

Quadratically Regularized Optimal Transport

Alberto González Sanz

Columbia University

Joint work with: Marcel Nutz, Andrés Riveros Valdevenito, Gilles Mordant and Alejandro Garriz-Molina

Quadratically Regularized Optimal transport

The quadratically regularized optimal transport is

$$\text{QOT}_\epsilon(P, Q) = \min_{\pi \in \Pi(P, Q)} \int_{\mathbb{R}^d \times \mathbb{R}^d} \frac{1}{2} \|x - y\|^2 d\pi(x, y) + \frac{\epsilon}{2} \underbrace{\left\| \frac{d\pi}{d(P \times Q)} \right\|_{L^2(P \times Q)}^2},$$

Quadratic entropy penalty

This talk we will see the some properties of QOT

Duality and shape of solutions (sparsity)

Differences between EOT and QOT

Rates of convergence

Open problems and conjectures

Optimal transport

Let P and Q be probability measures

$$\text{OT}(P, Q) = \min_{\pi \in \Pi(P, Q)} \int_{\mathbb{R}^d \times \mathbb{R}^d} \frac{1}{2} \|x - y\|^2 d\pi(x, y)$$

where $\Pi(P, Q)$ is the set of couplings between P and Q

The solution, π^* , is called **optimal transport plan**

$\text{OT}(P, Q)$, is called **optimal transport cost**

The optimal transport plan is concentrated on the graph of the sub-differential of a l.s.c. convex function

If P is absolutely continuous wrt Lebesgue: $\pi^* = (I \times \nabla \varphi) \# P$, where φ is a l.s.c. convex function

Duality

Let P and Q be probability measures

$$\text{OT}(P, Q) = \sup_{f(x)+g(y) \leq \|x-y\|^2} \int f(x)dP(x) + \int g(y)dQ(y),$$

A solution of the dual problem will be a pair (f, g) of functions where

$$(\varphi, \psi) = (\|\cdot\|^2/2 - f, \|\cdot\|^2/2 - g)$$

are conjugate l.s.c. convex functions and

If P is absolutely continuous wrt Lebesgue: $\pi^* = (I \times \nabla \varphi) \# P$

Finite sample approximation

Let $P_n = \frac{1}{n} \sum_{i=1}^n \delta_{X_i}$ and $Q_n = \frac{1}{n} \sum_{i=1}^n \delta_{Y_i}$ be empirical measures

$$\text{OT}(P_n, Q_n) = \frac{1}{n} \min_{\pi} \langle C, \pi \rangle_{Fr} : \quad \text{s.t.} \quad \pi \in \Omega_n$$

where Ω_n is the **Birkhoff polytope of doubly stochastic matrices**

$$\Omega_n = \left\{ \pi \in \mathbb{R}^{n \times n} : \sum_{i=1}^n \pi_{i,j} = 1, \sum_{j=1}^n \pi_{i,j} = 1, \pi_{i,j} \geq 0 \right\}$$

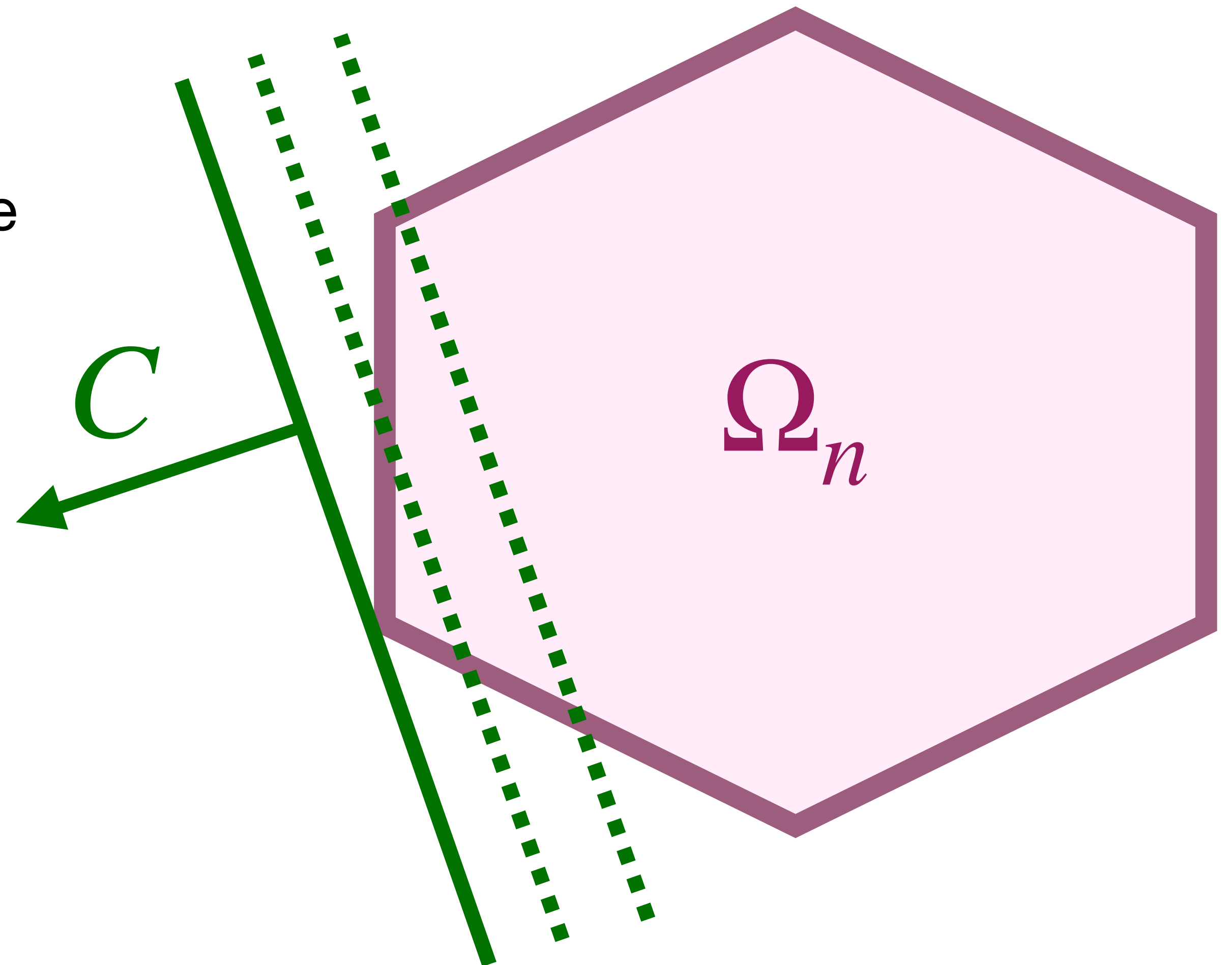
and $C = (\|X_i - Y_j\|^2)_{i,j}$ is the **cost matrix**

Finite sample approximation

OT is a **linear program**

The vertexes of the Birkhoff polytope are the **permutation matrices**

The empirical OT plans are **sparse**



Regularized Optimal transport

The entropy regularized optimal transport is

$$\text{EOT}_\epsilon(P, Q) = \min_{\pi \in \Pi(P, Q)} \int_{\mathbb{R}^d \times \mathbb{R}^d} \frac{1}{2} \|x - y\|^2 d\pi(x, y) + \underbrace{\epsilon H(\pi | P \times Q)}_{\text{Logarithmic entropy penalty}},$$

$$H(\alpha | \beta) = \begin{cases} \int \log\left(\frac{d\alpha}{d\beta}(x)\right) d\alpha(x) & \text{if } \alpha \ll \beta \\ +\infty & \text{otherwise} \end{cases}$$

Regularized Optimal transport

This problem can also be written in its dual formulation

$$\text{EOT}_\epsilon(P, Q) = \sup_{\substack{f \in L_1(P) \\ g \in L_1(Q)}} E \left(f(X) + g(Y) - \epsilon e^{\frac{f(X) + g(Y) - \frac{1}{2} \|X - Y\|^2}{\epsilon}} \right) + \epsilon,$$

We pass from a linear program to an strictly concave one

with $X \sim P$, $Y \sim Q$, where the two variables are independent. The solutions satisfy

$$f_{P,Q} = -\epsilon \log \left(\int e^{\frac{g_{P,Q}(y) - \frac{1}{2} \|\cdot - y\|^2}{\epsilon}} dQ(y) \right) \quad g_{P,Q} = -\epsilon \log \left(\int e^{\frac{f_{P,Q}(x) - \frac{1}{2} \|x - \cdot\|^2}{\epsilon}} dP(x) \right)$$

Regularized Optimal transport

Efficient computation **Sinkhorn algorithm**

Exponential convergence in the fixed point iterations
(Franklin and Lorenz (1989) Carlier (2022),

Instability when ϵ is small

The dual solutions are **smooth**

Some **regularity properties** of OT can be obtained via covariance inequalities (Chewi, Pooladian (2022))

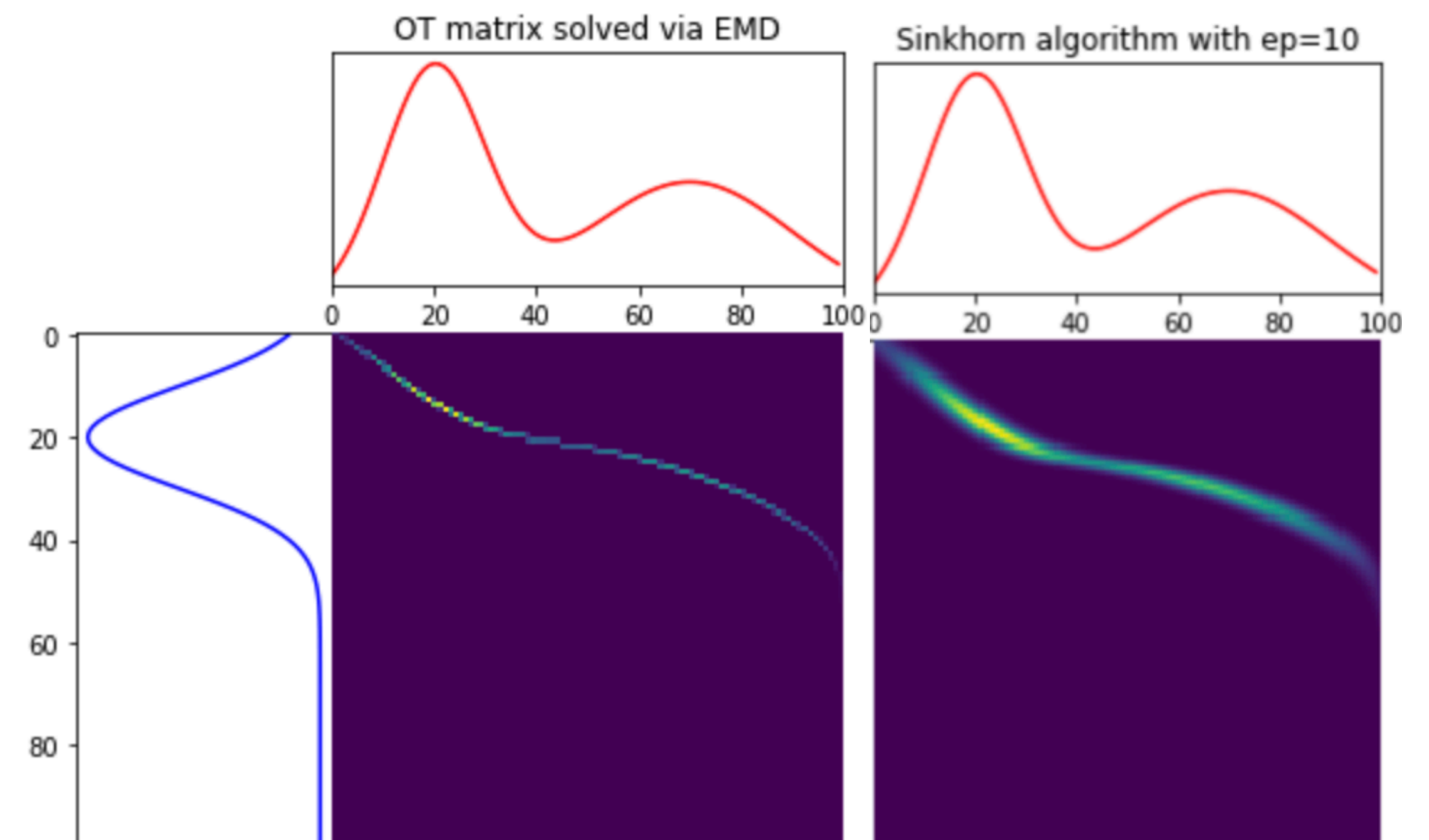
Reduction of **statistical complexity** (Donsker class)
Genevay et al. (2018), Mena and Weed (2019)

The regularized plan is a **noisy** approximation of the OT plan

Similar behaviour to a **Gaussian convolution** of OT (Pal (2019))

The regularized plan has **full support**

```
def sinkhorn(a, b, C, epsilon=0.1, max_iters=1000):  
    ## a,b=weights of measures C=Cost matrix  
    K=np.exp(-C/epsilon)  
    v=np.ones(b.shape[0])  
    for i in range(max_iters):  
        u=a/K.dot(v)  
        v=b/K.T.dot(u)  
    # Alternate projections  
  
    return np.diag(u).dot(K).dot(np.diag(v))
```



Quadratically Regularized Optimal transport

The quadratically regularized optimal transport is

$$\text{QOT}_\epsilon(P, Q) = \min_{\pi \in \Pi(P, Q)} \int_{\mathbb{R}^d \times \mathbb{R}^d} \frac{1}{2} \|x - y\|^2 d\pi(x, y) + \frac{\epsilon}{2} \underbrace{\left\| \frac{d\pi}{d(P \times Q)} \right\|_{L^2(P \times Q)}^2},$$

Quadratic entropy penalty

This talk we will see the some properties of QOT

Duality and shape of solutions (sparsity)

Differences between EOT and QOT

Rates of convergence

Open problems and conjectures

Quadratically Regularized Optimal transport

[A. Dessein et al. \(2018\)](#) considered optimal transport with convex regularization.

[Essid and Solomon \(2018\)](#) studied quadratic regularization for a minimum-cost flow problem on a graph, including discrete optimal transport as a special case.

[M. Blondel et al. \(2019\)](#) explored QOT in the discrete setting, with experiments highlighting the sparsity and theoretical results including convergence for small regularization parameter. They proposed an application to image processing.

[Li et al. \(2020\)](#) computes regularized Wasserstein barycenters using neural networks and finds that the quadratic penalty produces sharper results than the logarithmic.

The first work rigorously addressing a continuous setting is [Lorenz et al. \(2021\)](#). The authors derive duality results and present two algorithms, a nonlinear Gauss–Seidel method and a semismooth Newton method.

[Zhang et al. \(2023\)](#) uses quadratic regularization in a manifold learning task related to single cell RNA sequencing and notes that sparsity is crucial to avoid biasing the affinity matrix.

Dual formulation

This problem can also be written in its dual formulation

$$\text{QOT}_\epsilon(P, Q) = \sup_{a, b \in L^2(P) \times L^2(Q)} \int a(x) dP(x) + \int b(y) dQ(y) - \frac{1}{2\epsilon} \left(a(x) + b(y) - \frac{\epsilon}{2} \|x - y\|^2 \right)_+^2 d(P \times Q)(x, y)$$

Primal-dual relation

$$(a_\epsilon, b_\epsilon) \text{ solves Dual} \quad \longleftrightarrow \quad \frac{1}{\epsilon} \left(a_\epsilon(x) + b_\epsilon(y) - \frac{1}{2} \|x - y\|^2 \right)_+ d(P \times Q)(x, y) \text{ solves Primal}$$

Optimality Conditions

QOT

$$\left\{ \begin{array}{l} \int \left(a_\epsilon(x) + b_\epsilon(y) - \frac{1}{2} \|x - y\|^2 \right)_+ dQ(y) = \epsilon \quad P - \text{a.e. } x \\ \int \left(a_\epsilon(x) + b_\epsilon(y) - \frac{1}{2} \|x - y\|^2 \right)_+ dP(x) = \epsilon \quad Q - \text{a.e. } y \end{array} \right.$$

EOT

$$\left\{ \begin{array}{l} \int e^{\frac{a_\epsilon(x) + b_\epsilon(y) - \frac{1}{2} \|x - y\|^2}{\epsilon}} dQ(y) = 1 \quad P - \text{a.e. } x \\ \int e^{\frac{a_\epsilon(x) + b_\epsilon(y) - \frac{1}{2} \|x - y\|^2}{\epsilon}} dP(x) = 1 \quad Q - \text{a.e. } y \end{array} \right.$$

Empirical approximation

Let $P_n = \frac{1}{n} \sum_{i=1}^n \delta_{X_i}$ and $Q_n = \frac{1}{n} \sum_{i=1}^n \delta_{Y_i}$ be empirical measures

$$\text{QOT}_\epsilon(P_n, Q_n) = \frac{1}{n} \min_{\pi} \langle C, \pi \rangle_{Fr} + \frac{\epsilon}{2} \|\pi\|^2 : \quad \text{s.t.} \quad \pi \in \Omega_n$$

where Ω_n is the **Birkhoff polytope of doubly stochastic matrices**

$$\Omega_n = \left\{ \pi \in \mathbb{R}^{n \times n} : \sum_{i=1}^n \pi_{i,j} = 1, \sum_{j=1}^n \pi_{i,j} = 1, \pi_{i,j} \geq 0 \right\}$$

and $C = (\|X_i - Y_j\|^2)_{i,j}$ is the **cost matrix**

Stationary approximation

$$\text{QOT}_\epsilon(P_n, Q_n) = \frac{1}{n} \min_{\pi} \langle C, \pi \rangle_{Fr} + \frac{\epsilon}{2} \|\pi\|^2 : \quad \text{s.t.} \quad \pi \in \Omega_n$$

QOT is a **quadratically regularized linear program**

The QOT plan is the **projection** of $-\frac{C}{2\epsilon}$

(Mangasarian, Meyer, 1979) As $\epsilon \rightarrow 0$ the convergence of the QOT plans towards the OT plans is stationary:

There exists $\epsilon_0 > 0$ such that $(\text{QOT plan}) \in (\text{OT plans}) \quad \forall \epsilon \leq \epsilon_0$

Stationary convergence

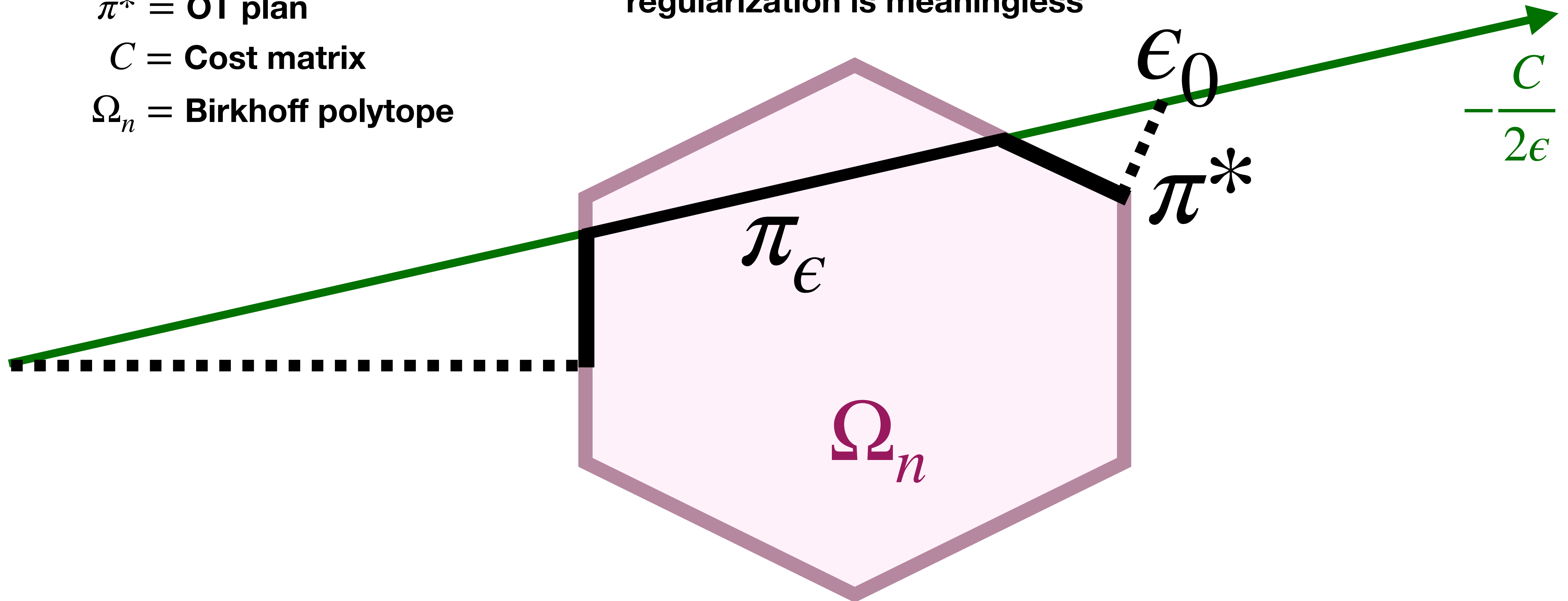
$\pi_\epsilon =$ QOT plan at ϵ

$\epsilon_0 =$ Regularization parameter where the effect of regularization is meaningless

$\pi^* =$ OT plan

$C =$ Cost matrix

$\Omega_n =$ Birkhoff polytope



Stationary convergence

Theorem (GS, Nutz, 2024)

$$(\text{QOT plan}) \in (\text{OT plans}) \iff \epsilon^{-1} \leq (\epsilon^*)^{-1} := 2N \cdot \max_{\pi \in \text{Permutations} \setminus \text{OTplans}} \frac{\langle \pi^*, \pi^* - \pi \rangle}{\langle C, \pi - \pi^* \rangle},$$

where π^* is the OT plan with smallest norm

Example rates as sample size increases If the points are a uniform grid of $[0,1]$

$$\epsilon^* = \frac{1}{2N^3}$$

Comparison with EOT

Theorem (Niles-Weed, 2021) The convergence of EOT plans is exponentially fast

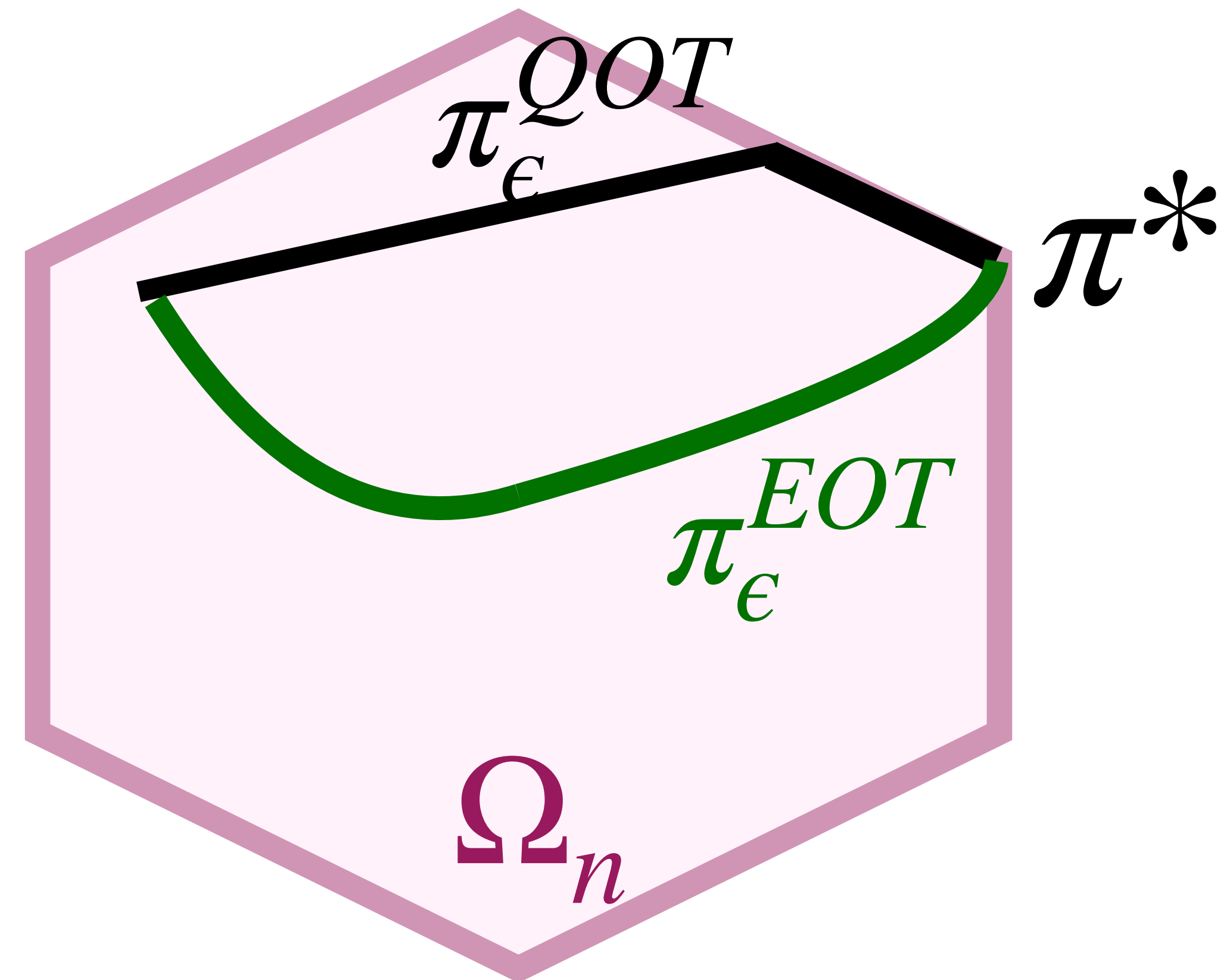
Why the convergence of EOT is not stationary?

EOT approaches OT from inside of the polytope

QOT approaches OT by changing to better faces

The relative interior of Ω_n are the couplings with full support

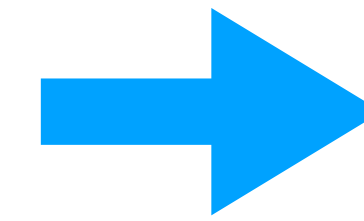
Changing from faces to surfaces of Ω_n is means creating new zeroes



Sparsity

EOT approaches OT from inside of the polytope

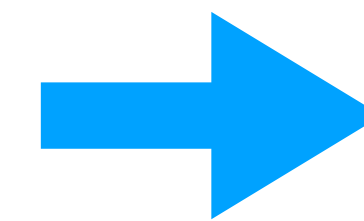
The relative interior of Ω_n are the couplings with full support



EOT plans full support
(All the entries of the matrix are strictly positive)

QOT approaches OT by changing to better faces

Changing from faces to subfaces of Ω_n means creating new zeroes



QOT plans sparse
(Progressively new zeroes are created)

Is the creation of zeroes monotone? That is, is the support of the QOT plan decreasing monotonously?

Non Monotonicity

Is the creation of zeroes monotone? That is, is the support of the QOT plan decreasing monotonously?



Each time the QOT plan enters in a face of the polytope it remains in that face



(GS, Nutz, Riveros Valdevenito, 2024)

The point of minimum norm of each face of the polytope belongs to the relative interior of that face

Theorem (GS, Nutz, Riveros Valdevenito, 2024)

The monotonicity of the support fails for n larger or equal than 5 and it is true otherwise. That is, for $n \geq 5$, there exists a configuration of points (or respectively a cost matrix) such that a zero created becomes positive for smaller regularization parameter

Non Monotonicity

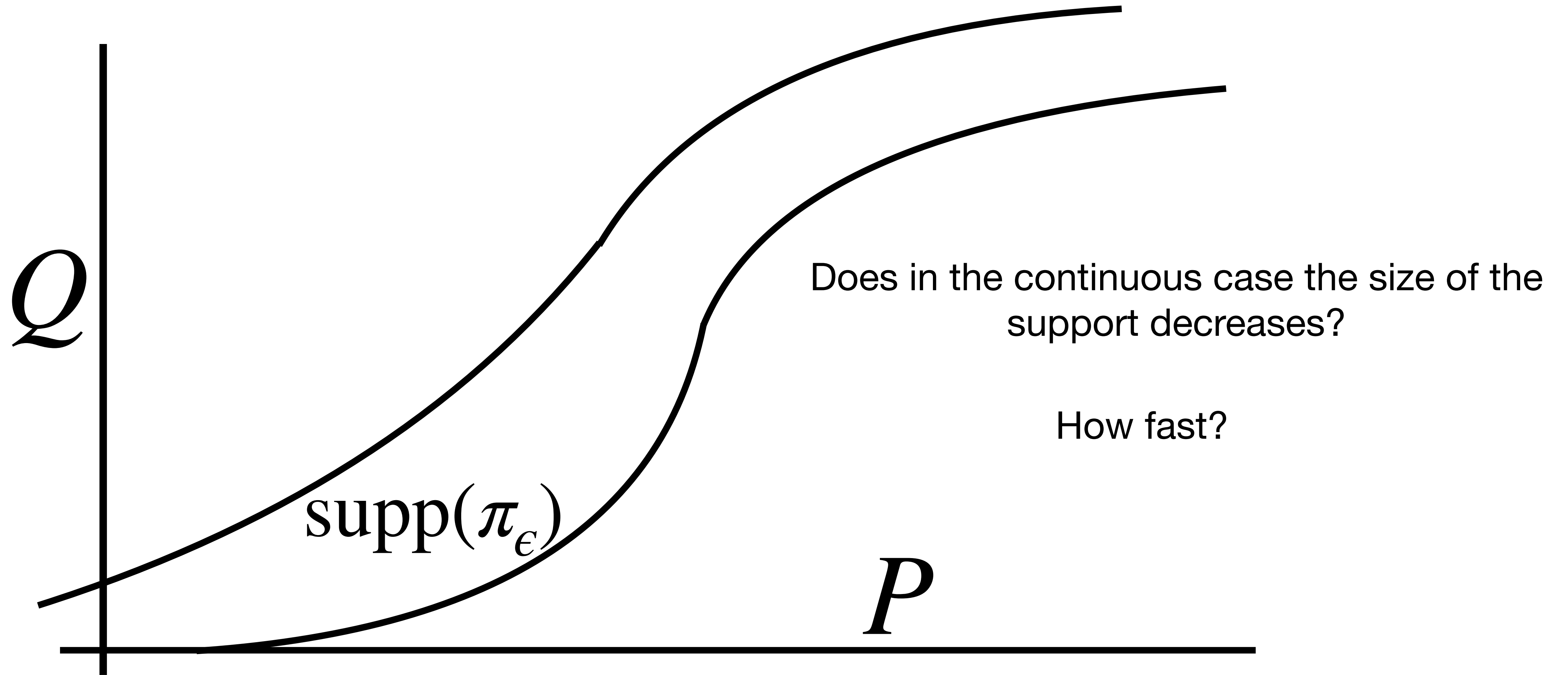
Theorem (GS, Nutz, Riveros Valdevenito, 2024)

The monotonicity of the support fails for n larger or equal than 5 and it is true otherwise. That is, for $n \geq 5$, there exists a configuration of points (or respectively a cost matrix) such that a zero created becomes positive for smaller regularization parameter.

$$\text{Cost} = \begin{bmatrix} -1.1 & -1 & -1 & -1 & -1 \\ -1 & -1.1 & 0 & 0 & 0 \\ -1 & 0 & -1.1 & 0 & 0 \\ -1 & 0 & 0 & -1.1 & 0 \\ -1 & 0 & 0 & 0 & -1.1 \end{bmatrix}.$$

$$\pi_{1/2.5} = \begin{bmatrix} 0 & 0.05 & 0.05 & 0.05 & 0.05 \\ 0.05 & 0.15 & 0 & 0 & 0 \\ 0.05 & 0 & 0.15 & 0 & 0 \\ 0.05 & 0 & 0 & 0.15 & 0 \\ 0.05 & 0 & 0 & 0 & 0.15 \end{bmatrix} \quad \pi_0 = \begin{bmatrix} 0.2 & 0 & 0 & 0 & 0 \\ 0 & 0.2 & 0 & 0 & 0 \\ 0 & 0 & 0.2 & 0 & 0 \\ 0 & 0 & 0 & 0.2 & 0 \\ 0 & 0 & 0 & 0 & 0.2 \end{bmatrix}$$

Continuous case



Qualitative result

Theorem (Nutz, 2024)

As $\epsilon \rightarrow 0$ the support of the QOT plan tends to the support of the OT plan (graph of OT map) in Hausdorff distance.

The proof is a consequence of the following facts

- The QOT potentials are uniformly Lipschitz
- and they converge uniformly to the OT potentials
- The shape of the conditional support

$$\mathcal{S}_x = \left\{ a_\epsilon(x) + b_\epsilon(y) - \frac{1}{2} \|x - y\|^2 \geq 0 \right\}$$

Explicit solutions

Let us solve it in a place where the solutions are explicit in order to gain some intuition

Consider the marginals $P = Q = \text{Uniform}[0,1]^d$ and the cost

$$d_T^2(x, y) = \frac{1}{2} \inf_{z \in \mathbb{Z}^d} \|x - y - z\|^2$$



The QOT plan has density (for small regularization)

$$\frac{1}{\epsilon} \left(C_d \epsilon^{\frac{2}{d+2}} - d_T^2(x, y) \right)_+$$

As $\epsilon \rightarrow 0$, the diameter of the support of the density tends to zero with rate $\epsilon^{\frac{1}{d+2}}$

As $\epsilon \rightarrow 0$, $\text{QOT}_\epsilon(P, Q) - \text{OT}(P, Q) = O(\epsilon^{\frac{2}{d+2}})$

Geometric properties of QOT

We rewrite things in a more proper way

$$\frac{1}{\epsilon} \left(a_\epsilon(x) + b_\epsilon(y) - \frac{1}{2} \|x - y\|^2 \right)_+ d(P \times Q)(x, y) = \frac{1}{\epsilon} \left(\langle x, y \rangle - f_\epsilon(x) - g_\epsilon(y) \right)_+ d(P \times Q)(x, y)$$

$$\text{where } f_\epsilon(x) = \frac{1}{2} \|x\|^2 - a_\epsilon(x), \quad g_\epsilon(y) = \frac{1}{2} \|y\|^2 - b_\epsilon(y)$$

Lemma (GS, Nutz, 2024)

Assume that P and Q are a.c. w.r.t. Lebesgue with bounded support. Then f_ϵ is a convex function with derivative

$$\nabla f_\epsilon(x) = \frac{\int_{\mathcal{S}_x} y dQ(y)}{Q(\mathcal{S}_x)} \quad \mathcal{S}_x = \{y : \langle x, y \rangle - f_\epsilon(x) - g_\epsilon(y) \geq 0\}$$

Geometric properties of QOT

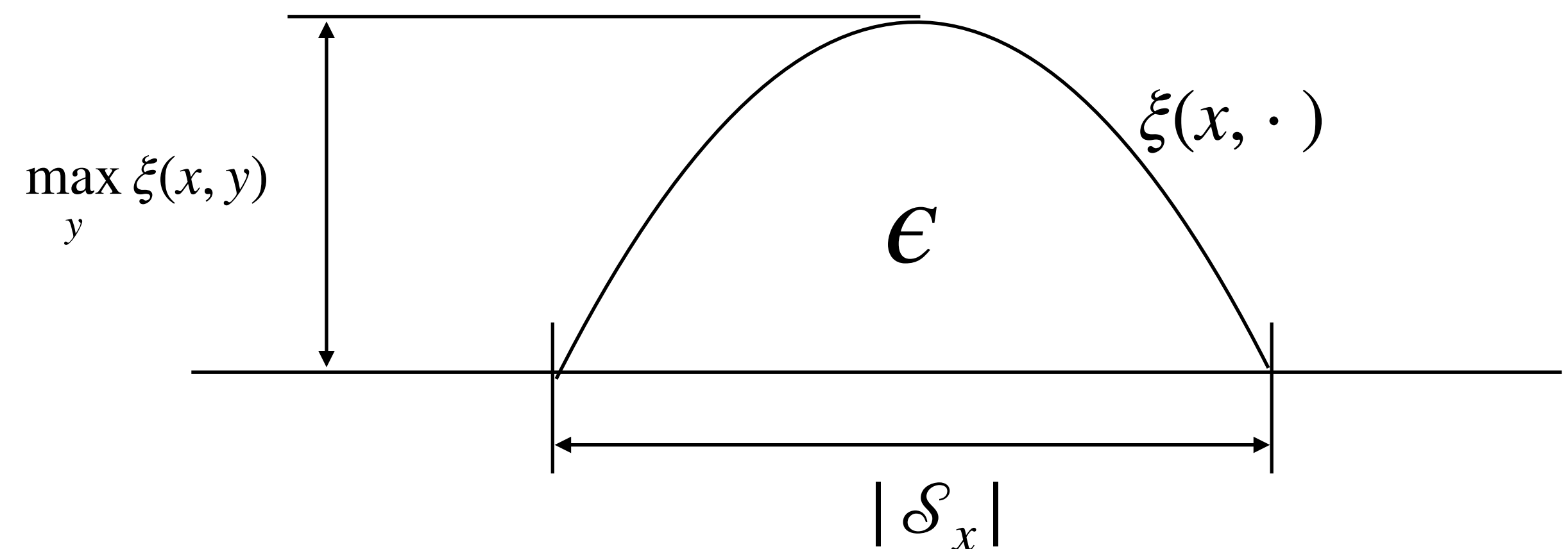
For fixed x the function $y \mapsto \xi(x, y) = \langle x, y \rangle - f_\epsilon(x) - g_\epsilon(y)$ is concave and integrates ϵ .

Therefore, if $P = p\mathbf{1}_{\Omega_0}dx$ and $Q = q\mathbf{1}_{\Omega_1}dx$ with p, q bounded away from zero and infinity, then

$$|\mathcal{S}_x| \max_{y \in \Omega_1} \xi(x, y) \approx \epsilon \quad \text{where} \quad |\mathcal{S}_x| = |\{y : \xi(x, y) \geq 0\}|$$

Then

- We need to understand the relation between the maximum of the paraboloid and its basis
- This involves controlling the second derivative of the paraboloid



Bound on the second derivative

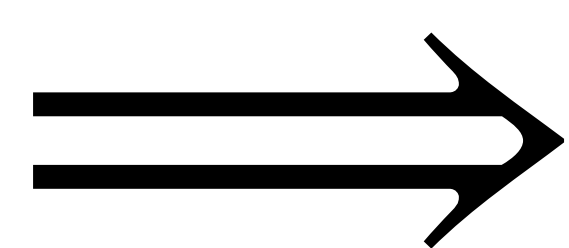
Step 1) Exact shape of the derivative

$$f''_{\varepsilon}(x) = q(y_m(x)) \frac{(f'_{\varepsilon}(x) - y_m(x))^2}{(x - g'_{\varepsilon}(y_m(x)))Q(\mathcal{S}_x)} \chi_{\Omega_0^{(2)} \cup \Omega_0^{(3)}}(x) + q(y_M(x)) \frac{(f'_{\varepsilon}(x) - y_M(x))^2}{(g'_{\varepsilon}(y_M(x)) - x)Q(\mathcal{S}_x)} \chi_{\Omega_0^{(1)} \cup \Omega_0^{(2)}}(x).$$

$$[y_m(x), y_M(x)] = \mathcal{S}_x$$

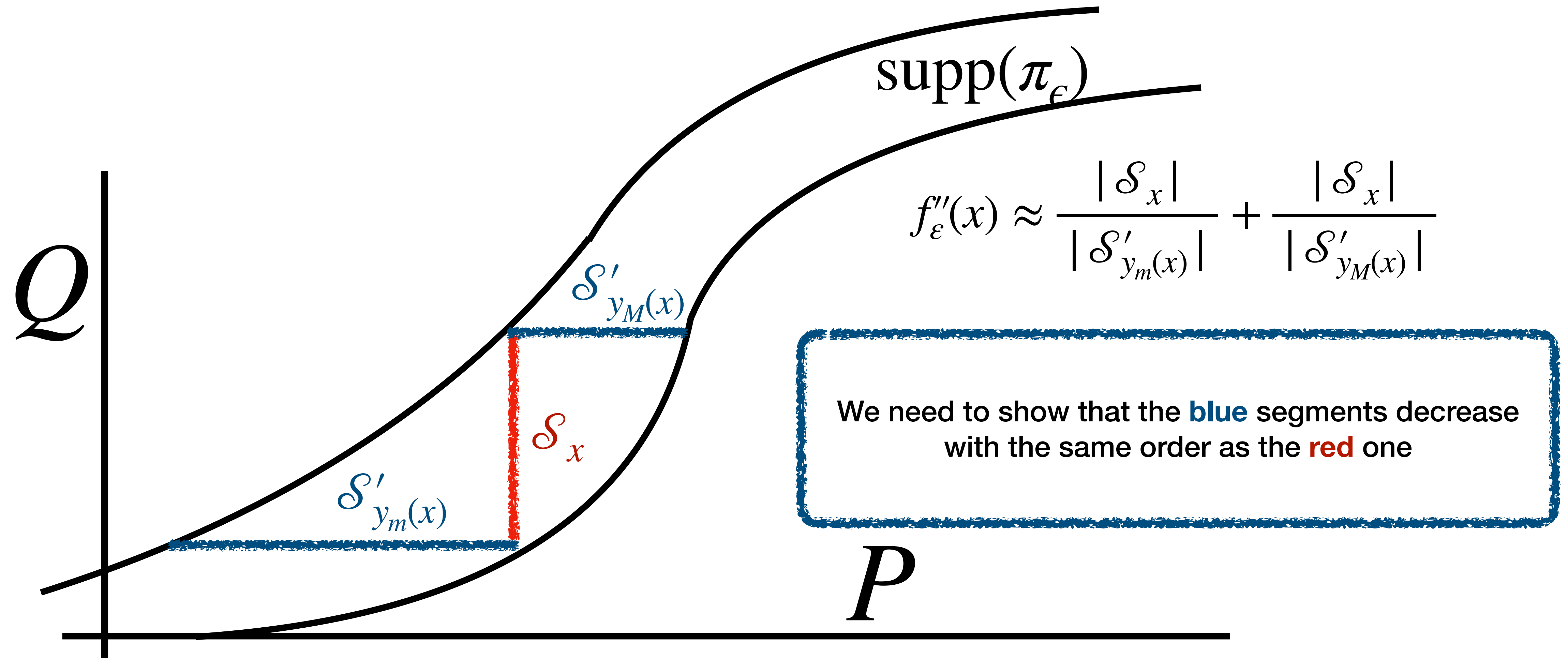
Step 2) A Bound on the derivative

$$\frac{1}{4} |\mathcal{S}'_{y_M(x)}| \leq |g'_{\varepsilon}(y_M(x)) - x| \leq |\mathcal{S}'_{y_M(x)}| \quad \frac{1}{4} |\mathcal{S}'_{y_m(x)}| \leq |g'_{\varepsilon}(y_m(x)) - x| \leq |\mathcal{S}'_{y_m(x)}|$$



$$f''_{\varepsilon}(x) \approx \frac{|\mathcal{S}_x|}{|\mathcal{S}'_{y_m(x)}|} + \frac{|\mathcal{S}_x|}{|\mathcal{S}'_{y_M(x)}|}$$

Bound on the second derivative



Bound on the second derivative

Call

$$\sigma_m(f_\varepsilon) := \inf_{x \in \Omega_0 \setminus \{x^{(m)}, x^{(M)}\}} f_\varepsilon''(x) > 0 \quad \sigma_M(f_\varepsilon) := \sup_{x \in \Omega_0 \setminus \{x^{(m)}, x^{(M)}\}} f_\varepsilon''(x) < +\infty.$$

Step 3) A Bound the derivative w.r.t. the maximum of the parabola

$$C^{-1}(\sigma_M(f_\varepsilon))^{-1/2} \max_{y \in [a_1, b_1]} (\xi(x, y))_+^{3/2} \leq \varepsilon \leq C(\sigma_m(f_\varepsilon))^{-1/2} \max_{y \in [a_1, b_1]} (\xi(x, y))_+^{3/2}.$$

$$\text{where } \xi(x, y) = \langle x, y \rangle - f_\varepsilon(x) - g_\varepsilon(y)$$

$$\text{and use } |\mathcal{S}_x| \max_{y \in [a_1, b_1]} \xi(x, y) \approx \varepsilon$$

$$\implies C^{-1} \left(\frac{\varepsilon}{\sigma_M(f_\varepsilon)} \right)^{\frac{1}{3}} \leq |\mathcal{S}_x| \leq C \left(\frac{\varepsilon}{\sigma_m(f_\varepsilon)} \right)^{\frac{1}{3}}.$$

Bound on the second derivative

Step 4) Use the estimates

$$C^{-1} \left(\frac{\varepsilon}{\sigma_M(f_\varepsilon)} \right)^{\frac{1}{3}} \leq |\mathcal{S}_x| \leq C \left(\frac{\varepsilon}{\sigma_m(f_\varepsilon)} \right)^{\frac{1}{3}} \quad \text{and} \quad f_\varepsilon''(x) \approx \frac{|\mathcal{S}_x|}{|\mathcal{S}'_{y_m(x)}|} + \frac{|\mathcal{S}_x|}{|\mathcal{S}'_{y_M(x)}|}$$

to get

$$\sigma_M(f_\varepsilon) \leq C \left(\frac{\sigma_M(g_\varepsilon)}{\sigma_m(f_\varepsilon)} \right)^{\frac{1}{3}} \quad \text{and} \quad \sigma_m(f_\varepsilon) \geq \frac{1}{C} \left(\frac{\sigma_m(g_\varepsilon)}{\sigma_M(f_\varepsilon)} \right)^{\frac{1}{3}}.$$

$$\sigma_M(g_\varepsilon) \leq C \left(\frac{\sigma_M(f_\varepsilon)}{\sigma_m(g_\varepsilon)} \right)^{\frac{1}{3}} \quad \text{and} \quad \sigma_m(g_\varepsilon) \geq \frac{1}{C} \left(\frac{\sigma_m(f_\varepsilon)}{\sigma_M(g_\varepsilon)} \right)^{\frac{1}{3}}$$

Bound on the second derivative

Step 5) Use

$$\sigma_M(f_\varepsilon) \leq C \left(\frac{\sigma_M(g_\varepsilon)}{\sigma_m(f_\varepsilon)} \right)^{\frac{1}{3}} \quad \text{and} \quad \sigma_m(f_\varepsilon) \geq \frac{1}{C} \left(\frac{\sigma_m(g_\varepsilon)}{\sigma_M(f_\varepsilon)} \right)^{\frac{1}{3}} .$$
$$\sigma_M(g_\varepsilon) \leq C \left(\frac{\sigma_M(f_\varepsilon)}{\sigma_m(g_\varepsilon)} \right)^{\frac{1}{3}} \quad \text{and} \quad \sigma_m(g_\varepsilon) \geq \frac{1}{C} \left(\frac{\sigma_m(f_\varepsilon)}{\sigma_M(g_\varepsilon)} \right)^{\frac{1}{3}}$$

to get

$$\sigma_M(f_\varepsilon) \leq C \sigma_M(f_\varepsilon)^{\frac{4}{49}} \quad \text{and} \quad \sigma_m(g_\varepsilon) \geq C^{-1} \sigma_m(g_\varepsilon)^{\frac{4}{49}}$$

Which yields the bound on the derivative

Bound on the second derivative

Theorem (GS, Nutz, 2024)

If dimension=1, if $P = p\mathbf{1}_{[a_0, b_0]}dx$ and $Q = q\mathbf{1}_{[a_1, b_1]}dx$ with p, q bounded away from zero and infinity,

- then f_ϵ is \mathcal{C}^2 in (a_0, b_0) except at two points,
- and there exists a constant $C > 0$ such that $C^{-1} \leq f''_\epsilon(x) \leq C$ for all $x \in \text{dom}(f''_\epsilon)$

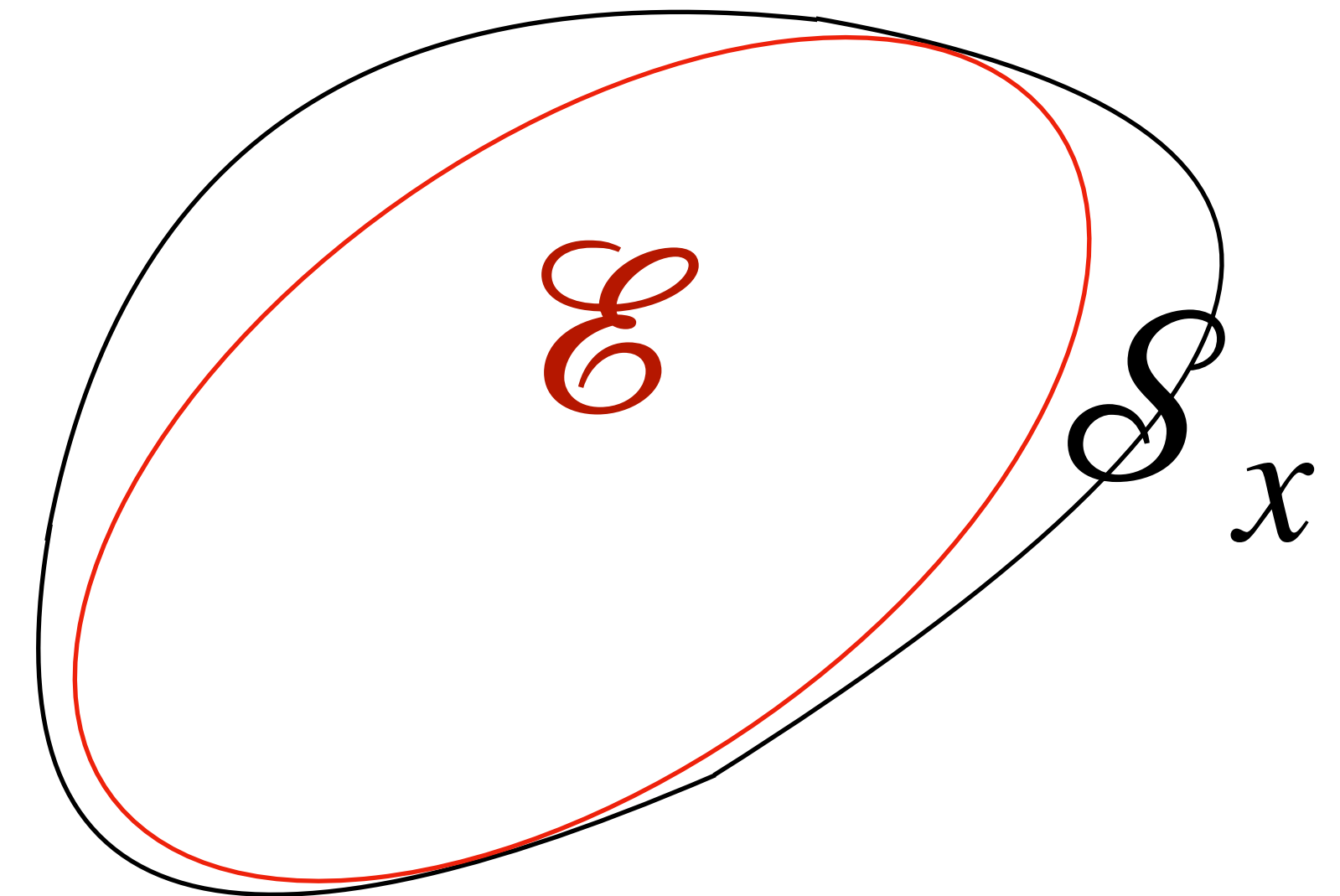
Corollary (GS, Nutz, 2024)

If dimension=1, if $P = p\mathbf{1}_{[a_0, b_0]}dx$ and $Q = q\mathbf{1}_{[a_1, b_1]}dx$ with p, q bounded away from zero and infinity, there exists a constant $C > 0$ such that

$$C^{-1}\epsilon^{\frac{1}{3}} \leq |\mathcal{S}_x| \leq C\epsilon^{\frac{1}{3}}, \quad \text{for all } x \in [a_0, b_0]$$

Difficulties on general dimension

- The sections are convex sets
- We can always introduce an ellipsoid \mathcal{E} of maximal volume (John ellipsoid) inside each section.
- To imitate the arguments of the 1D case, we need to ensure that the ellipsoid behaves like a ball. That is, the eigenvalues of the matrix defining the ellipsoid decrease to zero with the same order of convergence.



Sharp rates for the self-transport

Theorem (Wiesel, Xu, 2024)

If $P = Q = p\mathbf{1}_{\Omega_0}dx$ then

$$C^{-1}\epsilon^{\frac{1}{d+2}} \leq \text{diam}(\mathcal{S}_x) \leq C\epsilon^{\frac{1}{d+2}}, \quad \text{for all } x \in \Omega_0$$

If $P = p\mathbf{1}_{\Omega_0}dx$ and $Q = q\mathbf{1}_{\Omega_1}dx$ then

$$\text{diam}(\mathcal{S}_x) \leq C\epsilon^{\frac{1}{4(d+1)^2}}, \quad \text{for all } x \in \Omega_0$$

- The rates are sharp for the self-transport case (where the symmetry facilitates things a lot)
- The general case is far from the conjectured rate of $\epsilon^{1/(d+2)}$

Behaviour of QOT cost

We want to find the rate of convergence of the difference between the QOT and OT costs

$$\text{QOT}_\epsilon(P, Q) - \text{OT}(P, Q)$$

where

$$\text{QOT}_\epsilon(P, Q) = \min_{\pi \in \Pi(P, Q)} \int_{\mathbb{R}^d \times \mathbb{R}^d} \frac{1}{2} \|x - y\|^2 d\pi(x, y) + \frac{\epsilon}{2} \left\| \frac{d\pi}{d(P \times Q)} \right\|_{L^2(P \times Q)}^2,$$

$$\text{OT}(P, Q) = \min_{\pi \in \Pi(P, Q)} \int_{\mathbb{R}^d \times \mathbb{R}^d} \frac{1}{2} \|x - y\|^2 d\pi(x, y)$$

Behaviour of QOT cost

Theorem (Eckstein, Nutz, 2024)

If $P = p\mathbf{1}_{\Omega_0}dx$ and $Q = q\mathbf{1}_{\Omega_1}dx$ then

$$C^{-1} \leq \epsilon^{\frac{2}{d+2}}(\text{QOT}_\epsilon(P, Q) - \text{OT}(P, Q)) \leq C$$

- The proof is based on a quantization argument and an approximation by shadows instead to the block approximation of [Carlier et al. \(2017\)](#)
- The result holds for more general regularizations of OT
- In EOT the rate is

$$C^{-1} \leq \epsilon \log(\epsilon^{-1})(\text{EOT}_\epsilon(P, Q) - \text{OT}(P, Q)) \leq C$$

First order development

We want to find the exact limit

$$\lim_{\epsilon \rightarrow 0} \epsilon^{\frac{2}{d+2}} (\text{QOT}_\epsilon(P, Q) - \text{OT}(P, Q)) = ?$$

Strategy:

$$\liminf_{\epsilon \rightarrow 0} \epsilon^{\frac{2}{d+2}} (\text{QOT}_\epsilon(P, Q) - \text{OT}(P, Q)) \geq L$$

Lower bound (use dual QOT)

$$\limsup_{\epsilon \rightarrow 0} \epsilon^{\frac{2}{d+2}} (\text{QOT}_\epsilon(P, Q) - \text{OT}(P, Q)) \leq L$$

Upper bound (use primal QOT)

Lower bound

Strategy:

$$\liminf_{\epsilon \rightarrow 0} \epsilon^{\frac{2}{d+2}} (\text{QOT}_\epsilon(P, Q) - \text{OT}(P, Q)) \geq L \quad \text{Lower bound (use dual QOT)}$$

$$\Gamma(a, b) = \int a(x) dP(x) + \int b(y) dQ(y) - \frac{1}{2\epsilon} \left(a(x) + b(y) - \frac{\epsilon}{2} \|x - y\|^2 \right)_+^2 d(P \times Q)(x, y)$$

Since $\text{QOT}_\epsilon(P, Q) \geq \Gamma(a, b) \implies$ we need to find a correct candidate

$$C_\epsilon(x) := \frac{\epsilon^{\frac{2}{d+2}}}{C_d^{\frac{2}{d+2}} (p(x)q[\nabla f(x)])^{\frac{1}{d+2}}} \quad \text{for} \quad C_d := 2^{\frac{d+2}{2}} \mathcal{H}^{d-1}(\mathcal{S}^{d-1}) \frac{1}{d(d+2)}.$$

Lower bound

Since $\text{QOT}_\epsilon(P, Q) \geq \Gamma(a, b) \implies$ we need to find a correct candidate $(\tilde{a}_\epsilon, \tilde{b}_\epsilon)$

$$\tilde{a}_\epsilon(x) = f_0(x) + C_\epsilon(x), \quad C_\epsilon(x) := \frac{\epsilon^{\frac{2}{d+2}}}{C_d^{\frac{2}{d+2}} (p(x)q[T_{P \rightarrow Q}(x)])^{\frac{1}{d+2}}} \quad \text{for} \quad C_d := 2^{\frac{d+2}{2}} \mathcal{H}^{d-1}(\mathcal{S}^{d-1}) \frac{1}{d(d+2)}.$$

$$\tilde{b}_\epsilon(y) = g_0(y)$$

where (f_0, g_0) solves Dual OT and $T_{P \rightarrow Q}$ is the OT map from P to Q

$$\Gamma(\tilde{a}_\epsilon, \tilde{b}_\epsilon) = \text{OT}(P, Q) + \epsilon^{\frac{2}{d+2}} \frac{d^{\frac{d+4}{d+2}} (d+2)^{\frac{2}{d+2}}}{(\mathcal{H}^{d-1}(\mathcal{S}^{d-1}))^{\frac{2}{d+2}}} \int_{\Omega_0} (p(x)q[T_{P \rightarrow Q}(x)])^{-\frac{1}{d+2}} dP(x) + o(\epsilon^{\frac{2}{d+2}})$$

Upper bound

Strategy:

$$\limsup_{\epsilon \rightarrow 0} \epsilon^{\frac{2}{d+2}} (\text{QOT}_\epsilon(P, Q) - \text{OT}(P, Q)) \leq L \quad \text{Upper bound (use primal QOT)}$$

We need to find $\tilde{\pi}_\epsilon \in \Pi(P, Q)$ such that the functional

$$\Theta(\pi) = \int_{\mathbb{R}^d \times \mathbb{R}^d} \frac{1}{2} \|x - y\|^2 d\pi(x, y) + \frac{\epsilon}{2} \left\| \frac{d\pi}{d(P \times Q)} \right\|_{L^2(P \times Q)}^2,$$

evaluated at $\tilde{\pi}_\epsilon$ achieves the correct limit.

Upper bound

There are several strategies to find this candidate. For instance, in EOT:

- Block approximation ([Carlier et al., 2017](#))
- Gaussian approximation (or heat diffusion) ([Pal, 2019](#))
- Shadows ([Eckstein, Nutz, 2024](#))

For EOT Gaussian approximation is the correct approach because

- The Wasserstein gradient flow of the logarithmic entropy describes the flow of the heat equation ([Otto, 2001](#))
- The EOT potentials are created via iterative Gaussian convolutions

Upper bound

Porous media equation

$$\begin{cases} \partial_t u(t, x) = \Delta u(t, x)^2, & t > 0, x \in \mathbb{R}^d \\ u(0, x) = u_0(x) & x \in \mathbb{R}^d, \end{cases}$$

Fundamental solution (Barenblatt–Pattle)

$$\mathcal{B}(t, x) = \frac{1}{t^{\frac{d}{d+2}}} \left[C - \beta \frac{1}{4} \frac{\|x\|^2}{t^{\frac{2}{d+2}}} \right]_+,$$

For QOT a Barenblatt–Pattle type approximation is the correct approach because

- The Wasserstein gradient flow of the quadratic entropy describes the flow of the porous media equation ([Otto, 2001](#))
- The QOT potentials are created via an iterative modification of a Barenblatt–Pattle profile to create a coupling

Upper bound

Send Q to P via OT map and Find a nice coupling in $\Pi(P, P)$

Candidate
$$v(t, x; x') = \frac{1}{\epsilon} \left(\frac{C_d \epsilon^{\frac{2}{d+2}}}{(p(x')q(T_{P \rightarrow Q}(x')))^{\frac{1}{d+2}}} - \frac{1}{2} \|x - x'\|_{DT_{P \rightarrow Q}(x')}^2 \right)_+,$$

$$\|x - x'\|_{DT_{P \rightarrow Q}(x')}^2 = \langle x - x', DT_{P \rightarrow Q}(x')(x - x') \rangle,$$

This is not a
feasible coupling

$$v(t, x; x') := u \left(t, \left[T_{P \rightarrow Q}(x') \right]^{\frac{1}{2}} x; \left[T_{P \rightarrow Q}(x') \right]^{\frac{1}{2}} x \right) \begin{cases} \partial_t u(t, x; x') = \frac{1}{2(d+2)} \Delta_x u(t, x; x')^2, & t > 0, x \in \Omega_0 \\ u(0, x; x') = p(x')^{-\frac{1}{d+2}} \cdot \delta_{x'}(x) & x \in \mathbb{R}^d \end{cases}$$

Upper bound

Send Q to P via OT map and Find a nice coupling in $\Pi(P, P)$

Candidate

$$\nu(t, x; x') = \frac{1}{\epsilon} \left(\frac{C_d \epsilon^{\frac{2}{d+2}}}{(p(x')q(T_{P \rightarrow Q}(x')))^{\frac{1}{d+2}}} - \frac{1}{2} \|x - x'\|_{DT_{P \rightarrow Q}(x')}^2 \right)$$

$$\|x - x'\|_{DT_{P \rightarrow Q}(x')}^2 = \langle x - x', DT_{P \rightarrow Q}(x')(x - x') \rangle,$$

This is not a
feasible coupling

To make ν a feasible coupling:

- Normalize in order that it is a probability measure
- Send both marginals to P via the OT map
- Control the errors using Caffarelli's interior regularity theory

Upper bound

Theorem (Garriz Molina, GS, Mordant, 2024)

Assume $P = p\mathbf{1}_{\Omega_0}dx$ and $Q = q\mathbf{1}_{\Omega_1}dx$ with density bounded away from zero and infinity and supports with Lipschitz boundary. Then

$$\text{QOT}_\epsilon(P, Q) = \text{OT}(P, Q) + \epsilon^{\frac{2}{d+2}} \frac{d^{\frac{d+4}{d+2}}(d+2)^{\frac{2}{d+2}}}{(\mathcal{H}^{d-1}(\mathcal{S}^{d-1}))^{\frac{2}{d+2}}} \int_{\Omega_0} (p(x)q[T_{P \rightarrow Q}(x)])^{-\frac{1}{d+2}} dP(x) + o(\epsilon^{\frac{2}{d+2}})$$

Conclusions

- QOT represents a sparse alternative of EOT.
- In the discrete case the convergence is stationary.
- The contraction of the support is not monotone in general.
- In one dimension and the self-transport cases the rates of convergence of the support are $\epsilon^{1/(d+2)}$.
- The rate of the cost is $\epsilon^{2/(d+2)}$ and the first order limit is obtained by an approximation of the solutions via a modification of the fundamental solution of the Porous Media equation.

Open questions

- Sharp rates of convergence of the support in general dimension
- First order developments of the cost for other penalisations of OT
- Is there a variational formulation of QOT?
- Does a PL inequality hold?
- Statistical complexity of QOT (rates of convergence from the empirical to the population QOT)