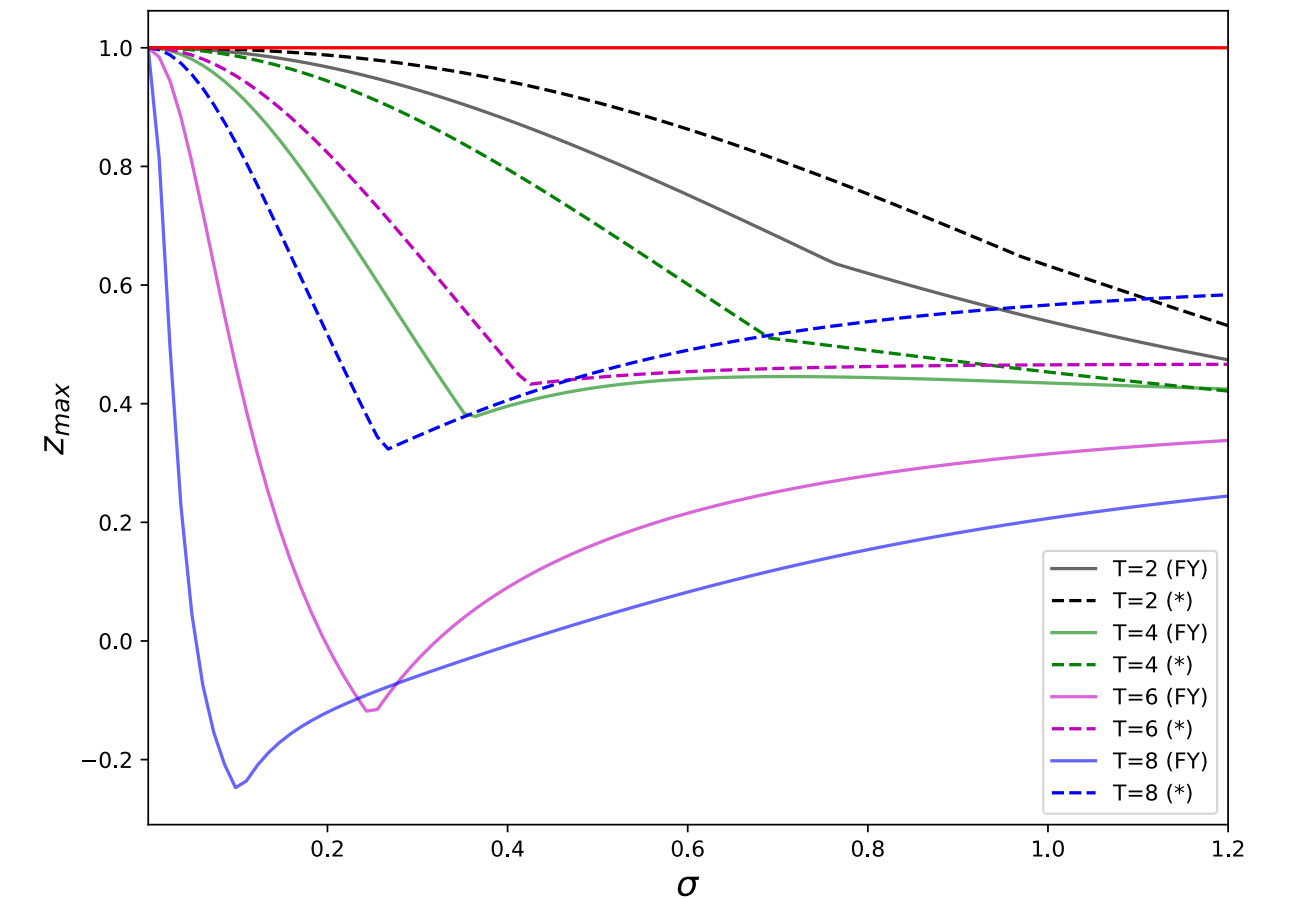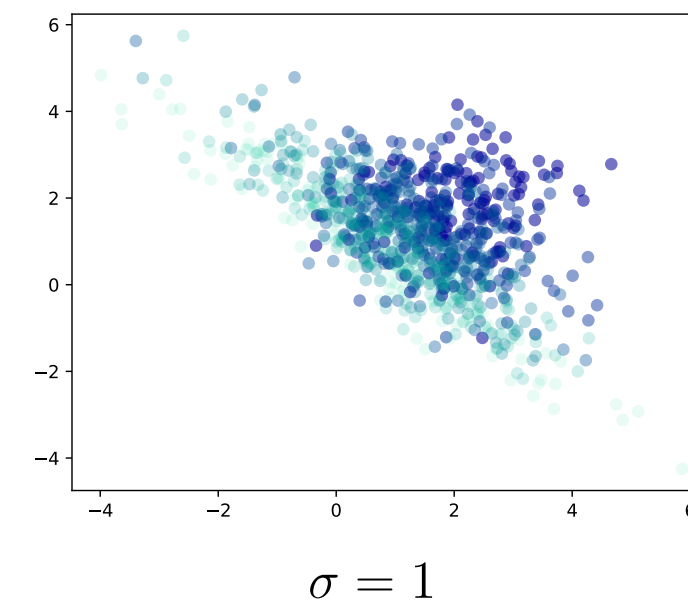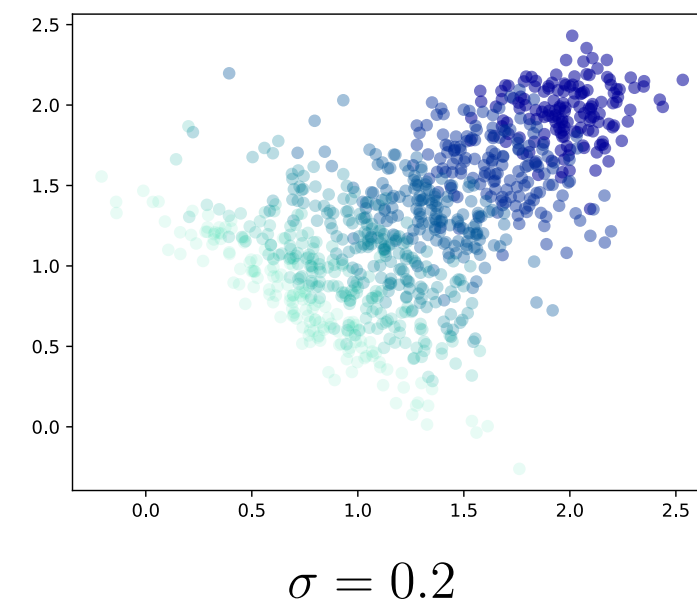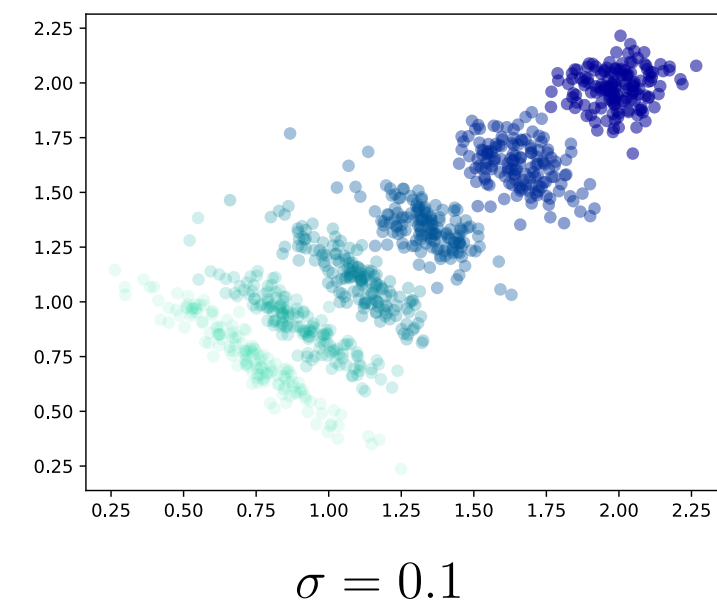
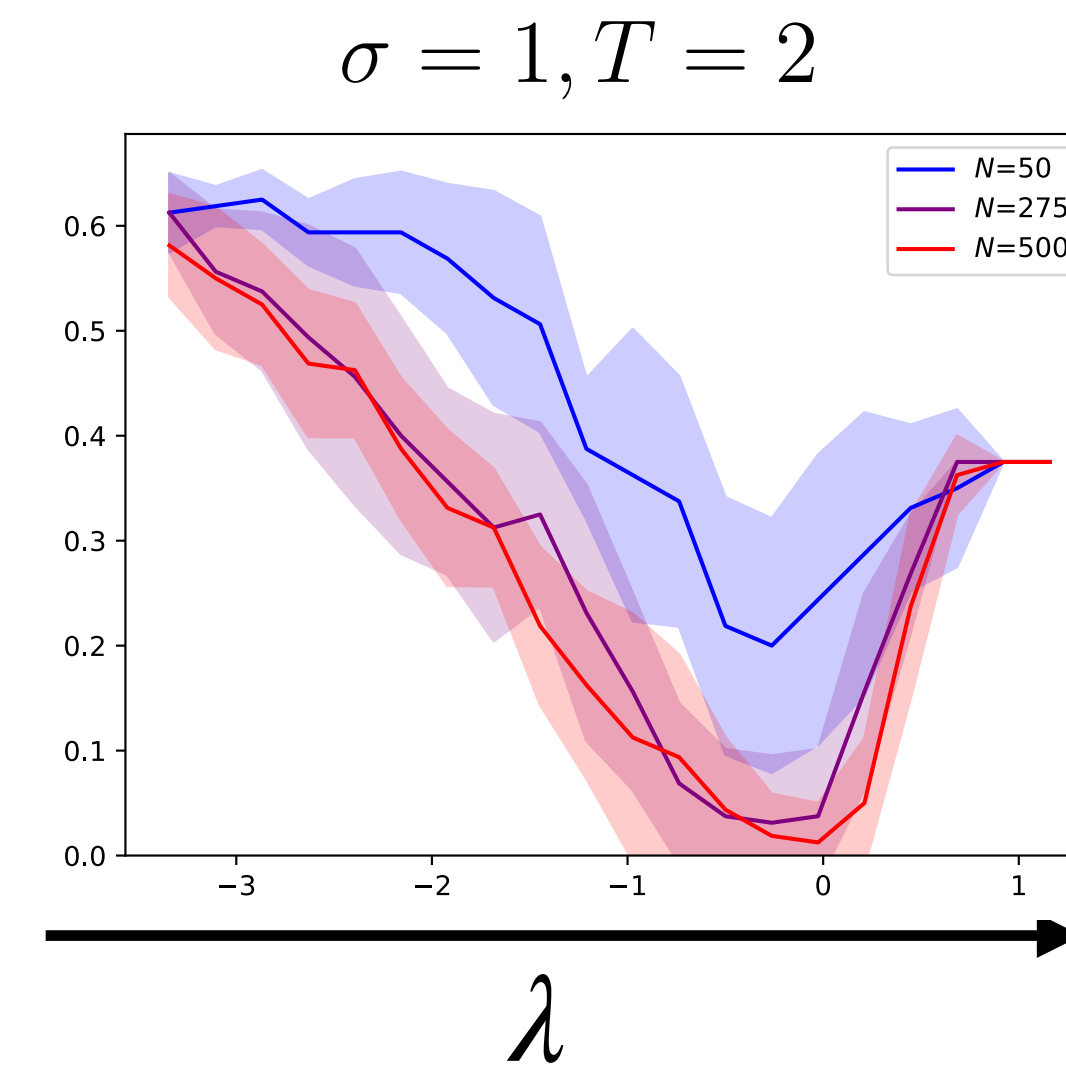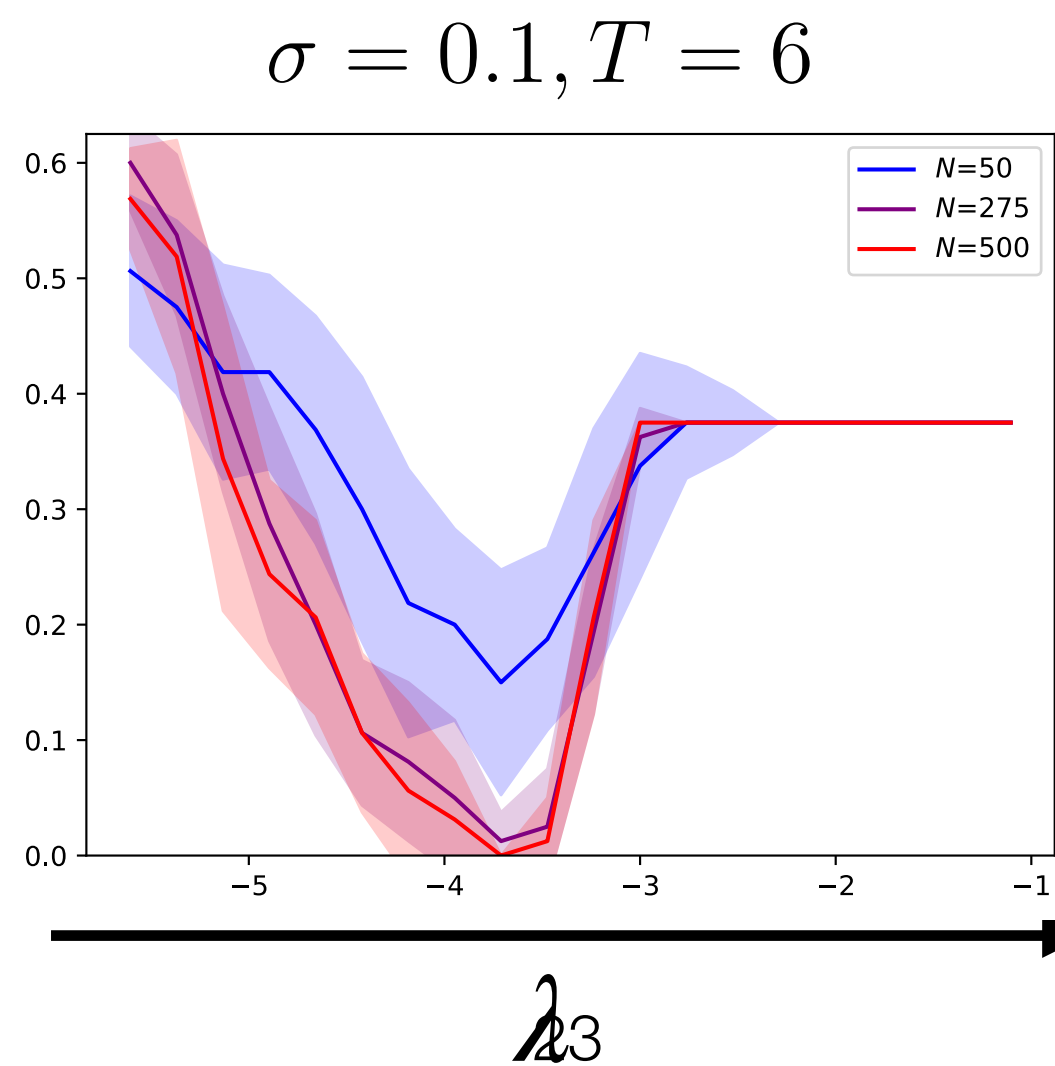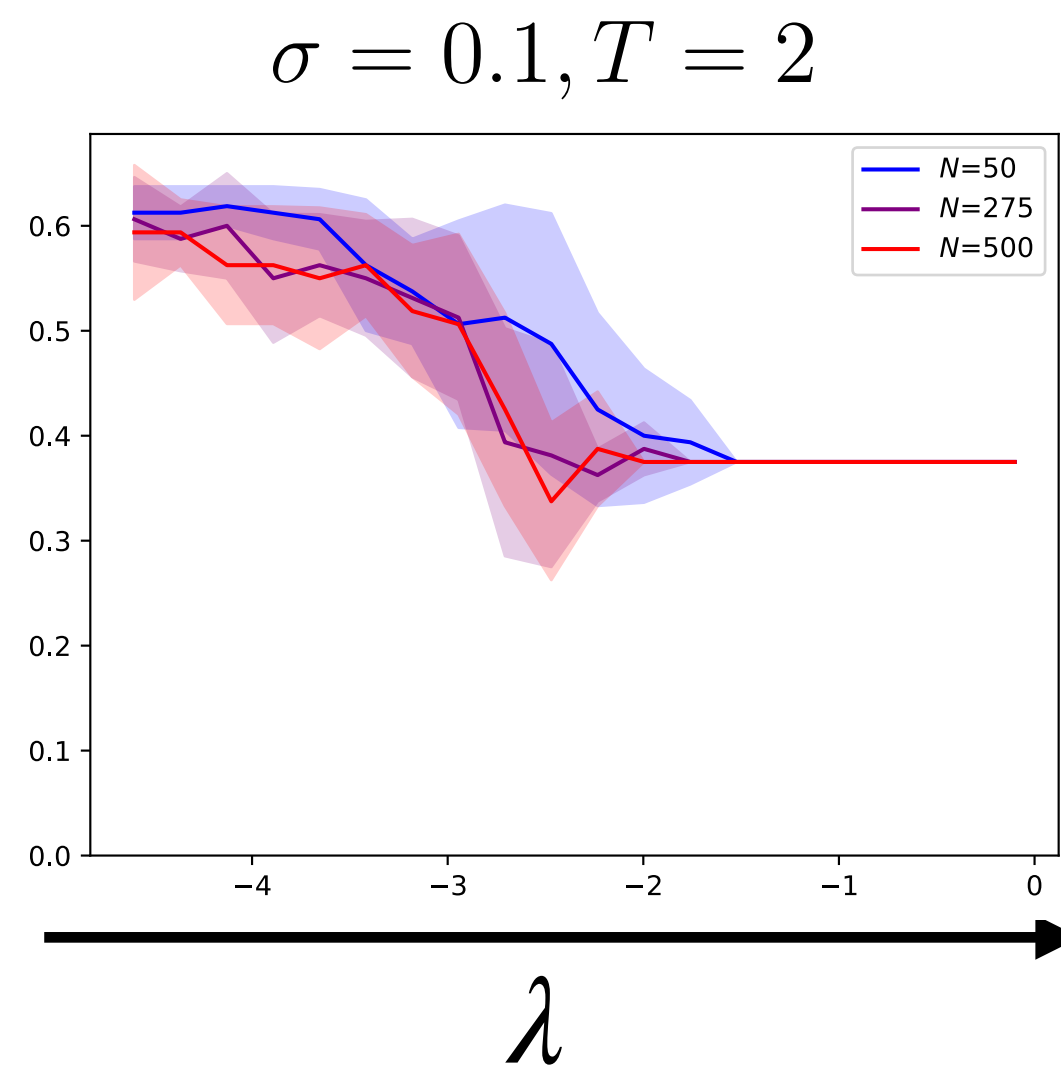Investigate using the closed form certificate at $r = \infty$

Evolution across 6 time points with $\tau = 0.1$, $m^\star = 2 \cdot \mathbf{1}$ and $\Sigma^\star = \sigma^2\mathrm{Id}$.
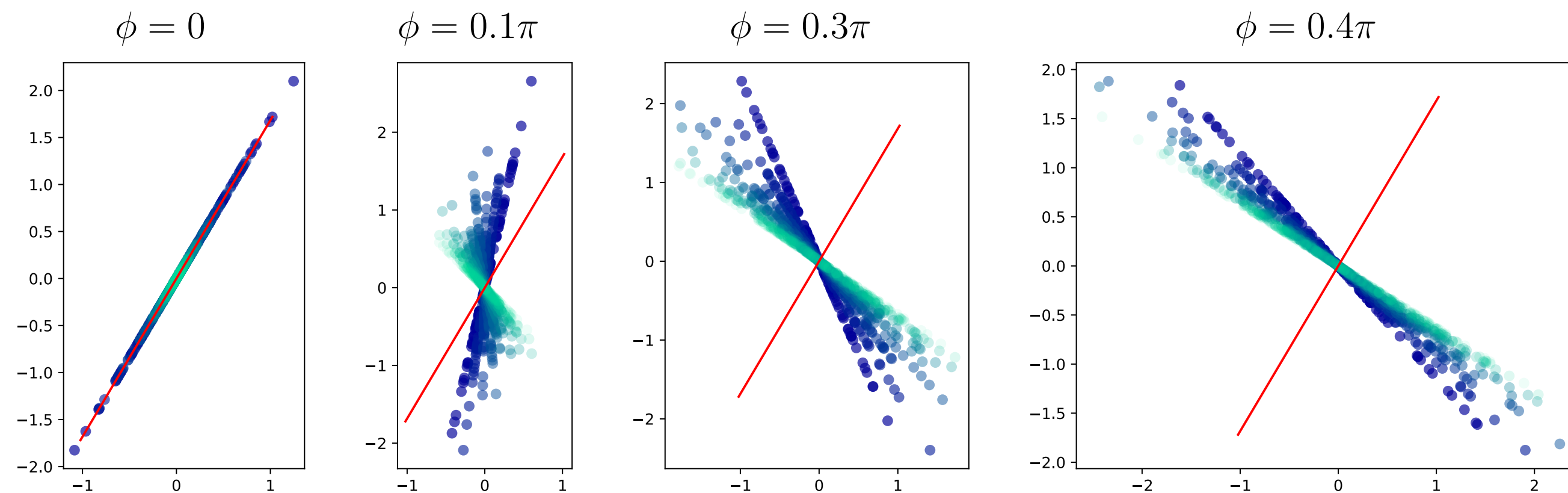
$\theta^\star =$

$\sigma = 0.1$

$\sigma = 0.2$

$\sigma = 1$

$\Sigma^\star = \sigma^2 I, \quad m^\star = 2 \cdot \mathbf{1}$

$z_{max}$

$\omega/\pi$

$\sigma$

Number of wrongly estimated positions

$\sigma = 0.1, T = 2$

$\sigma = 0.1, T = 6$

$\sigma = 1, T = 2$

# The low ra                  t



Nondegeneracy in the low rank setting.



$\phi = 0$  $\phi = 0.1\pi$  $\phi = 0.3\pi$  $\phi = 0.4\pi$

$\omega = 0$  $\omega = 0.1\pi$  $\omega = 0.4\pi$

$$\theta^\star = uu^\top$$

$$\Sigma^\star = \delta I + u_\phi u_\phi^\top \text{ with}$$

$$u_\phi = \cos(\phi)u + \sin(\phi)u^\perp$$

# Summary

- Optimal transport computes a coupling given two distributions and a cost metric.

- In some applications, the metric is unknown; or we might be interested in recovering certain dynamics given observations of probability distributions.

- Fenchel-Young losses allow us to construct convex losses to handle these inverse problems with probability measures.

- We derived a theoretical analysis of the sample complexity, and structural properties of regularization.

*Sparsistency for inverse optimal transport, ICLR 2024*

*Learning from Samples: Inverse Problems over measures via Sharpened Fenchel-Young Losses*