

Generalised Wasserstein Barycentres

Eloi Tanguy

A thesis presented for the degree of
Masters in Applied Mathematics

Under the supervision of:

Julie Delon, MAP5 - Université Paris-Cité;
Rémi Flamary, CMAP - Ecole Polytechnique.



Laboratoire MAP5
Université Paris-Cité
April-September 2022

Table of Contents

1	Abstract	2
2	Introduction to Generalised Wasserstein Barycentres	3
2.1	Notations	3
2.2	Reminders on Wasserstein Barycentres	3
2.3	The Generalised Barycentre Problem	4
2.4	Examples	5
3	Numerical Optimisation for GWB	8
3.1	Abstract Algorithms	8
3.2	Solving GWB with Gradient Descent	9
3.3	Solving GWB with Block Coordinate Descent	13
4	A particular GWB: the Reconstruction Problem	16
4.1	Discrete Reconstruction Theory	16
4.2	Consequences on Sliced Wasserstein Methods	24
4.3	Experimental Results on the Reconstruction Problem	25
4.4	Local Optima: $L = 2$ case	30
4.5	Local Optima: general case	39
5	The Blind Generalised Wasserstein Barycentre Problem	41
5.1	Problem Description	41
5.2	Theoretical Properties of BGWB	41
5.3	Extending the Gradient Descent Solvers	42
5.4	BCD resolution	45
6	Perspectives and Conclusion	47
6.1	Perspectives	47
6.2	Conclusion	47

1 Abstract

This work is the continuation of a previous article [11] by Julie Delon, Nathaël Gozlan, and Alexandre Saint-Dizier, which considers a generalised version of Wasserstein Barycentres - the latter were introduced by Martial Agueh and Guillaume Carlier [3]. Our objective is to further investigate the properties of this generalised Wasserstein Barycentre, to develop numerical solvers adapted to this new problem, and finally to introduce further generalisations.

To this end, we begin by reminding the context of generalised Wasserstein Barycentres in Section 2, as well as introduce the notations for this report, and provide some first visual experiments for our considered problems.

In Section 3, we introduce multiple numerical solvers for the generalised barycentre problem, which can be accessed at our [project's repository](#). Note that the gradient-descent based solver is an adaptation of the `free_support_barycentre` method [9] that we made [available in POT](#) [14]. The other solvers will also come to POT shortly.

Having noticed the connexions between Sliced Optimal Transport and generalised Wasserstein Barycentres in the case of projections onto lines, we considered a specific sub-problem of the GWB problem, the Reconstruction Problem. In this setting, computing a barycentre amounts to finding a reconstruction of a measure based on its projections by linear maps. In Section 4, we study this reconstruction problem and deduce some properties of generalised Wasserstein Barycentres and the Sliced Wasserstein distance in particular cases.

Finally, in Section 5, we introduce a further generalisation of GWB where the linear maps P_i are no longer inputs but unknown variables. We provide a summary analysis of the new difficulties of this problem and provide numerical solvers.

2 Introduction to Generalised Wasserstein Barycentres

2.1 Notations

- We denote $\mathcal{P}_2(\mathbb{R}^k)$ the set of Radon measures on \mathbb{R}^k admitting a second-order moment.
- Σ_n will denote the n -dimensional simplex: $\Sigma_n := \left\{ a \in \mathbb{R}_+^n : \sum_{i=1}^n a_i = 1 \right\}$.
- We consider p probability measures $\nu_1 \cdots, \nu_p$ respectively in $\mathcal{P}(\mathbb{R}^{d_1}) \cdots \mathcal{P}(\mathbb{R}^{d_p})$, as well as p weights $\lambda = (\lambda_1, \cdots, \lambda_p)^T$ such that $\lambda \in \Sigma_p$. In general we will assume that the $\lambda_i \neq 0$.
- Let $\Pi(\nu_1, \dots, \nu_p)$ the subset of measures in $\mathcal{P}_2(\mathbb{R}^{d_1} \times \cdots \times \mathbb{R}^{d_p})$ with marginals ν_1, \dots, ν_p .
- The push-forward measure of $T : \mathbb{R}^d \rightarrow \mathbb{R}^k$ on a measure μ on \mathbb{R}^d is defined such that for all borelians $B \subset \mathbb{R}^k$, $T\#\mu(B) = \mu(T^{\circ-}(B))$, where $T^{\circ-}(B)$ is the reciprocal image of B by the map T .
- In general, the index i will denote the i -th measure within the p input measures used to compute a barycentre.
- $\mathcal{M}_{m,n}(\Omega)$ will denote the set of $m \times n$ matrices with entries in Ω .
- $S_d(\mathbb{R})$, $S_d^+(\mathbb{R})$, $S_d^{++}(\mathbb{R})$ are respectively the symmetric, positive semi-definite and positive $d \times d$ matrices with real-valued entries.
- $\mathbf{1}$ will be a vector full of ones and \mathbb{S}^d the unit sphere for $\|\cdot\|_2$ in \mathbb{R}^d .
- \mathfrak{S}_n will denote the group of permutations of $\llbracket 1, n \rrbracket$.
- \mathbb{E} and \mathbb{V} denote respectively the expectancy and variance of a random variable (provided they exist).

2.2 Reminders on Wasserstein Barycentres

Definition 2.2.1 — Wasserstein Distance [17]

The 2-Wasserstein distance between two measures $\nu_1, \nu_2 \in \mathcal{P}_2(\mathbb{R}^d)$ is:

$$W_2^2(\nu_1, \nu_2) := \inf_{\pi \in \Pi(\nu_1, \nu_2)} \int_{\mathbb{R}^d \times \mathbb{R}^d} \|x_1 - x_2\|^2 d\pi(x_1, x_2). \quad (1)$$

The classical barycentre, as introduced in [8] then [3], is defined between measures in \mathbb{R}^d ($d_i = d$ here):

Definition 2.2.2 — Wasserstein Barycentre [3]

Given p points in \mathbb{R}^d , let $B_\lambda(x_1, \dots, x_p) := \sum_{i=1}^p \lambda_i x_i = \operatorname{argmin}_{y \in \mathbb{R}^d} \sum_{i=1}^p \lambda_i \|x_i - y\|^2$ (2).

The Wasserstein Barycentre problem is a Fréchet barycentre for W_2^2 :

$$\operatorname{argmin}_{\nu \in \mathcal{P}_2(\mathbb{R}^d)} \sum_{i=1}^p \lambda_i W_2^2(\nu_i, \nu). \quad (\text{WB})$$

Definition 2.2.3 — Multi-Marginal Wasserstein Barycentre [4]

Given $\nu_1, \dots, \nu_p \in \mathcal{P}_2(\mathbb{R}^d)$, we call MMWB of ν_1, \dots, ν_p any solution of:

$$\operatorname{argmin}_{\pi \in \Pi(\nu_1, \dots, \nu_p)} \int_{\mathbb{R}^d \times \dots \times \mathbb{R}^d} \sum_{i=1}^p \|x_i - B_\lambda(x_1, \dots, x_p)\|^2 d\pi(x_1, \dots, x_p). \quad (\text{MMWB})$$

As proven in [3], both formulations are equivalent:

Theorem 2.2.4 — WB-MMWB correspondence [3]

If π^* is a solution of (MMWB), then $\nu^* := B_\lambda \# \pi^*$ is a solution of (WB).

If at least one of the measures ν_i admits a density with respect to the Lebesgue measure, then both problems (WB) and (MMWB) have a unique solution.

2.3 The Generalised Barycentre Problem

Returning to the general case, given p linear applications $P_i \in \mathcal{M}_{d,d_i}(\mathbb{R})$, we present the generalised barycentre problem [11].

Let $A := \sum_{i=1}^p \lambda_i P_i^T P_i$ which we assume invertible (thus symmetric positive definite).

Problem 2.3.1 — Generalised Barycentre [11]

Given p points $x_i \in \mathbb{R}^{d_i}$, their generalised barycentre is

$$B_{P,\lambda}(x_1, \dots, x_p) = \operatorname{argmin}_{x \in \mathbb{R}^d} \sum_{i=1}^p \lambda_i \|P_i x - x_i\|^2 = A^{-1} \left(\sum_{i=1}^p \lambda_i P_i^T x_i \right). \quad (3)$$

$$\text{Let } F : \begin{cases} \mathcal{P}_2(\mathbb{R}^d) & \longrightarrow & \mathbb{R} \\ \gamma & \longmapsto & \sum_{i=1}^p \lambda_i W_2^2(\nu_i, P_i \# \gamma) \end{cases},$$

we call GWB of ν_1, \dots, ν_p with weights $\lambda_1, \dots, \lambda_p$ any solution of the problem:

$$\inf_{\gamma \in \mathcal{P}_2(\mathbb{R}^d)} F(\gamma). \quad (\text{GWB})$$

As proven by [11] §3.1, (GWB) can be re-written as a (WB)-type problem:

Theorem 2.3.2 — GWB reformulation [11]

$$\text{Let for } i \in \llbracket 1, p \rrbracket, \quad \tilde{\nu}_i := (A^{-1/2} P_i^T) \# \nu_i \text{ and } G : \begin{cases} \mathcal{P}_2(\mathbb{R}^d) & \longrightarrow & \mathbb{R} \\ \mu & \longmapsto & \sum_{i=1}^p \lambda_i W_2^2(\tilde{\nu}_i, \mu) \end{cases}.$$

There exists a constant $C \in \mathbb{R}$ such that for a given $\gamma \in \mathcal{P}_2(\mathbb{R}^d)$ and $\mu := A^{1/2} \# \gamma$, we have $F(\gamma) = G(\mu) + C$, yielding an equivalent formulation of (GWB):

$$\inf_{\mu \in \mathcal{P}_2(\mathbb{R}^d)} \sum_{i=1}^p \lambda_i W_2^2(\tilde{\nu}_i, \mu). \quad (\text{GWB}')$$

This reformulation amounts to a classical Wasserstein barycentre problem, which always admits a solution [3]. If the measures ν_i admit a density, then we do not have unicity, however if γ_1 and γ_2 are

solutions of (GWB) then $\forall i \in \llbracket 1, p \rrbracket$, $P_i \# \gamma_1 = P_i \# \gamma_2$ [11].

Note that in the discrete case, we do not have unicity for the classical or generalised versions of the Wasserstein barycentre.

Problem 2.3.3 — Multi-marginal formulation [11]

Like the classical barycentre problem, (GWB) can be re-written as a multi-marginal problem:

$$\operatorname{argmin}_{\pi \in \Pi(\nu_1, \dots, \nu_p)} \int_{\mathbb{R}^{d_1} \times \dots \times \mathbb{R}^{d_p}} \sum_{i=1}^p \lambda_i \|x_i - P_i B_{P, \lambda}(x_1, \dots, x_p)\|^2 d\pi(x_1, \dots, x_p). \quad (\text{MMGWB})$$

π^* is a solution of (MMGWB) iff $\gamma^* := B_{P, \lambda} \# \pi^*$ is a solution of (GWB).

Again this problem can be re-written as a multi-marginal classical barycentre problem:

$$\operatorname{argmin}_{\pi \in \Pi(\tilde{\nu}_1, \dots, \tilde{\nu}_p)} \int_{\mathbb{R}^{d_1} \times \dots \times \mathbb{R}^{d_p}} \sum_{i=1}^p \lambda_i \|x_i - B_\lambda(x_1, \dots, x_p)\|^2 d\pi(x_1, \dots, x_p). \quad (\text{MMGWB}')$$

π^* is a solution of (MMGWB') iff $\gamma^* := A^{-1/2} B_\lambda \# \pi^*$ is a solution of (GWB).

2.4 Examples

2.4.1 GWB between two measures

We consider two discrete 2D measures α and β which respectively draw an H and an O. We also consider the projections (which can be visualised as projecting onto two vertical sides of a cube):

$$P_1 = \begin{pmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \end{pmatrix}, \quad P_2 = \begin{pmatrix} 0 & 1 & 0 \\ 0 & 0 & 1 \end{pmatrix}.$$

In order to compute the barycentre, we reformulate the generalised barycentre problem as a classical barycentre problem. We let $A = \lambda_1 P_1^T P_1 + \lambda_2 P_2^T P_2$ and $\tilde{\alpha} = (A^{-1/2} P_1^T) \# \alpha$, $\tilde{\beta} = (A^{-1/2} P_2^T) \# \beta$.

The Generalised Barycentre Problem can be written as the following equivalent barycentre problem:

$$(\text{GWB}') : \quad \inf_{\mu \in \mathcal{P}_2(\mathbb{R}^3)} \lambda_1 W_2^2(\tilde{\alpha}, \mu) + \lambda_2 W_2^2(\tilde{\beta}, \mu).$$

Consider X, Y the respective supports of $\tilde{\alpha}, \tilde{\beta}$. Now we need to compute the barycentres between the two uniform measures $\tilde{\alpha} = \sum_{i=1}^n \frac{1}{n} \delta_{x_i}$ and $\tilde{\beta} = \sum_{j=1}^m \frac{1}{m} \delta_{y_j}$.

We compute π , the OT matrix solution to the discrete Kantorovitch problem:

$$\operatorname{argmin}_{\pi \in \mathcal{M}_{n, m}(\mathbb{R})} \sum_{i=1}^n \sum_{j=1}^m \pi_{ij} \|x_i - y_j\|^2, \quad \pi \mathbf{1} = \frac{1}{n} \mathbf{1}, \quad \pi^T \mathbf{1} = \frac{1}{m} \mathbf{1}, \quad \pi \geq 0$$

This can be done using `ot.emd` [14], from the Python OT library **POT**.

The associated transport plan is $\mu^* = \sum_{i=1}^n \sum_{j=1}^m \pi_{i,j} \delta_{x_i} \otimes \delta_{y_j}$.

Let $P_t(x, y) := (1-t)x + ty$. Finally, the $(1-t), t$ weighted barycentre is $P_t \# \mu^* = \sum_{i=1}^n \sum_{j=1}^m \pi_{i,j} \delta_{(1-t)x_i + ty_j}$.

Then our final desired transport plan is $\gamma^* = A^{-1/2} \# \mu^*$.

Now if we compare α with $P_1 \# \mu^*$ and β with $P_2 \# \mu^*$, the matching isn't perfect since the two measures α, β are not compatible, as shown on [Figure 1](#).

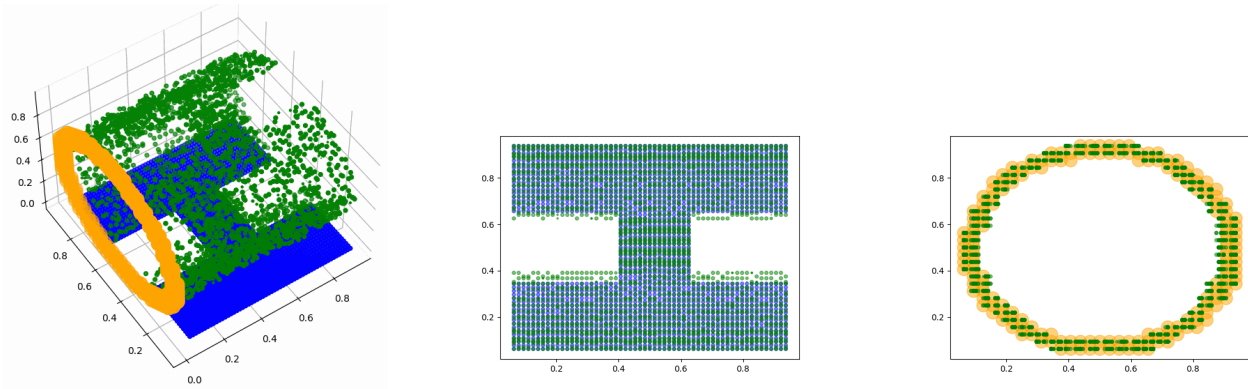


Figure 1: GWB resolution on a toy dataset with two 2D input measures.

2.4.2 GWB between three measures

We load 3 different logos as discrete measures and display them in the coupling space $\mathbb{R}^d = \mathbb{R}^3$.

Each logo is seen as a discrete uniform measure $\nu_i = \sum_{j=1}^{n_i} a_j^{(i)} \delta_{x_j^{(i)}}$ on \mathbb{R}^2 where $a_j^{(i)} = 1/n_i$, and we

consider the $P_i : \mathbb{R}^3 \rightarrow \mathbb{R}^2$ to be projections onto 3 faces of a cube.

Like the 2-measure example, we re-formulate the GWB problem as a Wasserstein Barycentre problem:

$$(GWB') : \inf_{\mu \in \mathcal{P}_2(\mathbb{R}^3)} \sum_{i=1}^3 \lambda_i W_2^2(\tilde{\nu}_i, \mu),$$

where the wanted optimal coupling is obtained using $\gamma^* = A^{-1/2} \# \mu^*$, and we let for $i \in [1, 3]$, $\tilde{\nu}_i := (A^{-1/2} P_i^T) \# \nu_i$. First, we solve this problem using the "free support barycentre" method [9] using `ot.lp.free_support_barycenter` [14], which optimises a fixed number of baycentre positions and leaves the barycentre weights as uniform throughout, as seen in Figure 2.

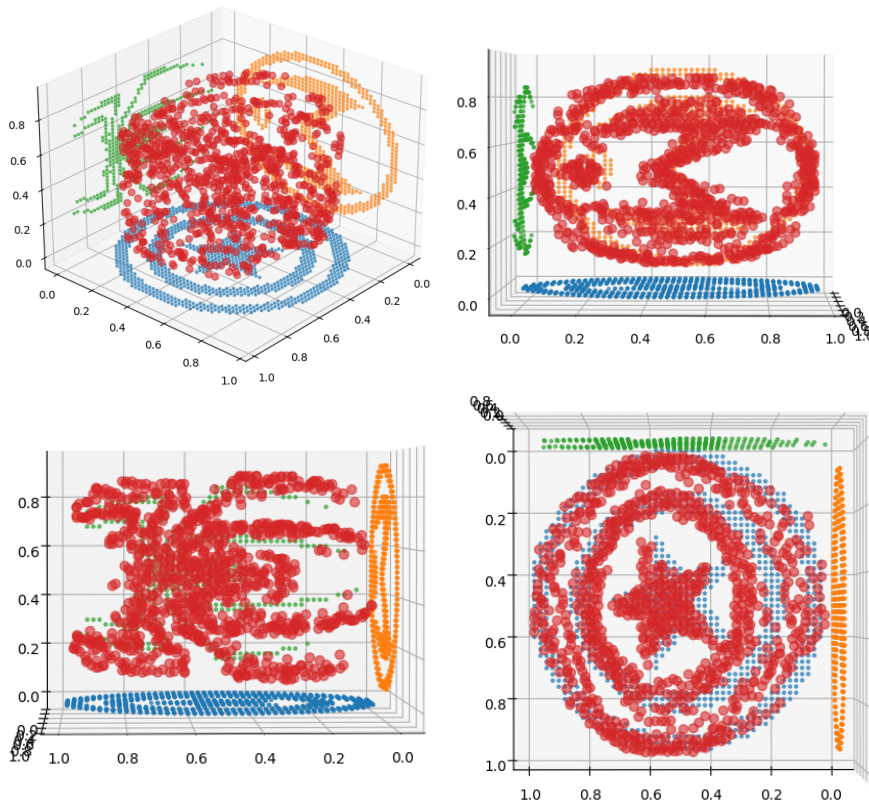


Figure 2: (GWB) resolution on a toy dataset with three 2D input measures.

Another way of solving this problem is the Multi-Marginal formulation:

$$(\text{MMGWB}') : \quad \underset{\pi \in \Pi(a^{(1)}, a^{(2)}, a^{(3)})}{\operatorname{argmin}} \quad \sum_{j_1, j_2, j_3} c_{j_1, j_2, j_3} \pi_{j_1, j_2, j_3},$$

where:

- $\Pi(a^{(1)}, a^{(2)}, a^{(3)})$ is the set of tensors of $\mathbb{R}_+^{n_1 \times n_2 \times n_3}$ with marginals $a^{(i)}$.
- The cost is $c_{j_1, j_2, j_3} = \sum_{i=1}^3 \lambda_i \left\| x_{j_i} - B_\lambda \left(x_{j_1}^{(1)}, x_{j_2}^{(2)}, x_{j_3}^{(3)} \right) \right\|^2$.
- $B_\lambda(x_1, x_2, x_3) := \lambda_1 x_1 + \lambda_2 x_2 + \lambda_3 x_3$ is the euclidian barycentre.

The associated measure is $\pi := \sum_{j_1, j_2, j_3} \pi_{j_1, j_2, j_3} \delta_{x_{j_1}^{(1)}} \otimes \delta_{x_{j_2}^{(2)}} \otimes \delta_{x_{j_3}^{(3)}}$ and $\gamma^* = A^{-1/2} B_\lambda \# \pi^*$ solves (GWB).

While very elegant mathematically, this method is prohibitively expensive in practice. Indeed, the multi-marginal problem is more general than the barycentre problem which is already extremely expensive [4].

3 Numerical Optimisation for GWB

3.1 Abstract Algorithms

In order to solve our problem (GWB), we consider different classical non-convex problems solvers.

We will present the algorithms in the following framework: minimising $f(x_1, \dots, x_n)$ for $(x_1, \dots, x_n) \in X_1 \times \dots \times X_p$, where each X_i is a (convex) subset of \mathbb{R}^{d_i} . Let Π_{X_i} the orthogonal projection onto $X_i \subset \mathbb{R}^{d_i}$.

Algorithm 1: Projected Alternated Gradient Descent

Data: Number of steps N , precision ε , learning rate η , l.r. decay ρ .

Result: $(x_1, \dots, x_n) \in X_1 \times \dots \times X_n$ minimising f .

```

1 Initialisation: Draw  $\forall i \in \llbracket 1, n \rrbracket$ , draw  $x_i^{(0)} \in X_i$ ;
2 for  $t \in \llbracket 1, N \rrbracket$  do
3    $\forall i \in \llbracket 1, n \rrbracket$ ,  $x_i^{(t)} = \Pi_{X_i} \left( x_i^{(t-1)} - \rho^t \eta \frac{\partial f}{\partial x_i} (x_1^{(t-1)}, \dots, x_n^{(t-1)}) \right)$ ;
4   if  $f(x_1^{(t-1)}, \dots, x_n^{(t-1)}) - f(x_1^{(t)}, \dots, x_n^{(t)}) < \varepsilon$  then
5     | Declare convergence and terminate.
6   end
7 end

```

An important variant is Algorithm 1's stochastic counterpart, where the indices i are chosen at random in batches of B following a discrete distribution \mathbb{D} on $\llbracket 1, n \rrbracket$.

Algorithm 2: Stochastic Projected Alternated Gradient Descent

Data: Number of steps N , precision ε , learning rate η , l.r. decay ρ , batch size B and sample law \mathbb{D} .

Result: $(x_1, \dots, x_n) \in X_1 \times \dots \times X_n$ minimising f .

```

1 Initialisation: Draw  $\forall i \in \llbracket 1, n \rrbracket$ , draw  $x_i^{(0)} \in X_i$ ;
2 for  $t \in \llbracket 1, N \rrbracket$  do
3   for  $\_ \in \llbracket 1, B \rrbracket$  do
4     | Draw  $i \sim \mathbb{D}$  and step:
4     |  $x_i^{(t)} = \Pi_{X_i} \left( x_i^{(t-1)} - \rho^t \eta \frac{\partial f}{\partial x_i} (x_1^{(t-1)}, \dots, x_n^{(t-1)}) \right)$ ;
5   end
6   if  $f(x_1^{(t-1)}, \dots, x_n^{(t-1)}) - f(x_1^{(t)}, \dots, x_n^{(t)}) < \varepsilon$  then
7     | Declare convergence and terminate.
8   end
9 end

```

Another classical solver is Block Coordinate Descent, which consists in optimising formally parameter by parameter, assuming that closed forms are known for the partial optimisations:

Algorithm 3: Block Coordinate Descent**Data:** Number of steps N , precision ε **Result:** $(x_1, \dots, x_n) \in X_1 \times \dots \times X_n$ minimising f .

```

1 Initialisation: Draw  $\forall i \in \llbracket 1, n \rrbracket$ , draw  $x_i^{(0)} \in X_i$ ;
2 for  $t \in \llbracket 1, N \rrbracket$  do
3    $\forall i \in \llbracket 1, n \rrbracket$ ,  $x_i^{(t)} = \operatorname{argmin}_{x_i \in X_i} f(x_1^{(t-1)}, \dots, x_{i-1}^{(t-1)}, x_i, x_{i+1}^{(t-1)}, \dots, x_n^{(t-1)})$ ;
4   if  $f(x_1^{(t-1)}, \dots, x_n^{(t-1)}) - f(x_1^{(t)}, \dots, x_n^{(t)}) < \varepsilon$  then
5     | Declare convergence and terminate.
6   end
7 end

```

For the stopping condition, it is also typical to use a stationarity criterion, for example:

$$\left\| (x_1^{(t-1)}, \dots, x_n^{(t-1)}) - (x_1^{(t)}, \dots, x_n^{(t)}) \right\| < \varepsilon$$

This condition is more natural for BCD than for GD and SGD, since BCD does not inherently require a computation of the energy f at each step.

3.2 Solving GWB with Gradient Descent

The idea is to solve the problem (GWB) by gradient descent (Algorithm 1 or Algorithm 2) on γ , where an arbitrary fixed number of points is considered.

We look for solutions $\gamma = \sum_{l=1}^L b_l \delta_{y_l}$, where L is fixed and $b \in \Sigma_L$ the simplex. This type of solution is similar to [9], which solves usual barycentres with gradient descent.

Let $i \in \llbracket 1, p \rrbracket$ and consider the discrete formulation of $J_i := W_2^2(\nu_i, P_i \# \gamma)$. For legibility we will temporarily drop the i indices. The computation of J is a Discrete Kantorovitch problem.

Writing $\nu = \sum_{k=1}^K a_k \delta_{x_k}$, $P \# \gamma = \sum_{l=1}^L b_l \delta_{P y_l}$, $M_{k,l} := \|x_k - P y_l\|_2^2$:

$$J = \min_{\substack{\pi \in \mathcal{M}_{K,L}(\mathbb{R}_+) \\ \pi \mathbf{1} = a \\ \pi^T \mathbf{1} = b}} M \cdot \pi = \max_{\substack{f \in \mathbb{R}^K, g \in \mathbb{R}^L \\ f \oplus g \leq M}} f \cdot a + g \cdot b, \quad (\text{DK})$$

where $f \oplus g = (f_l + g_k)_{(l,k) \in \llbracket 1, K \rrbracket \times \llbracket 1, L \rrbracket}$ and where the second equality is obtained through strong duality (in this admissible linear program), refer to [17] for details.

We are trying to compute the gradients of $J = J(b, Y)$, where $Y = (y_1, \dots, y_L)^T$. Note that this can be done via automatic differentiation, but it is important to consider the computations in order to understand the process.

3.2.1 Discussion on applying Shapiro's Theorem

We reproduce "Shapiro's Theorem" from [6]¹ for the sake of self-containedness:

Theorem 3.2.1 — Shapiro's Theorem [6]

Let X a metric space, U a normed space and $f : X \times U \rightarrow \mathbb{R}$ such that:

- $\forall x \in X$, $f(x, \cdot)$ is differentiable.
- $f(\cdot, \cdot)$ and the partial derivative $\partial_u f(\cdot, \cdot)$ are both continuous

Let $K \subset X$ compact and consider the optimum value $v(u) := \sup_{x \in K} f(x, u)$. We have:

- v is directionally differentiable
- If for $u_0 \in U$, $f(\cdot, u_0)$ has a unique minimiser x_0 on K then v is differentiable at u_0 , with $\boxed{dv(u_0) = \partial_u f(x_0, u_0)}$ (4)
- In general $D_h v(u_0) = \max_{x \in M(u_0)} \partial_u f(x, u_0) \cdot h$ (5) where $M(u_0) = \sup_{x \in K} f(x, u_0)$.

In the case of (DK) problem, let us check the hypotheses for the variables y_l, P_i, b :

- For the y_l , consider the primal problem $J(Y) = \min_{\substack{\pi \in \mathcal{M}_{K,L}(\mathbb{R}_+) \\ \pi \mathbf{1} = a \\ \pi^T \mathbf{1} = b}} M \cdot \pi$, where $Y := \begin{pmatrix} - & (y_1)^T & - \\ - & \vdots & - \\ - & (y_L)^T & - \end{pmatrix}$.

For $\pi \in \Pi(a, b)$, let $F(\pi, Y) := M \cdot \pi$. Since $M_{k,l} = \|x_k - P y_l\|_2^2$, the map $Y \mapsto F(\pi, Y)$ is differentiable (quadratic) and both maps $F(\cdot, \cdot)$ and $\partial_Y F(\cdot, \cdot)$ are continuous (polynomial).

Now we need to restrict π to a compact set. To that end, we can notice that each $\pi \in \Pi(a, b)$ satisfies $\|\pi\|_1 \leq 1$ ($\|A\|_1$ denoting here $\sum_{k,l} |A_{k,l}|$) since each $\pi_{k,l} \geq 0$ and $\sum_{k,l} \pi_{k,l} = \sum_k (\pi \mathbf{1})_k = \sum_k (a)_k = 1$. Thus $\Pi(a, b)$ is closed and bounded in a finite-dimensional space, thus compact.

- For b , we can directly use the sub-differentiability of $\alpha \mapsto W_2^2(\alpha, \beta)$, as is proven in detail in [18] §7.2.2, Proposition 7.17, which shows the desired result that a sub-gradient in b is given by an optimal dual potential for b .

Another difficulty is that we would want to apply the second point, which requires unicity. In theory, unicity is not guaranteed for a linear program such as (DK), and a possible counter example would be the following (Figure 3, from [17]), where the two possible transport plans are both optimal:

¹While [6] is generally cited, the original theorem is from [10] and proofs can be found in [5].

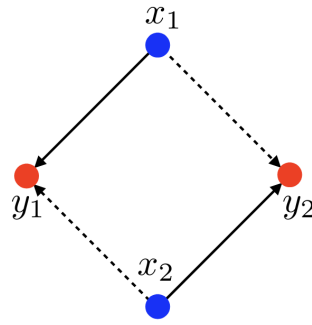


Figure 3: Example with two solutions to the (DK) problem.

In practice, there are several ways of dealing with this theoretical hurdle:

- One could apply the third point and compute the gradient as the maximum of all the found solutions (assuming there is a finite amount of them);
- One could replace (DK) problem with its Sinkhorn regularised version (see [17], §4.2), which is strictly convex, thus has unicity. This would be at the cost of the sparsity and exactitude of (DK) solution, but would provide the speed and stability advantages of Sinkhorn;
- If we keep an optimal solution instead of computing the maximum over optima, we obtain a sub-gradient (3.2.1, (3)) and can proceed with (projected) sub-gradient descent.

3.2.2 Gradient in b

Since $\nabla_b(f \cdot a + g \cdot b) = g$, by applying Santambrogio Proposition 7.17 [18] we obtain $\nabla_b J = g^*$ where (f^*, g^*) is a dual solution of (DK).

This allows optimisation on b using projected GD: $b_{t+1} = \Pi_{\Sigma_L}(b_t - \eta \nabla_b J(P_t, b_t, Y_t))$,

where $\Pi_{\Sigma_L}(v)$ is the projection on the simplex, computing $\operatorname{argmin}_{\substack{u \in \mathbb{R}_+^L \\ u \cdot \mathbf{1} = 1}} \|u - v\|_2$.

3.2.3 Gradient in M

Before computing the gradient in Y we need the gradient in M .

Using the primal formulation and Shapiro's theorem 3.2.1 we obtain $\nabla_M J = \pi^*$, where π^* is an optimal transport matrix for the (DK) problem.

3.2.4 Gradient in Y

We use mixed convention, where for $(l, \beta) \in \llbracket 1, L \rrbracket \times \llbracket 1, d \rrbracket$, $\frac{\partial J}{\partial M}$ and $\frac{\partial M}{\partial Y_{l,\beta}}$ have the shape (L, d) .

Using the chain rule² we can compute for $(l, \beta) \in \llbracket 1, L \rrbracket \times \llbracket 1, d \rrbracket$, $\frac{\partial J}{\partial Y_{l,\beta}} = \operatorname{Tr} \left(\left[\frac{\partial J}{\partial M} \right]^T \frac{\partial M}{\partial Y_{l,\beta}} \right)$,

where we compute $\frac{\partial M}{\partial Y_{l,\beta}}$ using $\frac{\partial M_{k,l}}{\partial y_{\nu}} = \delta_{l,\nu} (-2P^T x_k + 2P^T P y_l)$ (quadratic form in y_l).

²See *The Matrix Cookbook* [16], §2.8.1 (126)

3.2.5 GD for GWB Algorithm

We present our GD solver for the (GWB) problem. Note that technically, the method is Projected Alternated Gradient Descent.

Algorithm 4: (GWB) resolution with Gradient Descent

Data: Input measure points $(X_i)_{i \in \llbracket 1, p \rrbracket} \in \prod_{i=1}^p \mathcal{M}_{K_i, d_i}(\mathbb{R})$, weights $(a_i)_{i \in \llbracket 1, p \rrbracket} \in \prod_{i=1}^p \Sigma_{K_i}$,
 linear maps $(P_i)_{i \in \llbracket 1, p \rrbracket} \in \prod_{i=1}^p \mathbb{R}^{d_i}$, number of barycentre points L ,
 barycentric coefficients $\lambda \in \Sigma_p$, precision ε , iterations N , l.r. η and l.r. decay ρ .

Result: Barycentre positions $Y \in \mathcal{M}_{L, d}(\mathbb{R})$ and barycentre weights $b \in \Sigma_L$.

```

1 Initialisation: Draw  $Y \in \mathcal{M}_{L, d}(\mathbb{R})$  and  $b \in \Sigma_L$ ; let  $J_0 := +\infty$ ;
2 for  $t \in \llbracket 1, N \rrbracket$  do
3   for  $i \in \llbracket 1, p \rrbracket$  do
4     Compute  $J^{(i)} = \min_{\pi_i \in \Pi(a_i, b)} M_i \cdot \pi_i$  where  $M_{k, l}^{(i)} = \|x_k^{(i)} - P_i y_l\|_2^2$ 
5   end
6   Compute the loss  $J_t = \sum_{i=1}^p \lambda_i J^{(i)}$  and its gradients w.r.t.  $Y, b$ ;
7   Step the positions  $Y : Y \leftarrow Y - \rho^t \eta \frac{\partial J_t}{\partial Y}$ ;
8   Step the weights  $b : b \leftarrow \Pi_{\Sigma_L} \left( b - \rho^t \eta \frac{\partial J_t}{\partial b} \right)$ ;
9   if  $J_{t-1} - J_t < \varepsilon$  then
10    | Declare convergence and terminate.
11  end
12 end

```

The computation of the gradients at line 6 is facilitated by the use of the **PyTorch** back-end. Then the gradients of the Discrete Kantorovitch problem (DK) with respect to the distance matrix M and the weights b are computed by the `ot.emd2` function from **Python OT** [14]. PyTorch's **autograd** framework then computes automatically our desired gradients, thanks to the seamless backend integration done by the POT team.

3.2.6 SGD for BGWB Algorithm

Below is a stochastic variant of the GD solver for BGWB (Algorithm 8), which is technically Stochastic Projected Alternated Gradient Descent.

The modification is that instead of computing the complete loss J at every step, the loss is computed using p samples of a discrete probability distribution \mathbb{D} such that $\forall i \in \llbracket 1, p \rrbracket$, $\mathbb{D}(i) = \lambda_i$.

Algorithm 5: (GWB) resolution with Stochastic Gradient Descent

Data: Input measure points $(X_i)_{i \in \llbracket 1, p \rrbracket} \in \prod_{i=1}^p \mathcal{M}_{K_i, d_i}(\mathbb{R})$, weights $(a_i)_{i \in \llbracket 1, p \rrbracket} \in \prod_{i=1}^p \Sigma_{K_i}$,
 linear maps $(P_i)_{i \in \llbracket 1, p \rrbracket} \in \prod_{i=1}^p \mathbb{R}^{d_i}$, number of barycentre points L ,
 barycentric coefficients $\lambda \in \Sigma_p$, precision ε , iterations N , l.r. η and l.r. decay ρ .

Result: Barycentre positions $Y \in \mathcal{M}_{L, d}(\mathbb{R})$ and barycentre weights $b \in \Sigma_L$.

```

1 Initialisation: Draw  $Y \in \mathcal{M}_{L, d}(\mathbb{R})$  and  $b \in \Sigma_L$ ; let  $J_0 := +\infty$ ;
2 for  $t \in \llbracket 1, N \rrbracket$  do
3   for  $\_ \in \llbracket 1, p \rrbracket$  do
4     Draw  $i \sim \mathbb{D}$ ;
5     Compute  $J^{(i)} = \min_{\pi_i \in \Pi(a_i, b)} M_i \cdot \pi_i$  where  $M_{k, l}^{(i)} = \|x_k^{(i)} - P_i y_l\|_2^2$ 
6   end
7   Compute the loss  $J_t = \sum_{i=1}^p J^{(i)}$  and its gradients w.r.t.  $Y, b$ ;
8   Step the positions  $Y : Y \leftarrow Y - \rho^t \eta \frac{\partial J_t}{\partial Y}$ ;
9   Step the weights  $b : b \leftarrow \Pi_{\Sigma_L} \left( b - \rho^t \eta \frac{\partial J_t}{\partial b} \right)$ ;
10  if  $J_{t-1} - J_t < \varepsilon$  then
11    | Declare convergence and terminate.
12  end
13 end

```

Note that the SGD steps are done by batches of p , this allows better comparability to GD (Algorithm 4), in particular concerning the stopping criterion, which in the case of SGD consists in the average loss on the batch.

3.3 Solving GWB with Block Coordinate Descent

First of all let us define some notations. The matrices indexed in i (corresponding to the p input measures) will be denoted with an exponent (i) if we need to consider their entries.

We consider $\nu_i = \sum_{k=1}^{K_i} a_k^{(i)} \delta_{x_k^{(i)}}$, $\gamma = \sum_{l=1}^L b_l \delta_{y_l}$, where b_l is fixed (in practice $b_l = \frac{1}{L}$ is common).

Note that if we wanted to optimise in b , we could leverage 2.1.2 and do so by gradient descent, using the dual solution of the aforementioned (DK) problem in π . Sadly, there is no known closed form in b and thus no BCD optimisation in b .

This time we view the problem as minimising an energy $J(Y, (\pi^{(i)}))$ of variables $Y \in \mathcal{M}_{L, d}(\mathbb{R})$ and $\pi_i \in \Pi_i := \left\{ \pi \in \mathcal{M}_{K_i, L}(\mathbb{R}_+) : \pi \mathbf{1} = a_i, \pi^T \mathbf{1} = b \right\}$.

We define the energy J using the DK cost $J = \sum_{i=1}^p \lambda_i M_i \cdot \pi_i = \sum_{i=1}^p \lambda_i \sum_{(k, l) \in \llbracket 1, K_i \rrbracket \times \llbracket 1, L \rrbracket} \|x_k^{(i)} - P_i y_l\|_2^2 \pi_{k, l}^{(i)}$.

We will be optimising on each variable successively and independently, which leads us to computing the optimal values in closed form when possible.

In order to express the closed forms in matrix format, we define:

- For $i \in \llbracket 1, p \rrbracket$, the matrix of the K_i points of ν_i : $X_i := \begin{pmatrix} - & (x_1^{(i)})^T & - \\ - & \vdots & - \\ - & (x_{K_i}^{(i)})^T & - \end{pmatrix} \in \mathcal{M}_{K_i, d_i}(\mathbb{R})$
- The coupling measure points $Y := \begin{pmatrix} - & (y_1)^T & - \\ - & \vdots & - \\ - & (y_L)^T & - \end{pmatrix} \in \mathcal{M}_{L, d}(\mathbb{R})$

3.3.1 Closed form in Y

The structure of the problem allows us to optimise on each y_l separately. Let $l \in \llbracket 1, L \rrbracket$, we minimise the quadratic form:

$$J'(y_l) := \sum_{i=1}^p \lambda_i \sum_{k=1}^{K_i} \|x_k^{(i)} - P_i y_l\|_2^2 \pi_{k,l}^{(i)}.$$

By expanding the norm, we compute the gradient $\nabla_{y_l} J' = \sum_{i=1}^p \lambda_i \sum_{k=1}^{K_i} \pi_{k,l}^{(i)} (-2P_i^T x_k^{(i)} + P_i^T P_i y_l)$.

The minimising condition for the convex function J' then reads:

$$\left(\sum_{i=1}^p \lambda_i \sum_{k=1}^{K_i} \pi_{k,l}^{(i)} P_i^T P_i \right) y_l = \sum_{i=1}^p \lambda_i \sum_{k=1}^{K_i} \pi_{k,l}^{(i)} P_i^T x_k^{(i)}.$$

Which can be written $B_l y_l = c_l$, where due to the positivity of the coefficients $\lambda_i \pi_{k,l}^{(i)}$, the positive semi-definite matrix B_l is invertible iff $\sum_{i=1}^p \text{Im}(P_i) = \mathbb{R}^d$, which is a very weak assumption (see 4.1.5).

Finally we have a unique solution:

$$y_l^* = \left(\sum_{i=1}^p \lambda_i \sum_{k=1}^{K_i} \pi_{k,l}^{(i)} P_i^T P_i \right)^{-1} \left(\sum_{i=1}^p \lambda_i \sum_{k=1}^{K_i} \pi_{k,l}^{(i)} P_i^T x_k^{(i)} \right). \quad (6)$$

This expression can be expressed in matrix form with:

- $c_l = \sum_{i=1}^p \lambda_i \left(\pi_{\cdot, l}^{(i)} \right)^T X_i P_i \in \mathbb{R}^d$, where $\pi_{\cdot, l}^{(i)}$ is the l -th column of π_i .
- $B_l = \sum_{i=1}^p \lambda_i b_l P_i^T P_i$, using $\sum_{k=1}^{K_i} \pi_{k,l}^{(i)} = b_l$ since $\pi_i \in \Pi(a_i, b)$.

3.3.2 Minimising in π

Taking the π_i one at a time, minimising in π boils down to solving p (DK) problems.

3.3.3 Stopping Criterion

We first considered the following early stopping criterion:

$$\frac{1}{L} \sum_{l=1}^L \|y_l^{(t)} - y_l^{(t+1)}\|_2 < \varepsilon.$$

Where the (t) exponent denotes the value at iteration t during the BCD algorithm.

This method, similar to [9], stops the iterations when the barycentre positions can be considered a fixed point of the BCD iterator, which ensures that the current parameters constitute a local optimum.

However, in order to be a little more comparable to SGD and GD, we also experimented with a stopping criterion based on the energy to minimise:

$$J_t - J_{t+1} < \varepsilon$$

This adopted method has the disadvantage of requiring the computation of J at every step, which is not required in the BCD algorithm.

3.3.4 BCD for GWB Algorithm

Now that we have computed the closed forms explicitly, we can formulate a BCD algorithm:

Algorithm 6: (GWB) resolution with Block-Coordinate Descent

Data: Input measure points $(X_i)_{i \in \llbracket 1, p \rrbracket} \in \prod_{i=1}^p \mathcal{M}_{K_i, d_i}(\mathbb{R})$ and weights $(a_i)_{i \in \llbracket 1, p \rrbracket} \in \prod_{i=1}^p \Sigma_{K_i}$

Barycentre weights $b \in \Sigma_L$, linear maps $(P_i)_{i \in \llbracket 1, p \rrbracket} \in \prod_{i=1}^p \mathbb{R}^{d_i}$.

barycentric coefficients $\lambda \in \Sigma_p$, precision ε , iterations N .

Result: Barycentre positions $Y \in \mathcal{M}_{L, d}(\mathbb{R})$

```

1 Initialisation: Draw  $Y \in \mathcal{M}_{L, d}(\mathbb{R})$  and let  $J_0 := +\infty$ ;
2 for  $t \in \llbracket 1, N \rrbracket$  do
3   for  $i \in \llbracket 1, p \rrbracket$  do
4     Compute the OT distance matrix  $M_{k,l}^{(i)} = \|x_k^{(i)} - P_i y_l\|_2^2$ ;
5     Compute the OT map  $\pi_i$  by solving  $\min_{\pi_i \in \Pi(a_i, b)} M_i \cdot \pi_i$ ;
6   end
7   for  $l \in \llbracket 1, L \rrbracket$  do
8     Update  $y_l$ : compute  $y_l = \left( \sum_{i=1}^p \lambda_i b_l P_i^T P_i \right)^{-1} \left( \sum_{i=1}^p \lambda_i (\pi_{\cdot, l}^{(i)})^T X_i P_i \right)$ ;
9   end
10  Compute the energy  $J_t = \sum_{i=1}^p \lambda_i M_i \cdot \pi_i$ ;
11  if  $J_{t-1} - J_t < \varepsilon$  then
12    Declare convergence and terminate.
13  end
14 end

```

Note that the optimisation in the π_i is done first in order to avoid initialising them. Furthermore, one may consider a slight modification where the indices $i \in \llbracket 1, p \rrbracket$ (line 3) and $l \in \llbracket 1, L \rrbracket$ (line 7) are done in a random order instead of in separately and in sequence.

4 A particular GWB: the Reconstruction Problem

4.1 Discrete Reconstruction Theory

In this section, we consider a fixed probability measure $\gamma_0 = \sum_{l=1}^L b_l \delta_{z_l}$ on \mathbb{R}^d and fixed linear maps $P_i : \mathbb{R}^d \rightarrow \mathbb{R}^{d_i}$, which allows us to define the input measures $\nu_i := P_i \# \gamma_0$.

The associated GWB problem can be viewed as a reconstruction problem, since the original measure γ_0 is, by construction, an optimal solution. The question is now to determine the solution set of the Reconstruction Problem (RP):

$$\mathcal{M} := \operatorname{argmin}_{\gamma \in \mathcal{P}_2(\mathbb{R}^d)} \sum_{i=1}^p \lambda_i W_2^2(P_i \# \gamma_0, P_i \# \gamma). \quad (\text{RP})$$

Observe that $\mathcal{M} \neq \emptyset$ since $\gamma_0 \in \mathcal{M}$, by construction.

We will make the following assumptions:

- $\forall i \in \llbracket 1, p \rrbracket$, $\lambda_i \neq 0$ (note that in practice, it is even topical to have $\lambda_i = \frac{1}{p}$.)
- The (z_l) are distinct. For convenience, let $Z := (z_l)_{l \in \llbracket 1, L \rrbracket}$.

For the maps P_i , we will consider two settings, the first one being a rank condition:

$$\mathcal{H}_{\text{rank}} : \begin{cases} \forall i \in \llbracket 1, p \rrbracket, \operatorname{rank} P_i = d_i \\ \operatorname{Im} P_1^T + \dots + \operatorname{Im} P_p^T = \mathbb{R}^d \end{cases}$$

Note that the second condition in $\mathcal{H}_{\text{rank}}$ is equivalent to $\bigcap_{i=1}^p \operatorname{Ker} P_i = \{0\}$.

The second setting is where the rows of the maps P_i are drawn independently from a probability \mathbb{P} on \mathbb{R}^d that admits a density with respect to the Lebesgue measure.

$$\mathcal{H}_{\mathbb{P}} : \forall i \in \llbracket 1, p \rrbracket, P_i = \begin{pmatrix} - & (u_i^{(1)})^T & - \\ - & \vdots & - \\ - & (u_i^{(d_i)})^T & - \end{pmatrix} \quad u_i^{(i)} \sim \mathbb{P} \text{ i.i.d.}$$

We shall see that by Theorem 4.1.5, we have \mathbb{P} -almost surely $\mathcal{H}_{\mathbb{P}} \implies \mathcal{H}_{\text{rank}}$ if $\sum_{i=1}^p d_i \geq d$.

This problem is highly connected to tomography, and this section can be seen as a discrete version of [1], which attempts to reconstruct a continuous measure using projections onto lines, by minimising a functional with close ties to our reconstruction problem.

4.1.1 A first bound for hyperplanes

Still in the discrete case, if γ_1 and γ_2 have the same images by the P_i , then if p is large enough, we have $\gamma_1 = \gamma_2$. Indeed, a discrete measure of cardinality k is characterised by $k + 1$ projections on distinct hyperplanes [13]. This means that the linear maps can impose a unique faithful representation γ , but does not show the unicity of GWB. In particular, it is important to keep in mind that the problem is generally incompatible, in the sense that the measures ν_i and the maps P_i prohibit the existence of a faithful coupling γ .

Note that the condition $P_i \# \gamma_1 = P_i \# \gamma_2$ is implied if γ_1 and γ_2 are optimal in the continuous case with density, as shown by [11], Proposition 3.3.

4.1.2 Characterisation of \mathcal{M}

First of all, since W_2^2 is a distance on $\mathcal{P}_2(\mathbb{R}^d)$, a solution $\gamma \in \mathcal{M}$ must satisfy $\forall i \in \llbracket 1, p \rrbracket, P_i \# \gamma = P_i \# \gamma_0$. Conversely, any probability measure $\gamma \in \mathcal{P}_2(\mathbb{R}^d)$ such that $\forall i \in \llbracket 1, p \rrbracket, P_i \# \gamma = P_i \# \gamma_0$ zeros the energy and thus is optimal.

We thus have the property below, which shows that the set of solutions of the reconstruction problem (RP) is the set of measures that have the same push-forwards by the P_i as the original measure γ_0 (hence the name "reconstruction"). Note that this result does not require $\mathcal{H}_{\text{rank}}$.

Property 4.1.1 — Correspondence to a constrained measure set

$$\mathcal{M} = \left\{ \gamma \in \mathcal{P}_2(\mathbb{R}^d) \mid \forall i \in \llbracket 1, p \rrbracket : P_i \# \gamma = P_i \# \gamma_0 \right\}. \quad (7)$$

In order to avoid confusion with the inverse, we will define as $P^{-\circ}(B)$ the reciprocal image of a set B by a map P . In order to obtain more precise results, we will need the following Lemma:

Lemma 4.1.2 — Linear push-forward formula

Let $P \in \mathcal{M}_{d,h}(\mathbb{R})$ of rank $h \leq d$ and $B \subset \mathbb{R}^h$. Then $P^{-\circ}(B) = P^T(PP^T)^{-1}B + \text{Ker}P$

This yields the following useful result under $\mathcal{H}_{\text{rank}}$:

If B is a borelian of \mathbb{R}^{d_i} , $P_i \# \gamma(B) = \gamma\left(P_i^T(P_i P_i^T)^{-1}B + \text{Ker}P_i\right)$ (8)

Below is a visualisation of the set $P^T(PP^T)^{-1}B + \text{Ker}P$, first where B is comprised of two points of \mathbb{R}^2 and $\text{Ker}P$ is a horizontal plane in 3D, and second with B a measurable set of \mathbb{R}^2 : see Figure 4.

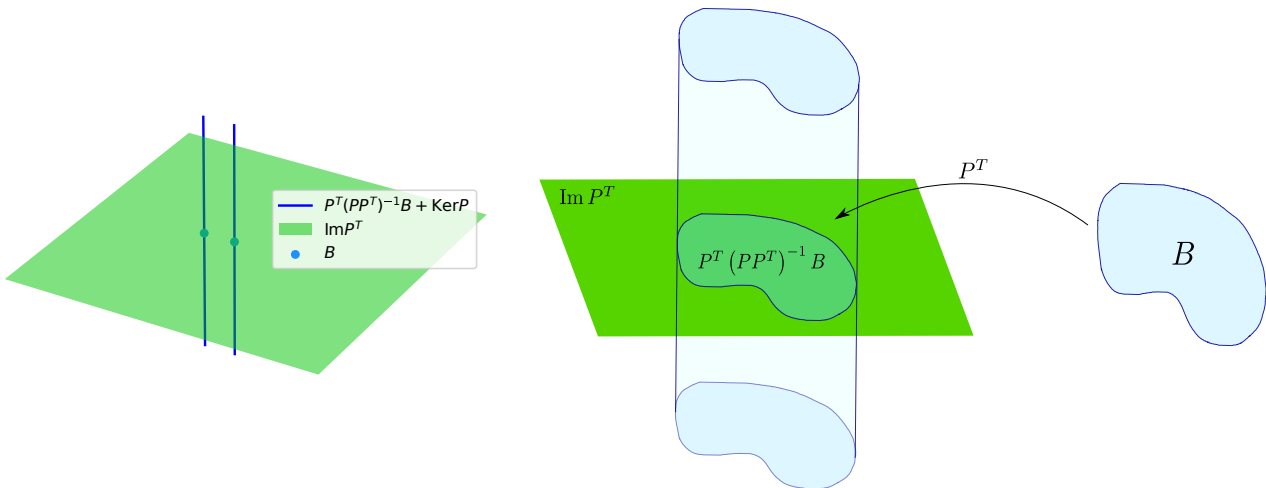


Figure 4: Two illustrations of the linear push-forward formula for a 3D to 2D projection.

Proof

If $a \in P^T(PP^T)^{-1}B + \text{Ker}P$, then by writing $a = P^T(PP^T)^{-1}b + x$ with $b \in B$ and $x \in \text{Ker}P$, we have $Pa = b \in B$, thus $a \in P^{-\circ}(B)$.

For the opposite inclusion, consider $a \in P^{-\circ}(B)$. Since P is of full rank h , we have the decomposition $\mathbb{R}^d = \text{Im}P^T \oplus \text{Ker}P$, with $Q := P^T(PP^T)^{-1}P$ the orthogonal projection on $\text{Im}P^T$.

Thus we can write $a = Qa + (I - Q)a = P^T(PP^T)^{-1}Pa + (I - Q)a$. Since $Pa \in B$, we conclude that $a \in P^T(PP^T)^{-1}B + \text{Ker}P$.

Now we can formulate a computational necessary condition for solutions in \mathcal{M} . Recall that the original measure is $\gamma_0 = \sum_{l=1}^L b_l \delta_{z_l}$, and that $Z := (z_l)_{l \in \llbracket 1, L \rrbracket}$.

Theorem 4.1.3 — Geometrical condition for solutions

Let $\gamma \in \mathcal{M}$. Under $\mathcal{H}_{\text{rank}}$, the support of γ is constrained to a discrete set:

$$\gamma \left(\bigcap_{i=1}^p (Z + \text{Ker} P_i) \right) = 1 \quad (9)$$

This can be re-written as a reunion of sets that are either singletons or empty:

$$S := \bigcap_{i=1}^p (Z + \text{Ker} P_i) = \bigcup_{(l_1, \dots, l_p) \in \llbracket 1, L \rrbracket^p} \bigcap_{i=1}^p (z_{l_i} + \text{Ker} P_i) \quad (10)$$

Conversely, for a measure supported by S to be in \mathcal{M} , it must satisfy a set of equations on the weights, which we do not describe in all generality.

Consider a setting with $p = 2$ projections onto lines in \mathbb{R}^2 , with $L = 3$ point $Z = (z_1, z_2, z_3)$. The final equation can be visualised as follows: the support of any solution is confined to the intersections between any two lines of the form $z_l + \text{Ker} P_i$. Here this corresponds to the intersecting points between an orange and a red line, allowing for 9 possible points, including the original 3: see Figure 5.

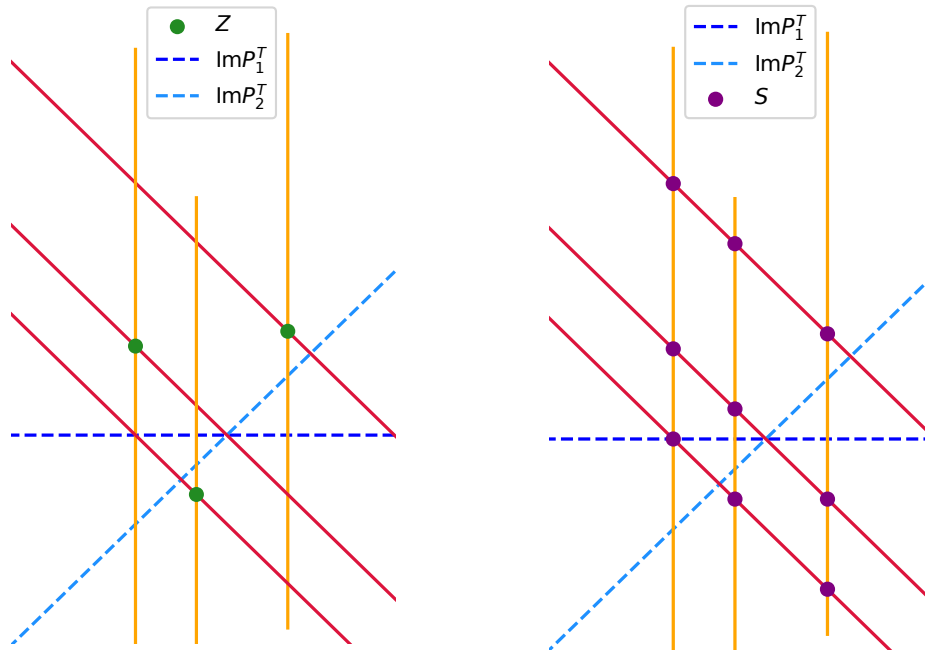


Figure 5: Illustration of the possible points for the support of a solution. On the left, Z is the original measure points, and on the right, S is the set of possible points for the support of a solution.

Proof

For the entire proof, consider $\gamma \in \mathcal{M}$ a solution.

- First, let $i \in \llbracket 1, p \rrbracket$ and A a borelian of \mathbb{R}^d such that $P_i A \cap P_i Z = \emptyset$.

Using again the decomposition $\mathbb{R}^d = \text{Im} P_i^T \oplus \text{Ker} P_i$ and $Q_i := P_i^T (P_i P_i^T)^{-1} P_i$ the orthogonal projection on $\text{Im} P_i^T$, we have $A \subset Q_i A + \text{Ker} P_i$.

Letting $B := P_i A$, $\gamma(A) \leq \gamma \left(P_i^T (P_i P_i^T)^{-1} B + \text{Ker} P_i \right)$. By our linear push-forward lemma 4.1.2, $\gamma(A) \leq P_i \# \gamma(B) = P_i \# \gamma_0(B) = 0$ since by hypothesis $B \cap P_i Z = \emptyset$.

This means that $\gamma(A) = 0$.

- Let $A_i := \text{Im}P_i^T \setminus Q_iZ + \text{Ker}P_i$, let us apply the previous result of this proof.

By construction, $P_iA_i \cap P_iZ = \emptyset$. Indeed, if there was such a $u \in P_iA_i \cap P_iZ$, then one could write $u = P_ia$ where $A_i \ni a = P_i^Tx + y$ with $x \in \mathbb{R}^{d_i}, y \in \text{Ker}P_i$, and also $u = P_iz, z \in Z$ thus $P_iP_i^Tx = P_iz$, yielding $x = (P_iP_i^T)^{-1}P_iz$ then $a = Q_iz + y$, contradicting $a \in A_i$.

Therefore using our result above, $\gamma(A_i) = 0$. Then $\gamma\left(\bigcap_{i=1}^p \overline{A_i}\right) = 1$.

Remarking that $\overline{A_i} = Q_iZ + \text{Ker}P_i = Z + \text{Ker}P_i$, we deduce the first result.

- For the second result, first note that since $\bigcap_{i=1}^p \text{Ker}P_i = \{0\}$, each affine set $\bigcap_{i=1}^p z_i + \text{Ker}P_i$ is either a single point or empty. The second equation can then be obtained by distributing the intersection over the union in $\bigcap_{i=1}^p Z + \text{Ker}P_i = \bigcap_{i=1}^p \bigcup_{l=1}^L z_l + \text{Ker}P_i$.

4.1.3 Some basic random geometry results

For this section, we consider \mathbb{P} a probability measure on \mathbb{R}^d which is absolutely continuous with respect to the Lebesgue measure. Note that these results also hold if \mathbb{P} is uniform over \mathbb{S}^d , the unit sphere.

Lemma 4.1.4 — Random Affine Hyperplanes

- Consider two affine hyperplanes $H = a + u^\perp$, $H' = b + v^\perp$, where (u, v) are mutually independent vectors of law \mathbb{P} and $a, b \in \mathbb{R}^d$. Then \mathbb{P} -almost surely, $\dim H \cap H' = d - 2$.

- Let $p \leq d$. If $\forall i \in \llbracket 1, p \rrbracket$, $H_i = a_i + u_i^\perp$ and the u_i are iid $\sim \mathbb{P}$ and $(a_i) \in \mathbb{R}^{d \times p}$,

then \mathbb{P} -almost surely, $\dim \bigcap_{i=1}^p H_i = d - p$.

First consider $u, v \sim \mathbb{P}$ iid.

Note that $\dim(u^\perp \cap v^\perp) = d - 1 \iff u^\perp = v^\perp \iff u \in \mathbb{R}v$ or $v = 0$

Since \mathbb{P} is absolutely continuous w.r.t. Lebesgue, $\mathbb{P}(v = 0) = 0$, thus we can neglect the case $v = 0$ in the following.

Then $\mathbb{P}(\dim u^\perp \cap v^\perp = d - 1) = \int_{\mathbb{R}^d} \mathbb{P}(u \in \mathbb{R}v \mid v) d\mathbb{P}(v)$

Yet given $v \in \mathbb{R}^d$ fixed, $\mathbb{P}(\mathbb{R}v) = 0$ since \mathbb{P} is absolutely continuous w.r.t. Lebesgue, thus $\mathbb{P}(\dim u^\perp \cap v^\perp = d - 1) = 0$. This proves point 1.

Point 2 is obtained using a similar reasoning to point 1 using induction. The idea is that if $\dim \bigcap_{i=1}^{p-1} H_i = d - p + 1$ and $\dim \left(\bigcap_{i=1}^{p-1} H_i \cap H_p \right) = d - p + 1$, then that would impose $u_p \in \text{Span}(u_1, \dots, u_{p-1})$, and this condition is false \mathbb{P} -a.s., since $\dim \text{Span}(u_1, \dots, u_{p-1}) < d$

Proof

In the random setting for our problem, a random linear map is almost-surely of full rank. Furthermore, given a list of p independent random linear maps $P_i : \mathbb{R}^d \rightarrow \mathbb{R}^{d_i}$, if their target dimensions are large enough, then their common kernel is trivial. Equivalently, their transposes span \mathbb{R}^d : $\sum_i \text{Im} P_i^T = \mathbb{R}^d$, a.s.. This is formalised in Property 4.1.5:

Property 4.1.5 — Kernel of random linear maps

Suppose the random map assumption $\mathcal{H}_{\mathbb{P}}$. Then \mathbb{P} -almost surely, $\forall i \in \llbracket 1, p \rrbracket$, $\text{rank} P_i = d_i$.

Furthermore, if $\sum_{i=1}^p d_i \geq d$, then \mathbb{P} -almost surely, $\bigcap_{i=1}^p \text{Ker} P_i = \{0\}$.

Proof

Let $i \in \llbracket 1, p \rrbracket$. We have $\text{rank} P_i = d - \dim \text{Ker} P_i$, and $\text{Ker} P_i = \bigcap_{j=1}^{d_i} (u_i^{(j)})^\perp$, which is of dimension $d - d_i$, \mathbb{P} -a.s. (by 4.1.4).

Assuming $\sum_{i=1}^p d_i \geq d$, we have $\bigcap_{i=1}^p \text{Ker} P_i = \bigcap_{i=1}^p \bigcap_{j=1}^{d_i} (u_i^{(j)})^\perp$, which is \mathbb{P} -a.s. of dimension $\alpha \leq d - \sum_{i=1}^p d_i$ thus is reduced to $\{0\}$.

4.1.4 Almost-sure unicity

With these new tools we can further restrict the condition on the set of solutions \mathcal{M} : Theorem 4.1.6 below shows that if the random linear maps P_i cover the original space \mathbb{R}^d with redundancy (i.e. the sum of their target space dimensions strictly exceed d), then almost surely, the reconstruction problem has a unique solution: the original measure.

Theorem 4.1.6 — Almost-sure unicity for reconstruction

If $\mathcal{H}_{\mathbb{P}}$ is verified, and further assuming $\sum_{i=1}^p d_i > d$:

\mathbb{P} -almost surely, $S = Z$ and thus $\mathcal{M} = \{\gamma_0\}$.

Proof

— *Step 1: $S \subset Z$*

Let $\mathbf{l} := (l_1, \dots, l_p) \in \llbracket 1, L \rrbracket^p$ and $S_{\mathbf{l}} := \bigcap_{i=1}^p (z_{l_i} + \text{Ker} P_i)$. We want to show $S_{\mathbf{l}} \subset Z$.

We have $\forall i \in \llbracket 1, p \rrbracket$, $\text{Ker} P_i = \bigcap_{j=1}^{d_i} (u_i^{(j)})^\perp$

Since $\sum_{i=1}^p d_i > d$, we can let $i_0 \in \llbracket 1, p \rrbracket$ and $j_0 \in \llbracket 1, d_{i_0} \rrbracket$ such that $\sum_{i=1}^{i_0} d_i + j_0 = d$, allowing us to separate $S_{\mathbf{l}} = F \cap G$ where:

$$\bullet F := \left[\bigcap_{i=1}^{i_0} (z_{l_i} + \text{Ker} P_i) \right] \cap \left[z_{l_{i_0}} + \bigcap_{j=1}^{j_0} (u_{i_0}^{(j)})^\perp \right]$$

$$\bullet G := \left[z_{l_{i_0}} + \bigcap_{j=j_0+1}^{d_{i_0}} (u_{i_0}^{(j)})^\perp \right] \cap \left[\bigcap_{i=i_0+1}^p (z_{l_i} + \text{Ker}P_i) \right]$$

By construction F is the intersection of d affine hyperplanes, hence a singleton $\{x\}$, $x \in \mathbb{R}^d$.

One may write $x = f\left((z_{l_1}, \dots, z_{l_{i_0}}), (u_i^{(j)})_{(i,j) \in T}\right)$,

where $T = \{(i, j) \in \llbracket 1, i_0 - 1 \rrbracket \times \llbracket 1, d \rrbracket : j \leq d_i\} \cup \{(i, 1), \dots, (i, j_0)\}$

which we will write concisely as $f(z_{l_1} \cdots z_{l_{i_0}}, U_0)$, U_0 denoting $(u_i^{(j)})_{(i,j) \in T}$.

Let $a + u^\perp$ one of the affine hyperplanes in the intersection defining G :

- if $j_0 + 1 \leq d_{i_0}$, let $a := z_{l_{i_0}}$ and $u := u_{i_0}^{(j_0+1)}$.
- if $j_0 + 1 > d_{i_0}$, let $a := z_{l_{i_0+1}}$ and $u := u_{i_0+1}^{(1)}$.

Let us now reason conditionally to U_0 , which in particular determines x but not $u \sim \mathbb{P}$.

We have $S_1 \subset \{x\} \cap a + u^\perp$.

- If $x = a$, then $S_1 \subset Z$ and this step is finished.
- If $x \neq a$ then let $y := x - a \neq 0$. We have $x \in a + u^\perp \Leftrightarrow u \in y^\perp$.
But $\mathbb{P}(u \in y^\perp) = 0$ thus $S_1 = \emptyset$, \mathbb{P} -a.s.

We conclude that $S_1 \subset Z$ and thus $S \subset Z$, \mathbb{P} -a.s. after marginalising over U_0 .

— *Step 2: $S \supset Z$*

To show $S \supset Z$, let $l \in \llbracket 1, L \rrbracket$ and $\mathbf{l} := (l, \dots, l)$. We have $S_{\mathbf{l}} = \bigcap_{i=1}^p (z_l + \text{Ker}P_i) = \{z_l\}$.

— *Step 3: $\mathcal{M} = \{\gamma_0\}$*

Now we have proven $S = Z$, and thus that any solution $\gamma \in \mathcal{M}$ is supported by Z , let us write it $\gamma = \sum_{l=1}^L a_l \delta_{z_l}$. In the following we let $i := 1$ and drop the i index for convenience.

Since $\gamma \in \mathcal{M}$, we have in particular $P\#\gamma = P\#\gamma_0$ ($P := P_1$). Thus $\sum_{l=1}^L a_l \delta_{Pz_l} = \sum_{l=1}^L b_l \delta_{Pz_l}$.

If the points $(Pz_l)_{l \in \llbracket 1, L \rrbracket}$ are distinct, then this implies $\forall l \in \llbracket 1, L \rrbracket$, $a_l = b_l$.

Let $l \neq l' \in \llbracket 1, L \rrbracket^2$. We have $Pz_l = Pz_{l'} \Leftrightarrow z_l - z_{l'} \in \text{Ker}P$.

Yet $x := z_l - z_{l'} \neq 0$ by hypothesis on Z . $x \in \text{Ker}P \Leftrightarrow \forall j \in \llbracket 1, d_1 \rrbracket$, $u^{(j)} \cdot x = 0$

Letting $u := u^{(1)}$, we have $\mathbb{P}(x \in \text{Ker}P) \leq \mathbb{P}(u \cdot x = 0) = \mathbb{P}(u \in x^\perp) = 0$

Finally, \mathbb{P} -a.s. the (Pz_l) are distinct and thus $\forall l \in \llbracket 1, L \rrbracket$, $a_l = b_l$

We can finally conclude that $\mathcal{M} = \{\gamma_0\}$, \mathbb{P} -a.s..

4.1.5 Discussing the general case

The previous theorem 4.1.6 only provides unicity almost-surely, however "improbable" counter examples do exist with excessive symmetry. Below we present a counter-example adapted from [13]. Let $d := 2$, $p := L > d$ and $\forall i \in \llbracket 1, p \rrbracket, d_i := 1$.

Consider $z_l := \left(\cos \left(\frac{(2l+1)\pi}{L} \right), \sin \left(\frac{(2l+1)\pi}{L} \right) \right)^T$, $P_l := \left(\cos \left(\frac{(2l+1)\pi}{2L} \right), \sin \left(\frac{(2l+1)\pi}{2L} \right) \right)$.

As can be seen below (Figure 6), for $L = 3$, this corresponds to placing the (z_l) on every other vertex of a regular $2L$ -gon, and defining the P_l such that $\text{Im}P_l^T$ is the l -th bisector of the $2L$ -gon.

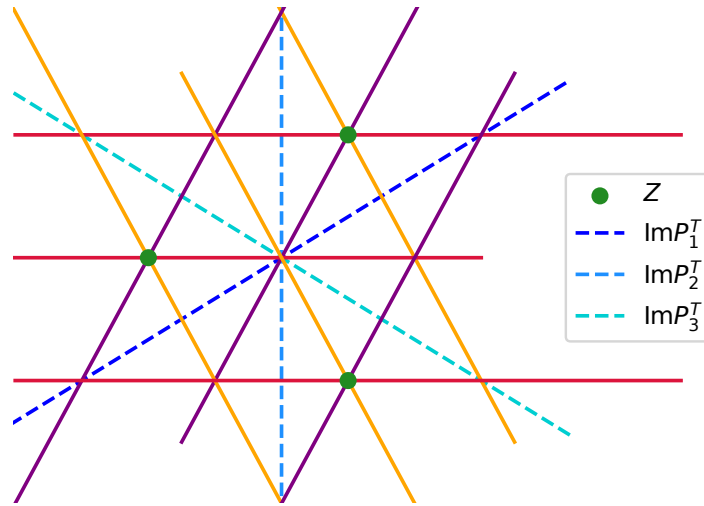


Figure 6: Illustration of a pathological sur-critical case without unicity for specific projections P_i .

The points of S are the points of the form $\bigcap_{i=1}^3 (z_{l_i} + \text{Ker}P_i)$, or visually the intersection points of a yellow line, a red line and a purple line. We can see that the remaining vertices of the polygon constitute another valid measure γ' whose push-forwards $P_i \# \gamma'$ are all the same as those of the original measure.

In the case of hyperplanes ($h = d - 1$), Theorem I.2 from [13] indicates that a necessary condition for unicity is $p > L$, whereas in our almost-sure setting (Theorem 4.1.6), we obtained a condition that is independent of L .

4.1.6 The critical case $\sum_i d_i = d$

In the theorem below, we show that the example below Theorem 4.1.3 is representative of the critical case: in this case, $d = 2$, $p = 2$, $d_1 = d_2 = 1$, $L = 3$ and there are $9 = 3^2 = L^p$ points in S .

Theorem 4.1.7 — Number of admissible points in the critical case

In the critical case $\sum_{i=1}^p d_i = d$ and under $\mathcal{H}_{\mathbb{P}}$, $\boxed{\#S = L^p}$ \mathbb{P} -almost surely.

— Step 1: finding other valid points

Proof

Recall 4.1.3: $S = \bigcup_{(l_1, \dots, l_p) \in \llbracket 1, L \rrbracket^p} \bigcap_{i=1}^p (z_{l_i} + \text{Ker}P_i)$.

Let $\mathbf{l} \in \llbracket 1, L \rrbracket^p$ non constant, and let $i \neq j \in \llbracket 1, L \rrbracket^2$ such that $l_i \neq l_j$.

Then $S_{\mathbf{l}} \subset (z_{l_i} + u^\perp) \cap (z_{l_j} + v^\perp)$, where $u := u_i^{(1)}$, $v := u_j^{(1)} \sim \mathbb{P} \otimes \mathbb{P}$.

Then since the (z_l) are distinct, $z_{l_i} \neq z_{l_j}$ and, using the same random geometry arguments as 4.1.6, step 3, we show separately that $\mathbb{P}(z_{l_i} \in z_{l_j} + v^\perp) = 0$ and $\mathbb{P}(z_{l_j} \in z_{l_i} + v^\perp) = 0$.

This proves that \mathbb{P} -a.s., z_{l_i} and z_{l_j} do not belong to S_1 .

Since we have shown this for any $i \neq j$ verifying $l_i \neq l_j$, we have proven that $S_1 \cap \{z_{l_1}, \dots, z_{l_p}\} = \emptyset$.

Then for the other points: let $z \in Z \setminus \{z_{l_1}, \dots, z_{l_p}\}$, we can use the same argument to prove that $\mathbb{P}\left(z \in z_{l_1} + \left(u_1^{(1)}\right)^\perp\right) = 0$, and conclude $S_1 \cap Z = \emptyset$, \mathbb{P} -a.s..

Finally, since S_1 is the intersection of exactly d random affine hyperplanes in dimension d , by Lemma 4.1.4, is it a point, \mathbb{P} -a.s..

— *Step 2: counting distinct the valid points*

Let $\mathbf{I} \neq \mathbf{I}' \in \llbracket 1, L \rrbracket^p \times \llbracket 1, L \rrbracket^p$ non constant. By Step 1, $S_1, S_{\mathbf{I}'}$ are points outside Z .

Let $i \in \llbracket 1, p \rrbracket$ such that $l_i \neq l'_i$. Then $S_1 \cap S_{\mathbf{I}'} \subset (z_{l_i} + \text{Ker}P_i) \cap (z_{l'_i} + \text{Ker}P_i) =: F$

Then since the (z_l) are distinct, we have $z_{l_i} \neq z_{l'_i}$, thus $F = \emptyset$.

We can conclude that the singletons $S_1, S_{\mathbf{I}'}$ are distinct. Finally, $\#S = \#\llbracket 1, L \rrbracket^p = L^p$.

To finish with the critical case, we will now count the number of optimal measures.

We assume that $\forall i \in \llbracket 1, p \rrbracket, \forall l \neq l' \in \llbracket 1, L \rrbracket^2, P_i z_l \neq P_i z_{l'}$ (which is true \mathbb{P} -a.s. under $\mathcal{H}_{\mathbb{P}}$).

— *Counting the optimal measures in the uniform case*

If we impose both γ and γ_0 to be uniform discrete measures with L points (following the framework of [9]), then there is a finite (but combinatorial) amount of optimal measures.

Indeed, choosing an optimal solution corresponds to choosing L points x_1 , with the only restriction of having for every pair of points x_1, x'_1 verify $\forall i \in \llbracket 1, p \rrbracket, l_i \neq l'_i$ (otherwise they would have the same projection by P_i and excessive weight on $P_i z_{l_i}$, contradicting the requirements for \mathcal{M} : 4.1.1).

We then have $\prod_{n=0}^{L-1} \#\llbracket 1, L - n \rrbracket^p = (L!)^p$ possibilities.

— *Counting the optimal measures without restrictions on the weights*

Now we take a look at the number of optimal measures, so let $\gamma \in \mathcal{M}$. We have \mathbb{P} -a.s.:

By Theorem 4.1.7, γ is of the form $\gamma = \sum_{\mathbf{l} \in \llbracket 1, L \rrbracket^p} a_{\mathbf{l}} \delta_{x_{\mathbf{l}}}$. Note that by construction, $P_i x_{\mathbf{l}} = P_i z_{l_i}$.

By Property 4.1.1, we have $\forall i \in \llbracket 1, p \rrbracket, P_i \# \gamma = P_i \# \gamma_0$.

Let $i \in \llbracket 1, p \rrbracket$. We have $P_i \# \gamma = \sum_{\mathbf{l} \in \llbracket 1, L \rrbracket^p} a_{\mathbf{l}} \delta_{P_i z_{l_i}}$. which equals $P_i \# \gamma_0 = \sum_{k=1}^L b_k \delta_{P_i z_k}$.

Since the $(P_i z_l)_l$ are assumed distinct, this entails for all $k \in \llbracket 1, L \rrbracket$: $\sum_{\mathbf{l}_{-i} \in \llbracket 1, L \rrbracket^{p-1}} a_{\mathbf{l}_{-i}, \dots, l_{i-1}, k, l_{i+1}, \dots, l_p} = b_k$

Where \mathbf{l}_{-i} indicates that we index this $(p-1)$ -tuple on $\llbracket 1, L \rrbracket \setminus \{k\}$.

We can re-write this condition as $a \in \Pi_L^p(b)$, the set of L -dimensional p -tensors on \mathbb{R}_+ ($\mathbb{R}_+^{L^p}$) with all p marginals equal to b .

Conversely, if γ is of the form $\gamma = \sum_{\mathbf{l} \in \llbracket 1, L \rrbracket^p} a_{\mathbf{l}} \delta_{x_{\mathbf{l}}}$ with $a \in \Pi_L^p(b)$, then we have by construction

$\forall i \in \llbracket 1, p \rrbracket, P_i \# \gamma = P_i \# \gamma_0$ and thus $\gamma \in \mathcal{M}$.

In particular, there is an *infinite* amount of solutions to the reconstruction problem in the critical case, \mathbb{P} -a.s..

4.2 Consequences on Sliced Wasserstein Methods

4.2.1 Reminders on Sliced Wasserstein Methods

Before drawing conclusions on Sliced OT, we provide a swift reminder of said methods [7, 17].

Definition 4.2.1 — Sliced Wasserstein Distance [7]

The Sliced Wasserstein Distance between two measures $\alpha, \beta \in \mathcal{P}_2(\mathbb{R}^d)$ is defined as:

$$\text{SW}_2^2(\alpha, \beta) := \int_{\mathbb{S}^d} W_2^2(P_\theta \# \alpha, P_\theta \# \beta) d\sigma(\theta) \quad (11)$$

Where \mathbb{S}^d is the unit sphere of \mathbb{R}^d , σ is the uniform measure on \mathbb{S}^d and $P_\theta := \theta^T$.

In practice, the integral defining SW_2^2 is estimated via a Monte-Carlo method.

The main algorithm used for computing the sliced OT mapping is the following. It relies on the fact that discrete 1D OT boils down to sorting the input points (see [17], §2.6 for example).

Algorithm 4.2.2 — Lagrangian SGD for sliced OT [7]

Let $\alpha := \frac{1}{L} \sum_{l=1}^L \delta_{x_l}$, $\beta := \frac{1}{L} \sum_{l=1}^L \delta_{y_l}$ two discrete probability measures on \mathbb{R}^d .

Consider the closed-form $W_2^2(P_\theta \# \alpha, P_\theta \# \beta) = \sum_{l=1}^L \|\langle x_{\sigma(l)} - y_{\tau(l)}, \theta \rangle\|_2^2$ (12),

where $\theta \in \mathbb{S}^d$ and $\sigma, \tau \in \mathfrak{S}_L$ respectively sort the numbers $(x_l \cdot \theta)_l$ and $(y_l \cdot \theta)_l$.

The mapping from α to β is done using SGD by sampling $\theta \in \mathbb{S}^d$ and descending the gradient of the closed form above with respect to $(x_l)_l$.

4.2.2 Consequence of the reconstruction result on SW

Given the previous reconstruction results 4.1.6 and 4.1.7, we can describe a pathological behaviour of SW with insufficient projections. We consider \mathbb{P} a probability over \mathbb{S}^d admitting a density (in particular, this applies to the logical example \mathbb{P} uniform).

Property 4.2.3 — SW with insufficient projections

Let $\gamma_0 := \sum_{l=1}^L b_l \delta_{z_l}$, where the (z_l) are distinct and $b \in \Sigma_L$. Assume that $\theta_1, \dots, \theta_p$ i.i.d. $\sim \mathbb{P}$.

Consider the MC estimator $\widehat{\text{SW}}_p(\alpha, \beta) := \frac{1}{p} \sum_{i=1}^p W_2^2(P_{\theta_i} \# \alpha, P_{\theta_i} \# \beta)$ and assume $p \leq d$.

\mathbb{P} -a.s., there exists an infinity of measures $\gamma \neq \gamma_0 \in \mathcal{P}_2(\mathbb{R}^d)$ verifying $\widehat{\text{SW}}_p(\gamma_0, \gamma) = 0$.

Note that, as can be seen visually in 3.3: *Discussing the general case*, the other measures γ can be at an arbitrarily large W_2^2 distance of γ_0 :

- In the case $p < d$, this is easy to see since the set of admissible points is infinite and one may consider measures with points that are arbitrarily far away.
- In the case $p = d$, the distance can be grown by scaling the points of γ_0 further away from the origin.

In practice, sliced-Wasserstein Generative Models compute SW in the data space or in the data encoding space ([15], [12]), which yields high values of d , in particular for images. Note that the necessity behind having a large p was already hinted at in [15], §3.3.

Conversely, in the sur-critical case, the Sliced distance does have the desired property:

Property 4.2.4 — SW with sufficient projections

In the same framework as the previous property, assume now $p > d$.

Then by 4.1.6, \mathbb{P} -almost surely, $\{\gamma_0\} = \operatorname{argmin}_{\gamma \in \mathcal{P}_2(\mathbb{R}^d)} \widehat{\text{SW}}_p(\gamma_0, \gamma)$.

4.3 Experimental Results on the Reconstruction Problem

4.3.1 The Spiral Dataset

These tests can be reproduced using [release v1.0 of our repository](#). We build a spiral dataset, where a discrete uniform measure μ is projected p times using randomly drawn linear maps $P_i : \mathbb{R}^d \rightarrow \mathbb{R}^h$ (with normalised lines): the input measures are defined as $P_i \# \mu$.

We define our 2D spiral as a sampling of $x(t) = \frac{t}{\pi} \begin{pmatrix} \cos(t) \\ \sin(t) \end{pmatrix}$, $t \in [0, 2N\pi]$, (N loops).

In higher dimensions $d = 2n$ we stack n 2D spirals, for $d = 2n + 1$ we stack n 2D spirals and a linear dimension $[t]$

In this setting, we can view the (GWB) problem as a reconstruction problem: the estimated barycentre γ can be compared with the ground truth μ (the spiral) using the 2-Wasserstein distance.

For the entire study, we consider measures made of $K_i = L = 30$ points. In order to get a visual understanding of the values, below is an example with $p = 15$ projections, which yields an imperfect reconstruction the original spiral, as can be seen in Figure 7. Given the reconstruction results (in particular Theorem 4.1.6), we conclude that in this case, BCD converges towards a local optimum that is different to the unique global optimum (the original spiral).

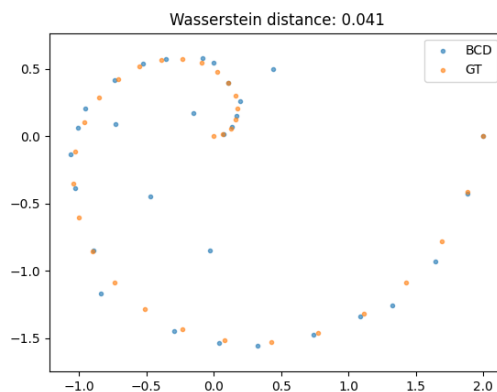


Figure 7: (GWB) resolution on the spiral dataset using BCD with $p = 15$ projections.

Note that with a Wasserstein distance of 0.041, the spiral is imperfectly reconstructed, yet recognisable.

As discussed in §4.1, with $L = 30$ points and $p = 40$ projections, there is a unique global optimum (which is the original measure γ_0), however there are local optima, which are unfortunately reached by the algorithms:

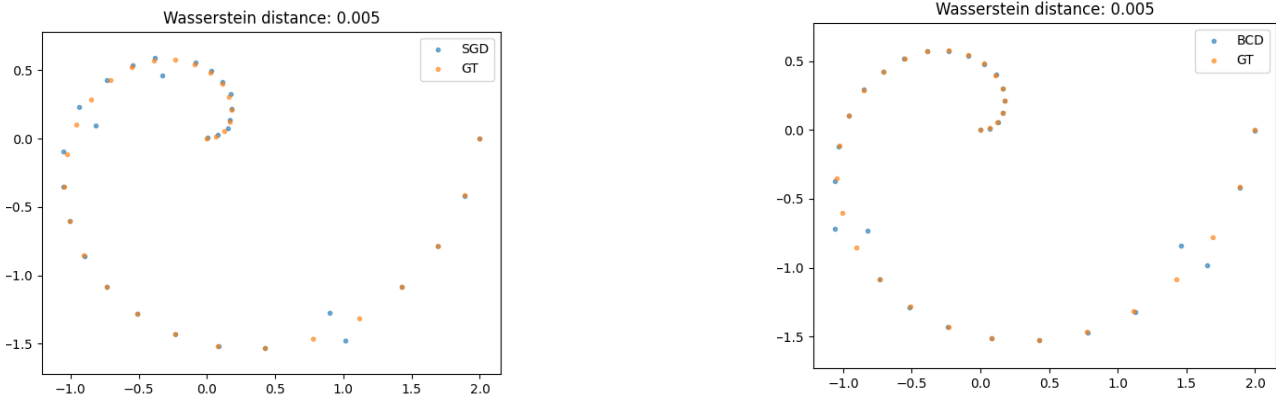


Figure 8: (GWB) resolution on the spiral dataset using SGD and BCD, with $p = 40$ projections.

In this section, we test our three solvers (GD: Algorithm 4, SGD: Algorithm 5 and BCD: Algorithm 6) on the Reconstruction Problem (RP) using the toy spiral dataset.

4.3.2 Impact of the precision ε

This test can be reproduced using [version 1.3 of our repository](#). For the same reconstruction setting, we test our different algorithms while varying the target precision ε and monitor the impact on the W_2^2 reconstruction error and on the computation time.

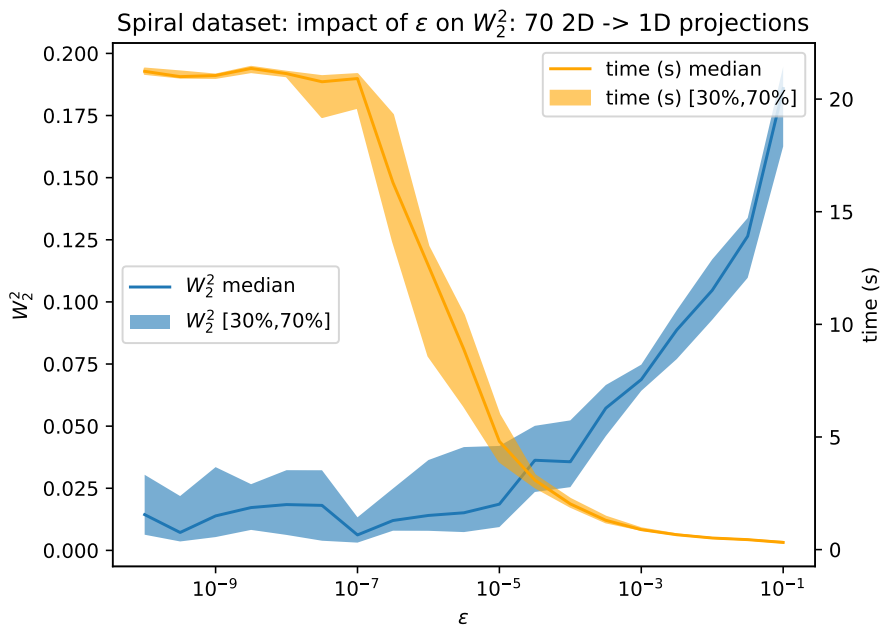


Figure 9: GD resolution of the reconstruction problem, varying the precision ε . The original measure is a 30-point 2D spiral, projected 70 times on lines, and the solver is capped at 300 iterations.

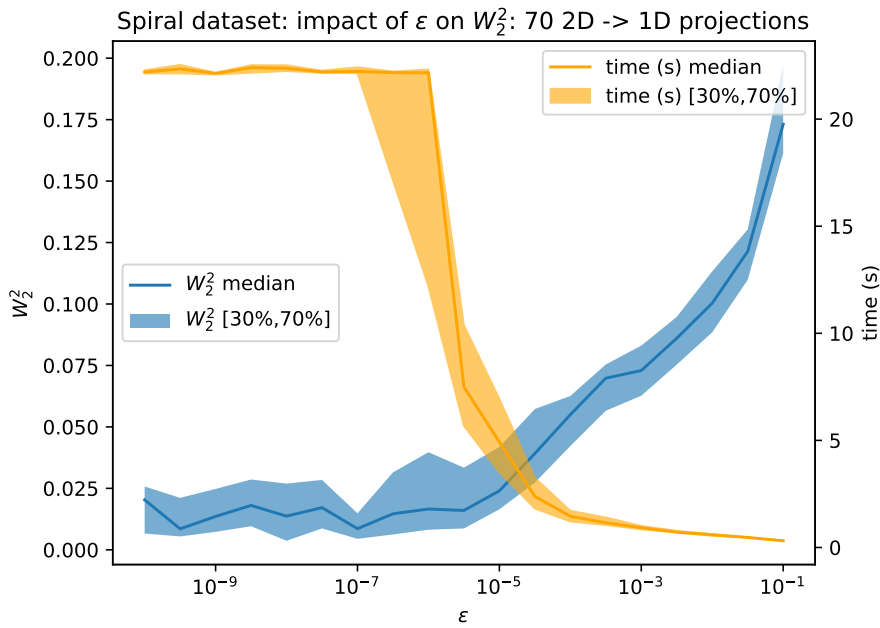


Figure 10: SGD resolution of the reconstruction problem, varying the precision ϵ . The original measure is a 30-point 2D spiral, projected 70 times on lines and the solver is capped at 300 iterations.

Neither Gradient Descent nor Stochastic Gradient Descent manage to reconstruct the original spiral exactly (see Figure 9 and Figure 10), whatever the target precision: the median W_2^2 distance is around 0.015, which corresponds to an imperfect but recognisable spiral.

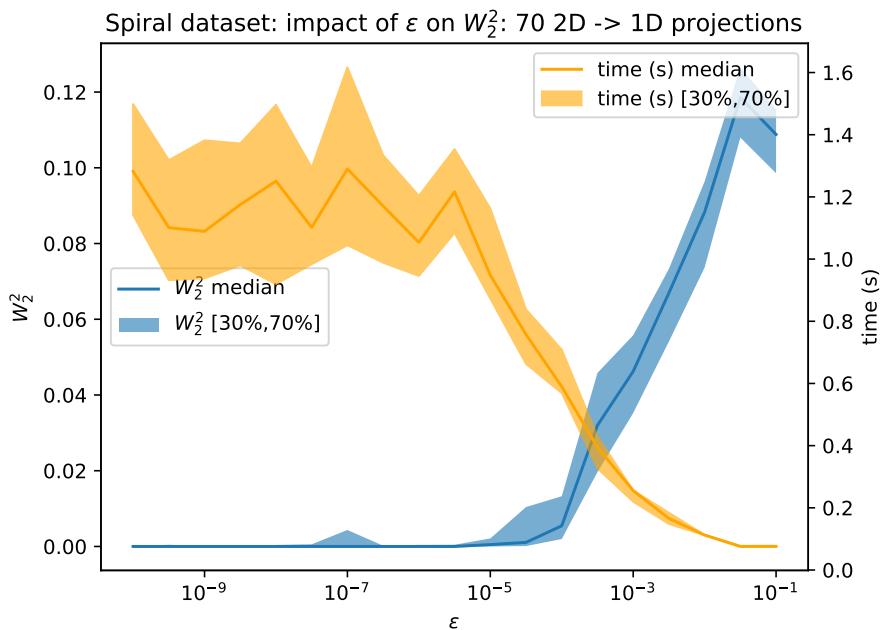


Figure 11: BCD resolution of the reconstruction problem, varying the precision ϵ . The original measure is a 30-point 2D spiral, projected 70 times on lines and the solver is capped at 300 iterations.

However, not only is BCD over ten times faster, with a target precision of under 10^{-5} , the original spiral is properly reconstructed at numerical precision (see Figure 11).

Theoretically (Theorem 4.1.6), with $p = 70$ projections the setting is well past the minimum required projections ($p = 3$ in 2D for 1D projections), however experiments show that a significantly greater amount of projections is required in order to limit the number of local optima and stationary points.

4.3.3 Impact of the number of projections p

Now we let the number of projections p vary and observe the median W_2 distance (over 30 samples to account for algorithm instability) and also display the .3 and .7 quantiles as an area around the $W_2(p)$ curve. This test can be reproduced using [version 1.3 of our repository](#).

If the optimisation was perfect, as soon as $p \geq 3$ in dimension 2, we should converge to the unique global optimum (see Theorem 4.1.6). However, neither SGD nor GD converge towards the original measure (see Figure 12), while BCD converges most of the time for $p \geq 55$ with significant instability in 2D.

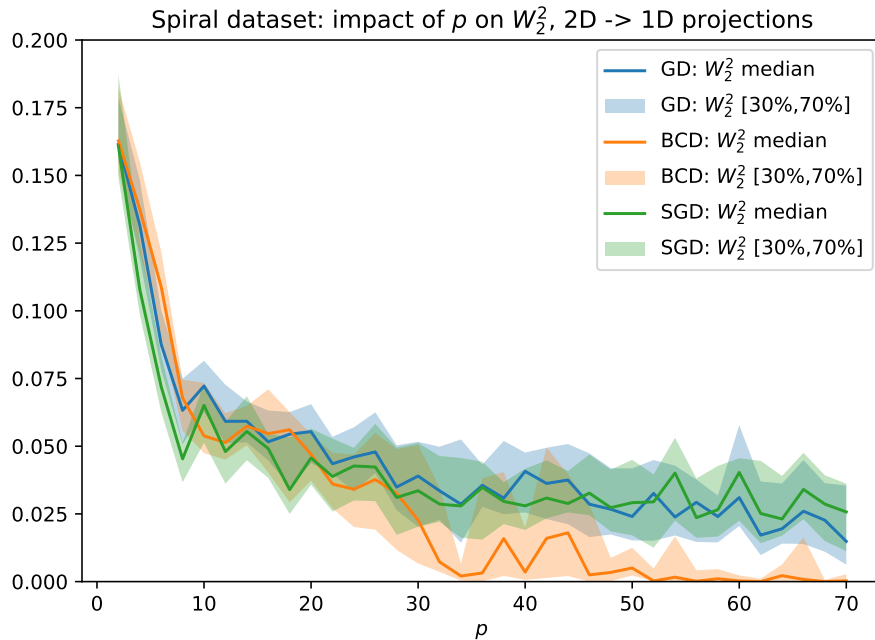


Figure 12: Impact of the number of projections p on the W_2^2 reconstruction error for the (GWB). The original measure is a 30-point 2D spiral, projected on lines and the solver is capped at 300 iterations with a precision of 10^{-5} .

In the 3D case with projections on lines (Figure 13), the discrepancy between the Gradient Descent Methods and the Block Coordinate Method is clear, with essential convergence of BCD for $p \geq 60$.

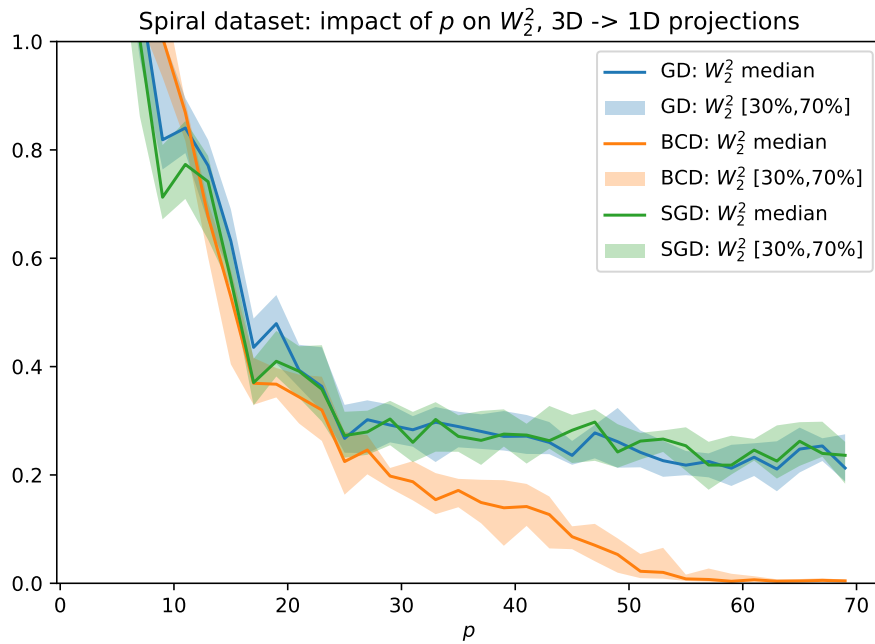


Figure 13: Impact of the number of projections p on the W_2^2 reconstruction error for the (GWB) resolution. The original measure is a 30-point 3D spiral, projected on lines and the solver is capped at 300 iterations with a precision of 10^{-5} .

An important factor to keep in mind is the importance of the projection dimension: in 3D, with projections on planes, all algorithms reconstruct perfectly for $p \geq 10$. This shows that the non-convex optimisation problem becomes significantly simpler with higher-dimensional projections, even though there is theoretical unicity with less information.

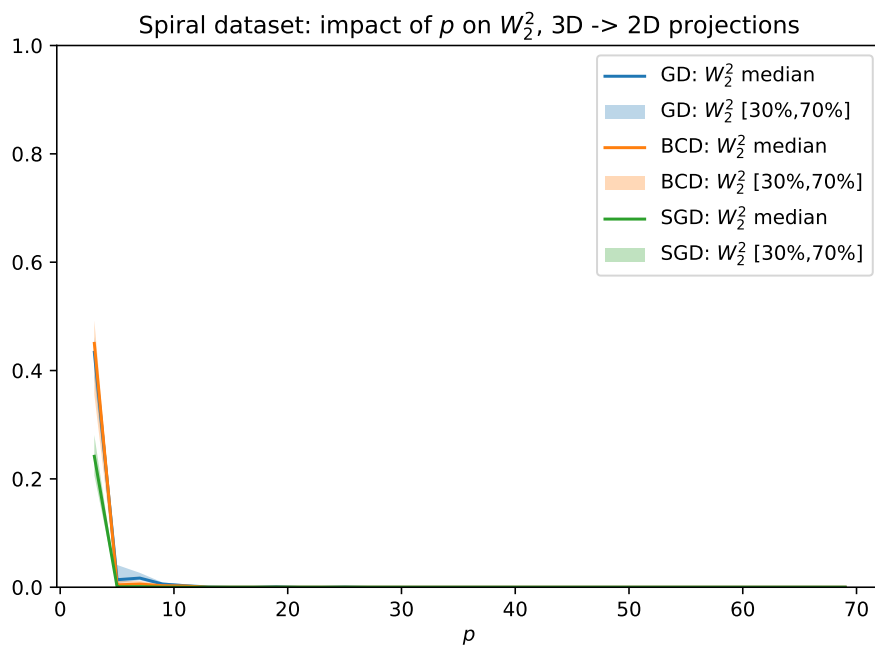


Figure 14: Impact of the number of projections p on the W_2^2 reconstruction error for the (GWB) resolution. The original measure is a 30-point 3D spiral, projected on planes and the solver is capped at 300 iterations with a precision of 10^{-5} .

This suggests that reconstruction from lines in any dimension higher than 2 is fairly unreliable and would require a prohibitive amount of projections, contrary to intuition. This is due to optimisation woes and not to a lack of unicity.

4.4 Local Optima: $L = 2$ case

In order to study the local optima for the reconstruction problem (RP), we first study the case $L = 2$. For simplicity, we consider the measures weights to be fixed and uniform.

We optimise the coupling positions y_1, y_2 by minimising

$$E(y_1, y_2) := \sum_{i=1}^p \lambda_i W_2^2 \left(\frac{1}{2} \delta_{P_i z_1} + \frac{1}{2} \delta_{P_i z_2}, \frac{1}{2} \delta_{P_i y_1} + \frac{1}{2} \delta_{P_i y_2} \right).$$

Since E is non convex, we use alternated optimisation with the energy:

$$J(Y, \pi) = \sum_{i=1}^p \lambda_i \sum_{(k,l) \in \llbracket 1,2 \rrbracket^2} \pi_{k,l}^{(i)} \|P_i z_k - P_i y_l\|_2^2 \quad (13)$$

With $Y = (y_1, y_2) \in \mathbb{R}^d \times \mathbb{R}^d$ and $\pi = (\pi_1, \dots, \pi_p) \in \mathbb{U}_2^p$, where $\mathbb{U}_2 := \Pi \left(\left(\frac{1}{2}, \frac{1}{2} \right), \left(\frac{1}{2}, \frac{1}{2} \right) \right)$.

We begin with a result regarding the problem structure:

Theorem 4.4.1 — Cell Structure for 2-point reconstruction

The assignment map $\mathcal{A} : \begin{cases} \mathbb{R}^d \times \mathbb{R}^d & \longrightarrow & \mathcal{P}(\mathbb{U}_2^p) \\ (y_1, y_2) & \longmapsto & \operatorname{argmin}_{\pi \in \mathbb{U}_2^p} J(Y, \pi) \end{cases}$ is piece-wise constant, defining a cell structure on \mathbb{R}^{2d} .

Precisely, let $(y_1, y_2) \in \mathbb{R}^{2d}$ and $i \in \llbracket 1, p \rrbracket$. There are three cases:

- Under (\mathcal{C}_e^i) : $\|P_i(z_1 - y_1)\|_2^2 + \|P_i(z_2 - y_2)\|_2^2 = \|P_i(z_2 - y_1)\|_2^2 + \|P_i(z_1 - y_2)\|_2^2$,
any $\pi_i = \frac{1}{2} \begin{pmatrix} a & 1-a \\ 1-a & a \end{pmatrix}$, $a \in [0, 1]$ is optimal.
- Under (\mathcal{C}_m^i) : $\|P_i(z_1 - y_1)\|_2^2 + \|P_i(z_2 - y_2)\|_2^2 < \|P_i(z_2 - y_1)\|_2^2 + \|P_i(z_1 - y_2)\|_2^2$,
 $\pi_i = \frac{1}{2} \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix}$ is optimal.
- Under (\mathcal{C}_s^i) : $\|P_i(z_1 - y_1)\|_2^2 + \|P_i(z_2 - y_2)\|_2^2 > \|P_i(z_2 - y_1)\|_2^2 + \|P_i(z_1 - y_2)\|_2^2$,
 $\pi_i = \frac{1}{2} \begin{pmatrix} 0 & 1 \\ 1 & 0 \end{pmatrix}$ is optimal.

Notes on the names of the three cases:

- (\mathcal{C}_e) ("e" for "equal") corresponds to ambiguity between matching y_1, y_2 to z_1, z_2 or z_2, z_1 . It is the case shown in [Figure 6](#).
- (\mathcal{C}_m) ("m" for "match"): matching y_1, y_2 to z_1, z_2 is less costly than z_2, z_1 .
- (\mathcal{C}_s) ("s" for "swap"): matching y_1, y_2 to z_2, z_1 is less costly than z_1, z_2 .

Furthermore, note that by expanding the square distances, one may re-write the conditions. One recovers conditions of the form:

$$(\mathcal{C}_m^i) : P_i^T P_i (z_2 - z_1) \cdot (y_2 - y_1) > 0 \quad (14)$$

Notably, this condition is linear in $y_2 - y_1$ and only depends on this difference.

Proof

First of all, by combining the three independent linear equations that derive from $\pi_i \in \mathbb{U}_2$: $\pi_i(\frac{1}{2}, \frac{1}{2})^T = (\frac{1}{2}, \frac{1}{2})^T$ and $\pi_i^T(\frac{1}{2}, \frac{1}{2})^T = (\frac{1}{2}, \frac{1}{2})^T$, one may prove than an optimal coupling is always of the form

$$\pi_i = \frac{1}{2} \begin{pmatrix} a & 1-a \\ 1-a & a \end{pmatrix}, \quad a \in [0, 1].$$

Then for such a coupling π_i , we write $\sum_{(k,l) \in \llbracket 1,2 \rrbracket^2} \pi_{k,l}^{(i)} \|P_i z_k - P_i y_l\|_2^2 = at + c$,

with $t = \frac{1}{2} (\|P_i(z_2 - y_1)\|_2^2 + \|P_i(z_1 - y_2)\|_2^2 - \|P_i(z_1 - y_1)\|_2^2 - \|P_i(z_2 - y_2)\|_2^2)$,

and $c = \frac{1}{2} (\|P_i(z_2 - y_1)\|_2^2 + \|P_i(z_1 - y_2)\|_2^2)$.

We solve $\operatorname{argmin}_{a \in [0,1]} aw + c$, yielding the announced three cases depending on t .

We define the *matching configuration* of a position (y_1, y_2) as $\mathbf{m} \in \{-1, 0, 1\}^p$ such that $\forall i \in \llbracket 1, p \rrbracket$, $m_i = 0$ under C_e^i , 1 under C_m^i and -1 under C_s^i . These configurations encode the cell structure of the (y_1, y_2) space, where a *cell* is defined as the set $C_{\mathbf{m}} = \{(y_1, y_2) \in \mathbb{R}^{2d} \mid \mathbf{m}(y_1, y_2) = \mathbf{m}\}$, for a given non-degenerate configuration $\mathbf{m} \in \{\pm 1\}^p$. A cell is determined by the conjunction of linear conditions (see (14)), thus it is either empty or a polytope.

By Theorem 4.4.1, a matching configuration $\mathbf{m} \in \{-1, 1\}^p$ (assumed with no ambiguities, i.e. each $\mathbf{m}_i \neq 0$) defines a unique optimal coupling $\pi_{\mathbf{m}}$ where each $\pi_{\mathbf{m}}^{(i)}$ satisfies the corresponding case in 4.4.1.

For theoretical purposes, if we are in a toy situation with known positions z_1, z_2 , we propose the following algorithm that computes the configurations of all the non-empty cells. Algorithm 7 is a refinement over the brute-force solution which checks the emptiness of all 2^p cells:

Algorithm 7: Recursive computation of valid cell configurations

Data: True positions $z_1, z_2 \in \mathbb{R}^d$, linear maps $\forall i \in \llbracket 1, p \rrbracket$, $P_i \in \mathcal{M}_{d_i, d}(\mathbb{R})$ and tolerance $t > 0$.

Result: List M of all the non-empty cell configurations.

- 1 **Initialisation:** $M = \emptyset$;
 - 2 **if** $p = 1$ **then**
 - 3 **Return** $M := \{(1), (-1)\}$;
 - 4 **end**
 - 5 Compute recursively M_{-1} the non-empty cell configurations for $i \in \llbracket 1, p-1 \rrbracket$, i.e. the output of Algorithm 7 for $(P_i)_{i \in \llbracket 1, p-1 \rrbracket}$;
 - 6 Let $\widetilde{M} := \{(\mathbf{m}, \varepsilon) \mid \mathbf{m} \in M_{-1}, \varepsilon \in \{\pm 1\}\}$;
 - 7 **for** $\mathbf{m} \in \widetilde{M}$ **do**
 - 8 **If**
 - 8 $\{x \in \mathbb{R}^d \mid \forall i \in \llbracket 1, p \rrbracket : m_i P_i^T P_i (z_2 - z_1) \cdot x > 0\} \neq \emptyset$, $M \leftarrow M \cup \{\mathbf{m}\}$;
 - 9 **end**
 - 10 **Return** M ;
-

Note that line 8, is done in practice by solving the linear feasibility problem
$$\min_{x \in \mathbb{R}^d} \max_{\forall i \in \llbracket 1, p \rrbracket : m_i P_i^T P_i (z_2 - z_1) \cdot x \geq t} 1.$$

Any non-empty cell is an unbounded polytope, since its is defined by an intersection of linear inequations (see (14)) with the common second term 0. Therefore, the value of the tolerance parameter does not matter, since a non-empty cell is stable by multiplication by a positive factor.

Theorem 4.4.2 — Cell movement

Let $(y_1, y_2) \in \mathbb{R}^d \times \mathbb{R}^d$ and let \mathbf{m} its configuration. Assume $\mathbf{m} \in \{\pm 1\}^p$.

Then the next optimal positions $y_1^*, y_2^* \in \underset{(y_1, y_2) \in \mathbb{R}^{2d}}{\operatorname{argmin}} J(y_1, y_2, \pi_{\mathbf{m}})$ are

$$y_1^* = A^{-1}(B_{\mathbf{m}}z_1 + C_{\mathbf{m}}z_2) \text{ and } y_2^* = A^{-1}(B_{\mathbf{m}}z_2 + C_{\mathbf{m}}z_1),$$

$$\text{where } B_{\mathbf{m}} = \sum_{\substack{i=1 \\ m_i=1}}^p \lambda_i P_i^T P_i, \quad C_{\mathbf{m}} = \sum_{\substack{i=1 \\ m_i=-1}}^p \lambda_i P_i^T P_i \text{ and } A = \sum_{i=1}^p \lambda_i P_i^T P_i.$$

(y_1^*, y_2^*) is a local optimum iff its cell is stable, i.e. if its configuration \mathbf{m}' equal \mathbf{m} . This is equivalent to the condition, with $v := z_2 - z_1$:

$$\forall i \in \llbracket 1, p \rrbracket, \quad \begin{cases} \text{if } m_i = 1 : B_{\mathbf{m}}A^{-1}P_i^T P_i v \cdot v > C_{\mathbf{m}}A^{-1}P_i^T P_i v \cdot v \\ \text{if } m_i = -1 : C_{\mathbf{m}}A^{-1}P_i^T P_i v \cdot v > B_{\mathbf{m}}A^{-1}P_i^T P_i v \cdot v \end{cases}.$$

More generally, given (y_1, y_2) of configuration \mathbf{m} , the configuration (y_1^*, y_2^*) of the next step (optimising in (y_1, y_2) with $\pi = \pi_{\mathbf{m}}$ fixed), is given by:

$$\forall i \in \llbracket 1, p \rrbracket, \quad m_i^* = \operatorname{sign} \left(v^T (B_{\mathbf{m}} - C_{\mathbf{m}}) A^{-1} P_i^T P_i v \right). \quad (15)$$

Notice that in particular, $y_1^* + y_2^* = z_1 + z_2$ since $B_{\mathbf{m}} + C_{\mathbf{m}} = A$.

Noticing that $J(y_1, y_2, \pi_{\mathbf{m}})$ is jointly convex and quadratic in y_1, y_2 , we first write:

$$\begin{aligned} J(y_1, y_2, \pi_{\mathbf{m}}) &= \sum_{\substack{i=1 \\ m_i=1}}^p \lambda_i \left(\frac{1}{2} \|P_i(z_1 - y_1)\|_2^2 + \frac{1}{2} \|P_i(z_2 - y_2)\|_2^2 \right) \\ &\quad + \sum_{\substack{i=1 \\ m_i=-1}}^p \lambda_i \left(\frac{1}{2} \|P_i(z_2 - y_1)\|_2^2 + \frac{1}{2} \|P_i(z_1 - y_2)\|_2^2 \right). \end{aligned}$$

Then

$$\begin{aligned} \nabla_{y_1} J(y_1, y_2, \pi_{\mathbf{m}}) &= \sum_{\substack{i=1 \\ m_i=1}}^p \lambda_i P_i^T P_i (y_1 - z_1) + \sum_{\substack{i=1 \\ m_i=-1}}^p \lambda_i P_i^T P_i (y_1 - z_2) \\ &= (B_{\mathbf{m}} + C_{\mathbf{m}})y_1 - B_{\mathbf{m}}z_1 - C_{\mathbf{m}}z_2. \end{aligned}$$

Using $B_{\mathbf{m}} + C_{\mathbf{m}} = A$, $y_1^* = A^{-1}(B_{\mathbf{m}}z_1 + C_{\mathbf{m}}z_2)$ and similarly $y_2^* = A^{-1}(B_{\mathbf{m}}z_2 + C_{\mathbf{m}}z_1)$.

The local optimality condition and the next step's configuration are obtained by expanding the squares in 4.4.1 and substituting y_1^*, y_2^* .

Theorem 4.4.2 shows that in order to study the local optima of E , the energy of (RP), one can study the graph of cell movements, i.e. the graph where the nodes are the admissible configurations \mathbf{m} and the edges are the next steps in the BCD algorithm. If a path ends at a configuration $\mathbf{m} \ni \{(1, \dots, 1), (0, \dots, 0)\}$, then the associated BCD algorithm will convergence towards a (strict) local optimum.

Below (Figure 15) is an example of the cell structure. The numerical experiments below can be reproduced using [version 1.4 of our repository](#).

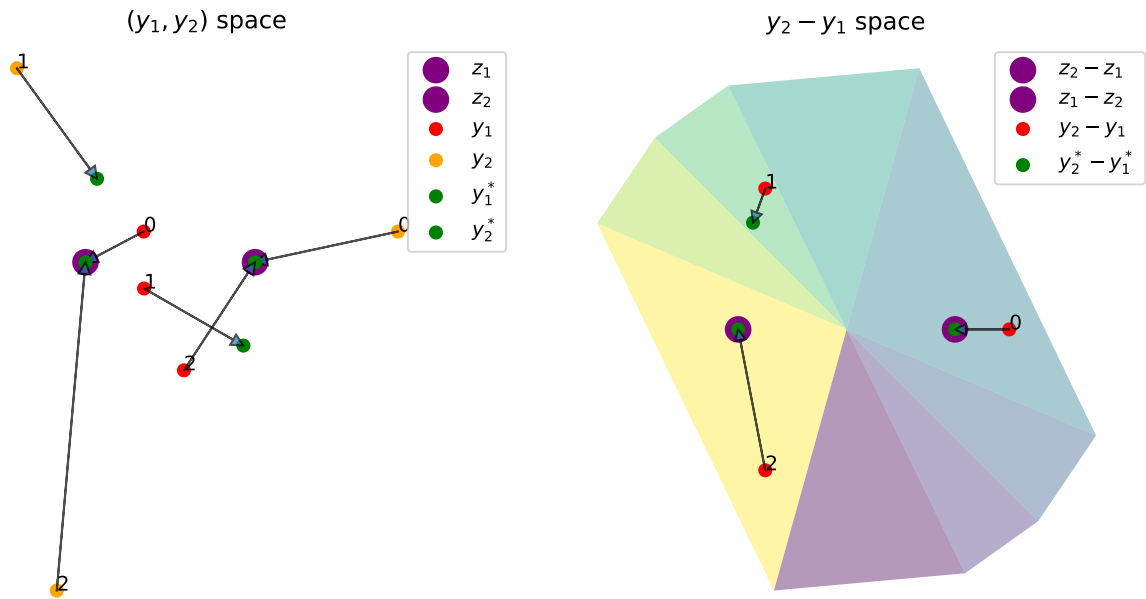


Figure 15: Visualisation of the cell structure for $p = 4$ in dimension 2. On the left, we view different initialisations (y_1, y_2) (in red and orange) and their corresponding BCD steps (y_1^*, y_2^*) , which should be compared to the original points (z_1, z_2) in purple. On the right, we view the cell structure depending on the position of $y_2 - y_1 \in \mathbb{R}^2$, since the cell conditions only depend on this difference (see (14)). We can see that in this example all cells are stable, thus there are three strict local optima of E in addition to the global optimum. The (y_1, y_2) pair number 0 is sent in one iteration to (z_2, z_1) , while the pair "1" is sent to a local optimum, and the pair "2" is sent to (z_1, z_2) .

A point of crucial importance is the number of stable cells (i.e. the number of local optima in the reconstruction problem) and the proportion of stable cells amongst all the cells.

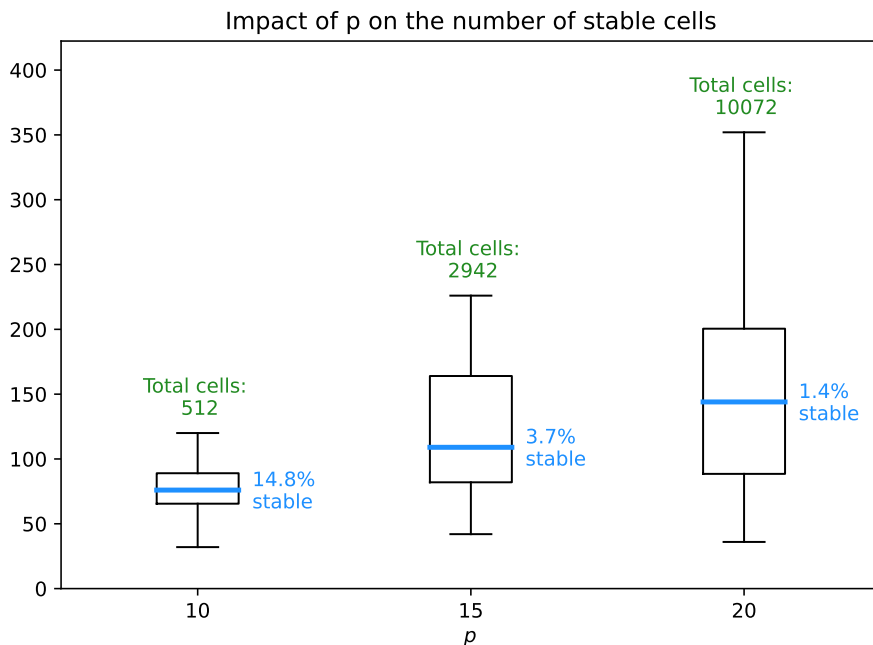


Figure 16: For $p = 10, 15, 20$, we count the number of local optima for a 2-point, 5-dimensional reconstruction problem (projecting on lines: $h = 1$). For each of the three settings we sampled the maps P_i and depicted the samples stable cells counts as box plots, with in blue the median ratio of stable cells. In green, above the box plots is the total number of cells (which only depends on p)

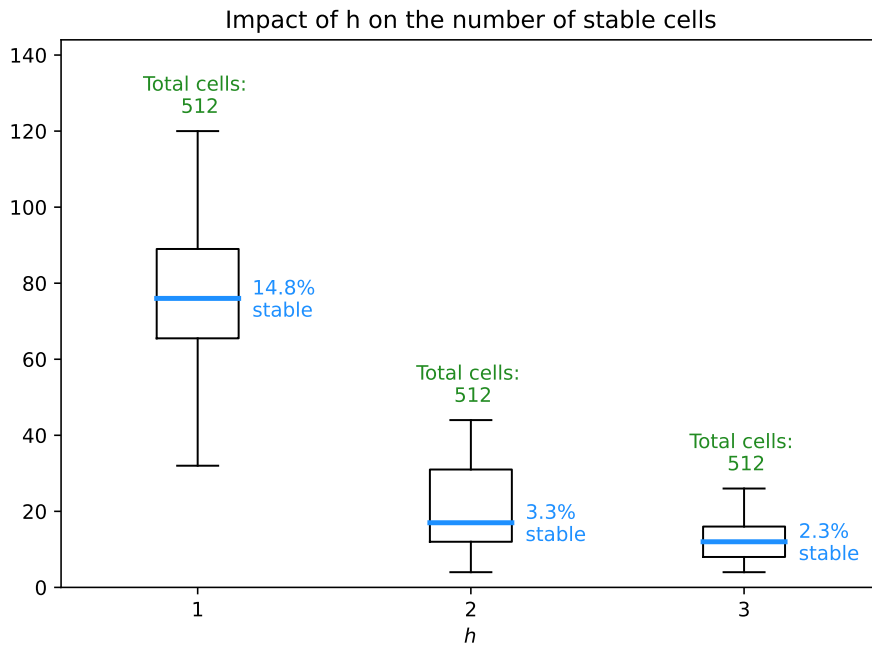


Figure 17: For $p = 10$ and the projection dimension $h = 1, 2, 3$, we count the number of local optima for a 2-point, 5-dimensional reconstruction problem. For each of the three settings we sampled the maps P_i and depicted the samples stable cells counts as box plots, with in blue the median ratio of stable cells. In green, above the box plots is the total number of cells (which only depends on p)

Furthermore, an important question is how many cells lead to the global optimum. In order to test this numerically, we fix a sample of the maps P_i from a previous experiment (mapping from \mathbb{R}^5 to \mathbb{R}), and compute all the next steps from each cell (thereby computing the cell movement graph):

- For a sample in dimension $d = 5$ and $p = 10$ projections, 70 cells out of 512 are stable, and taking a starting cell at random, the probability of reaching the global optimum is 3.5%.
- For a sample in dimension $d = 5$ and $p = 20$ projections, 190 cells out of 10 072 are stable, and taking a starting cell at random, the probability of reaching the global optimum is 33.5%.

Before establishing a more precise result on the movement between cells, we need a technical concentration lemma for random matrices. Indeed, in the case $d_i = 1$ with uniform weights $\lambda_i = 1/p$, the matrices A and even $B_{\mathbf{m}}$ can be seen as empirical covariance matrices: for instance $A = \frac{1}{p} \sum_{i=1}^p u_i u_i^T$, where $P_i = u_i^T$. Then with some random matrix theory, we can estimate how close the spectrum of A is to the real covariance matrix $\frac{1}{d} I_d$.

Lemma 4.4.3 — Concentration of the spectra of A and $B_{\mathbf{m}}$

If the $P_i = u_i^T$ are drawn using $u_i \sim \mathbb{P}$ uniform on \mathbb{S}^d , then with probability at least $1 - 4 \exp(-c\sqrt{d})$, with $c > 0$ a constant:

$$\forall \mu \in \text{sp}(A), \quad \frac{1}{d} - \frac{C}{\sqrt{dp}} \leq \mu \leq \frac{1}{d} + \frac{C}{\sqrt{dp}},$$

with $C > 0$ a numerical constant (empirically estimated to satisfy $0 < C \leq 5$), and similarly:

$$\forall \mu \in \text{sp}(B_{\mathbf{m}}), \quad \frac{1}{d} - \frac{C}{\sqrt{dp^+}} \leq \frac{p}{p^+} \mu \leq \frac{1}{d} + \frac{C}{\sqrt{dp^+}},$$

with $p^+ := \#\{i \in \llbracket 1, p \rrbracket \mid m_i = +1\}$ is the number of terms in the sum defining $B_{\mathbf{m}}$.

In order to prove this lemma we will use the following theorem, reproduced from [2]:

Theorem 4.4.4 — Adamczak et al. [2], Corollary 1

Let $A = \frac{1}{p} \sum_{i=1}^p u_i u_i^T$, with the $u_i \in \mathbb{R}^d$ isotropic random vectors.

Suppose further that for some $\psi, K > 0$ numerical constants:

- a) $\max_{i \in [1, p]} \sup_{y \in \mathbb{S}^d} \|u_i \cdot y\|_{\psi_1} \leq \psi$, with $\|X\|_{\psi_1} = \inf \left\{ C > 0 \mid \mathbb{E} \left[e^{|X|/C} \right] \leq 2 \right\}$.
- b) $\mathbb{P} \left(\max_{i \in [1, p]} (|u_i|/\sqrt{d}) > K \max(1, (p/d)^{1/4}) \right) \leq e^{-\sqrt{d}}$.

Then, with probability at least $1 - 2e^{-c\sqrt{d}}$, $c > 0$ being a numerical constant:

$$\forall \mu \in \text{sp}(A), \quad 1 - C_0(\psi + K)^2 \sqrt{\frac{d}{p}} \leq \mu \leq 1 + C_0(\psi + K)^2 \sqrt{\frac{d}{p}}.$$

A sufficient condition for a) and b) is for the u_i to satisfy a Poincaré inequality of constant L : $\mathbb{V}[f(u_i)] \leq L^2 \mathbb{E}[\|\nabla f(u_i)\|^2]$, for any compactly supported smooth function f .

We now use Theorem 4.4.4 to prove our Lemma 4.4.3:

We check that the Poincaré inequality is satisfied for $u_i = \sqrt{d}n_i$, where $n_i \sim \sigma$, the uniform law on the sphere \mathbb{S}^d . We drop the i index, considering that our variables are i.i.d..

We use the *spherical Poincaré-Wirtinger* inequality, which states that for any smooth function f such that $\mathbb{E}[f(n)] = 0$:

$$\int_{\mathbb{S}^d} f^2(x) d\sigma(x) \leq \frac{1}{d} \int_{\mathbb{S}^d} \|\nabla f(x)\|^2 d\sigma(x).$$

Let g smooth on \mathbb{R}^d and let $\bar{g} := \mathbb{E}[g(\sqrt{d}n)]$. We have: $\mathbb{V}[g(u)] = \mathbb{E}[(g(\sqrt{d}n) - \bar{g})^2]$,

and letting $f := x \mapsto g(\sqrt{d}x) - \bar{g}$, f is smooth and verifies $\mathbb{E}[f(n)] = 0$.

Thus by Poincaré-Wirtinger: $\int_{\mathbb{S}^d} f^2(x) d\sigma(x) \leq \frac{1}{d} \int_{\mathbb{S}^d} \|\nabla f(x)\|^2 d\sigma(x)$,

which is to say: $\mathbb{V}[g(u)] \leq \int_{\mathbb{S}^d} \|\nabla g(\sqrt{d}x)\|^2 d\sigma(x) = \mathbb{E}[\|\nabla g(u)\|^2]$.

Thus we can apply Theorem 4.4.4 with $L = 1$, yielding the first equation of Lemma 4.4.3, after dividing by d . Let $\mu \in \text{sp}(dA)$: by Theorem 4.4.4, we have with high probability $\mu \in \left[1 \pm C \sqrt{\frac{d}{p}}\right]$. Then for any $\nu \in \text{sp} A$, dividing by d yields $\nu \in \left[\frac{1}{d} \pm \frac{C}{\sqrt{dp}}\right]$.

The second equation is obtained by renormalising $B_{\mathbf{m}}$ using $\frac{p^+}{p}$.

By union bound, the probability of the intersection of the two events exceeds $1 - 2e^{-c\sqrt{d}} - 2e^{-c\sqrt{d}} = 1 - 4e^{-c\sqrt{d}}$.

Proof

Theorem 4.4.5 — Proportion of cells leading to the optimum

Suppose that the linear maps P_i are $P_i = u_i^T$, where the $u_i \sim \sigma$, the uniform law on \mathbb{S}^d . Further assume that the barycentre weights are uniform: $\lambda_i = \frac{1}{p}$.

Then there exists two constant $c, C > 0$ (taken from 4.4.4) such that with probability exceeding $1 - 4e^{-c\sqrt{d}}$, with $t = C\sqrt{\frac{d}{p}}$ we have:

If $t < 1/5$, and $\sqrt{\frac{p^+}{p}}$ or $\sqrt{\frac{p^-}{p}} > \frac{t + \sqrt{5t^2 - 6t + 2}}{2 - 4t}$, then the next BCD step will reach the optimal solution: $(y_1^*, y_2^*) = (z_1, z_2)$ or $(y_1^*, y_2^*) = (z_2, z_1)$ respectively.

In order to get more insight on the result, let us look at the first two terms of the asymptotic expansion in $t \rightarrow 0$ of our condition, where we consider d fixed and $p \rightarrow +\infty$ for our interpretation:

- At the limit $t \rightarrow 0$, all configurations \mathbf{m} such that $p^+ > p/2$ or $p^- > p/2$ are such that their next step (y_1^*, y_2^*) will be the global optimum. This means that any majority of $+$ or $-$ in the configuration will lead to a next step with only $+$ or $-$ respectively, which is to say the matching $+$ solution $(y_1^*, y_2^*) = (z_1, z_2)$ or the swap $-$ solution $(y_1^*, y_2^*) = (z_2, z_1)$.
- At the first order, the condition is $\sqrt{\frac{p^+}{p}}$ or $\sqrt{\frac{p^-}{p}} > \frac{1}{\sqrt{2}} + \frac{2 + \sqrt{2}}{4}t$. This means that in reality, the "vote" must have a slightly stricter majority, with a condition of the form $p^+ > p/2 + \alpha\sqrt{p}$ for a next step with only $+$'s and $p^- > p/2 + \alpha\sqrt{p}$ for a next iteration with only $-$'s.

We remind that $p^+ := \#\{i \in \llbracket 1, p \rrbracket \mid m_i = +1\}$ is the number of terms in the sum defining $B_{\mathbf{m}}$, with a similar definition of $p^- = p - p^+$ for $C_{\mathbf{m}}$.

Let us take a look at the trend of the necessary discriminative ratio:

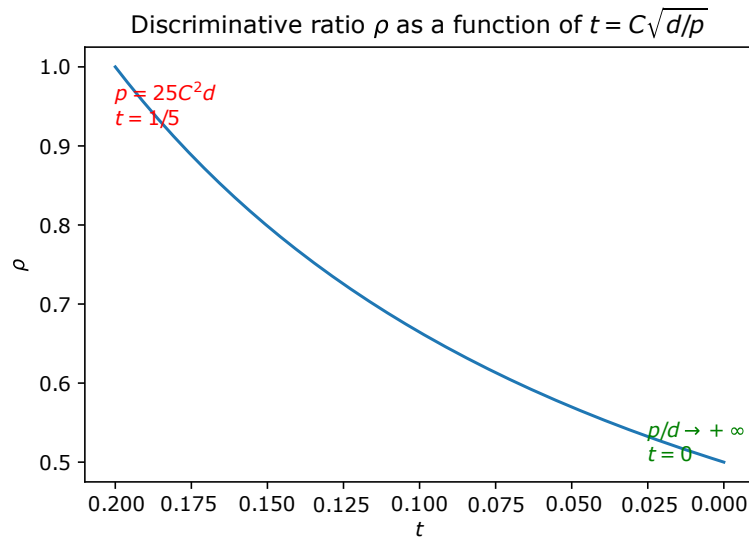


Figure 18: Trend of the discriminative ratio $\rho(t)$: ρ is the minimum value of p^+/p for the cell to lead directly to the global optimum. Precisely, we have $\rho(t) = \rho(C\sqrt{\frac{d}{p}}) = \left(\frac{t + \sqrt{5t^2 - 6t + 2}}{2 - 4t}\right)^2$. Notice the limit for $p/d \rightarrow +\infty$, the sufficient ratio is $1/2$, which is to say that a simple majority is enough for the cell to lead directly to the optimum. On the other end, the case $p = 25C^2d$ yields a necessary ratio of 1, thus only the cells corresponding to the real solution lead to the global optimum.

—*Step 1: Formulating the cell destination as an eigenvalue problem*

Let $(y_1, y_2) \in \mathbb{R}^d$ such that its configuration $\mathbf{m} \in \{\pm 1\}^p$. With $v = z_2 - z_1$, we study the configuration \mathbf{m}^* of the next step $(y_1^*, y_2^*) = \underset{Y=(y_1, y_2) \in \mathbb{R}^d \times \mathbb{R}^d}{\operatorname{argmin}} J(Y, \pi_{\mathbf{m}})$ (see 4.4.1). By (15):

$$\forall i \in \llbracket 1, p \rrbracket, \quad m_i^* = \operatorname{sign} \left(v^T (B_{\mathbf{m}} - C_{\mathbf{m}}) A^{-1} P_i^T P_i v \right)$$

This proof and the use of RMT relies on the intuition that the above condition depends essentially on the sign of the eigenvalues of $(B_{\mathbf{m}} - C_{\mathbf{m}}) A^{-1} P_i^T P_i$.

Let $i \in \llbracket 1, p \rrbracket$ and $p^+ := \#\{j \in \llbracket 1, p \rrbracket \mid m_j = +1\}$. We will consider the case $p^+ > \lfloor \frac{p}{2} \rfloor$, the case $p^+ < \lceil \frac{p}{2} \rceil$ can be studied symmetrically by focusing on $C_{\mathbf{m}}$ instead of $B_{\mathbf{m}}$. (We do not study the pathological case $p^+ = \frac{p}{2}$ for p even in this theorem).

Now re-write $M := (B_{\mathbf{m}} - C_{\mathbf{m}}) A^{-1} = 2B_{\mathbf{m}} A^{-1} - I$ using $B_{\mathbf{m}} - C_{\mathbf{m}} = A$. Then since A is real-symmetric, there exists $\Omega \in O_d(\mathbb{R})$ such that $\Omega^T A \Omega =: D$ is diagonal.

We operate this change of variable, first let $\tilde{P}_i := P_i \Omega$, notice $D = \frac{1}{p} \sum_{i=1}^p \tilde{P}_i^T \tilde{P}_i$.

Let $\widetilde{B}_{\mathbf{m}} := \Omega^T B_{\mathbf{m}} \Omega = \frac{1}{p} \sum_{m_j=1} \tilde{P}_i^T \tilde{P}_i$ and $\widetilde{M} := \Omega^T M \Omega$, we have $\widetilde{M} = 2\widetilde{B}_{\mathbf{m}} D^{-1} - I$.

The idea for what follows will be to estimate $\operatorname{sp} M = \operatorname{sp} \widetilde{M}$ using a quantitative approximation $D \approx I/d$, the true covariance matrix of $P_i^T \sim \sigma$.

Our objective will be to find some conditions under which $\forall \mu \in \operatorname{sp}(M), \mu > 0$. Assume that this is true, then since $u_i := P_i^T \in \mathbb{S}^d$, we have $\operatorname{sp}(M u_i u_i^T) = (u_i^T M u_i, \mathbf{0}_{d-1})$. Then under our assumption, $u_i^T M u_i > 0$, have almost-surely $v^T M P_i^T P_i v > 0$ and thus $m_i^* = +1$.

—*Step 2: Applying Random Matrix Theory*

Using Lemma 4.4.3, with probability exceeding $1 - 4e^{-c\sqrt{d}}$, one has:

$$\forall \mu \in \operatorname{sp}(A), \quad \mu \in \left[\frac{1}{d} \pm \frac{C}{\sqrt{dp}} \right] \quad \text{and} \quad \forall \mu \in \operatorname{sp}(B_{\mathbf{m}}), \quad \frac{p}{p^+} \mu \in \left[\frac{1}{d} \pm \frac{C}{\sqrt{dp^+}} \right].$$

We write $D = I/d - D_{\varepsilon}$, where $D_{\varepsilon} = \operatorname{diag}(\varepsilon_1, \dots, \varepsilon_d)$ and $|\varepsilon_j| \leq \frac{C}{\sqrt{dp}}$.

Then $D = \frac{1}{d}(I - dD_{\varepsilon})$, thus if $\|dD_{\varepsilon}\|_{\infty} < 1$, we have $D^{-1} = d \sum_{k=0}^{+\infty} d^k D_{\varepsilon}^k$.

Since $d \frac{C}{\sqrt{dp}} < 1 \Rightarrow \|dD_{\varepsilon}\|_{\infty} < 1$, we suppose our first condition $\mathcal{H}_1 := p > C^2 d$, i.e. $t < 1$.

At the first order in $d\varepsilon$, one therefore has $D^{-1} = dI + d^2 D_{\varepsilon}$. We will omit the error term in $d^3 \varepsilon^2$, which is negligible compared to the first term in $d^2 \varepsilon$. For better interpretation, we perform the computations at the first order in $t = C\sqrt{d}/p \rightarrow 0$, then start over with the complete non-asymptotic version of R in Step 4).

Now we can re-write $\widetilde{M} = 2\widetilde{B}_{\mathbf{m}}(dI + d^2 D_{\varepsilon}) - I$ and decompose with $N := 2d\widetilde{B}_{\mathbf{m}} - I$ and $R := 2d^2 \widetilde{B}_{\mathbf{m}} D_{\varepsilon}$. We have $\lambda_{\min}(M) \geq \lambda_{\min}(N) + \lambda_{\min}(R)$.

Since $\forall \mu \in \operatorname{sp}(B_{\mathbf{m}}), \frac{p}{p^+} \mu \in \left[\frac{1}{d} \pm \frac{C}{\sqrt{dp^+}} \right]$, we have $\forall \nu \in \operatorname{sp}(N), \nu \in \left[2\frac{p^+}{p} - 1 \pm 2C \frac{\sqrt{dp^+}}{p} \right]$,

and then $\lambda_{\min}(N) \geq 2\frac{p^+}{p} - 1 - 2C \frac{\sqrt{dp^+}}{p} =: \nu_0$.

Now for R : since $|\varepsilon_j| \leq \frac{C}{\sqrt{dp}}$, $\lambda_{\min}(R) \geq \lambda_{\min}\left(-2d^2 \widetilde{B}_{\mathbf{m}} \frac{C}{\sqrt{dp}}\right) = -\frac{2d^{3/2}C}{\sqrt{p}} \lambda_{\min}(B_{\mathbf{m}})$.

Finally, $\lambda_{\min}(R) \geq -2C\sqrt{d} \frac{p^+}{p^{3/2}} - 2dC^2 \frac{\sqrt{p^+}}{p^{3/2}} =: -r$.

Below Figure 19 is a representation of the intervals that we are studying:

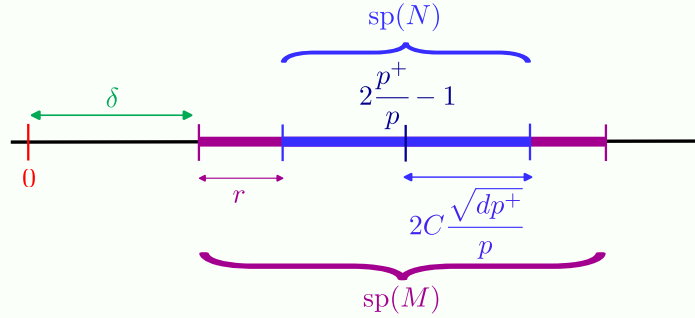


Figure 19: Visualisation of the spectra of the matrices N and M . The likely interval for the eigenvalues of N is drawn in blue, with the width of interval corresponding to the sampling error for the empirical covariance matrix $B_{\mathbf{m}}$. Similarly, when taking into account the added uncertainty (of width r) due to the empirical covariance matrix A , we obtain in purple the likely interval for the eigenvalues of $M = N + R$.

—Step 3: Conditions for $\text{sp}(M) \subset]0, +\infty[$

Our objective is to find a set of conditions under which the eigenvalues of M are all strictly positive, i.e. $\delta > 0$ for $\delta := \nu_0 - r$, our best lower bound.

We have for $x := \sqrt{\frac{p^+}{p}}$ and $t := C\sqrt{\frac{d}{p}}$, $\delta(x) = 2(1-t)x^2 - 2(t+t^2)x - 1$.

At the first order in t , we have $\delta(x) > 0$ for $x > x_+ := \frac{1}{\sqrt{2}} + \frac{2 + \sqrt{2}}{4}t + O(t^2)$.

This means that at the limit $t \rightarrow 0$, all values of $\frac{p^+}{p}$ will allow the desired result $\lambda_{\min}(M) > 0$, since the condition $x > 1/\sqrt{2}$ corresponds to $p^+ > \frac{p}{2}$.

However, at the first order, the condition is stricter: $\sqrt{\frac{p^+}{p}} > \frac{1}{\sqrt{2}} + \frac{2 + \sqrt{2}}{4}C\sqrt{\frac{d}{p}}$

—Step 4: Non-asymptotic condition

Only assuming $t < 1$, we can do our computations again with $R = 2d\widetilde{B}_{\mathbf{m}} \sum_{k=1}^{+\infty} d^k D_{\varepsilon}^k$, thus

$$\lambda_{\min}(R) \leq -2d \frac{t}{1-t} \lambda_{\max}(\widetilde{B}_{\mathbf{m}}).$$

Then since $\lambda_{\max}(\widetilde{B}_{\mathbf{m}}) \leq \frac{x^2}{d} + x \frac{t}{d}$, we have $\delta(x) = \frac{2-4t}{1-t}x^2 - \frac{2t}{1-t}x - 1$.

Finally, the condition $\delta(x) > 0$ is equivalent to $x > x^+ := \frac{t + \sqrt{5t^2 - 6t + 2}}{2 - 4t}$.

Now since $x \in]1/\sqrt{2}, 1]$, in order to have some x satisfying this condition, we need $\mathcal{H}_2 : x^+(t) \leq 1$, which is reached for $t \leq 1/5$.

4.5 Local Optima: general case

In this section we generalise the findings in the previous section to the case $L > 2$. This amounts to minimising the following energy:

$$J(Y, \pi) = \sum_{i=1}^p \lambda_i \sum_{(k,l) \in \llbracket 1, L \rrbracket^2} \pi_{k,l}^{(i)} \|P_i z_k - P_i y_l\|_2^2. \quad (16)$$

Where $Y \in \mathbb{R}^{L \times d}$ is the measure positions and $\pi \in \mathbb{U}_L^p$ is a list of valid OT matrices with uniform marginals.

Theorem 4.5.1 — Cell structure for any $L > 2$

Let $Y = (y_1, \dots, y_L)^T$ a current coupling position.

For each $i \in \llbracket 1, p \rrbracket$, the associated optimal transport matrix π_i is a normalised permutation matrix P_{σ_i}/L (recall $(P_{\sigma_i})_{k,l} = \mathbf{1}(\sigma_i(k) = l)$), where σ_i is determined by:

$$\sigma_i \in \operatorname{argmin}_{\sigma \in \mathfrak{S}_L} \sum_{k=1}^L \|P_i z_k - P_i y_{\sigma(k)}\|_2^2. \quad (17)$$

Note that if there are multiple solutions, any convex combination of them is also optimal.

Assuming that Y is such that each of the σ_i are unique, we define the *matching configuration* of Y by the list of its permutations: $\mathbf{m} = (\sigma_1, \dots, \sigma_p) \in \mathfrak{S}_L^p$.

Like the 2-point case, the mapping $\begin{cases} \mathbb{R}^{L \times d} & \longrightarrow & \mathfrak{S}_L^p \\ Y & \longmapsto & \mathbf{m}(Y) \end{cases}$ is constant by parts where the regions of $\mathbb{R}^{L \times d}$ with the same configuration are determined by a set of linear constraints: given $\mathbf{m} = (\sigma_1, \dots, \sigma_p) \in \mathfrak{S}_L^p$, the associated set of Y positions with configuration \mathbf{m} is solution to the set of linear inequations:

$$\forall i \in \llbracket 1, p \rrbracket, \quad \forall \sigma \in \mathfrak{S}_L \setminus \{\sigma_i\}, \quad \sum_{k=1}^L P_i^T P_i z_k \cdot y_{\sigma_i(k)} > \sum_{k=1}^L P_i^T P_i z_k \cdot y_{\sigma(k)} \quad (18)$$

In order to prove (17), we apply the Birkhoff-Von-Neumann theorem, which states the solutions $(\pi_1^*, \dots, \pi_p^*)$ minimising the linear cost $\pi \mapsto J(Y, \pi)$ over valid OT matrices are convex combinations of (normalised) permutation matrices.

Let $i \in \llbracket 1, p \rrbracket$. One may write $\pi_i^* = \sum_{\sigma \in \mathfrak{S}_L} a_\sigma P_\sigma / L$, with $a \in \Sigma_{L!}$.

Then one may re-write the i -th term in the energy J as $\frac{\lambda_i}{L} \sum_{\sigma \in \mathfrak{S}_L} a_\sigma \sum_{k=1}^L \|P_i z_k - P_i y_{\sigma(k)}\|_2^2$, hence (17).

Note that if a minimiser of (17) is not unique, then any $a \in \Sigma_{L!}$ whose support corresponds to the minimisers defines an optimal OT matrix $\pi_i^* = \sum_{\sigma \in \mathfrak{S}_L} a_\sigma P_\sigma / L$.

Finally, (18) is obtained by re-writing (17) by expanding the square and considering the case of a unique minimiser.

Similarly to the two-point case, we can predict the next BCD operation based on the measure positions' configuration:

Theorem 4.5.2 — Cell movement

Consider $Y \in \mathbb{R}^{L \times d}$ of configuration $\mathbf{m} = (\sigma_1, \dots, \sigma_p)$, and let $\pi_{\mathbf{m}}$ the associated OT matrices.

The next optimal positions $Y^* \in \underset{(y_1, \dots, y_L) \in \mathbb{R}^{L \times d}}{\operatorname{argmin}} J(Y, \pi_{\mathbf{m}})$ are:

$$y_k^* = A^{-1} \left(\sum_{i=1}^p P_i^T P_i z_{\sigma_i^{-1}(k)} \right) \quad (19)$$

The proof uses the same computations as 4.4.2, deriving the gradients of sums of quadratic forms in the y_k .

Unfortunately, unlike the 2-point case, there is no simplification that can be done, and the computation of the next configuration $\mathbf{m}(Y^*)$ relies on computing p OT problems: one for each permutation σ_i in the configuration.

5 The Blind Generalised Wasserstein Barycentre Problem

5.1 Problem Description

We now consider a generalised version of the GWB problem:

$$\operatorname{argmin}_{\substack{\gamma \in \mathcal{P}_2(\mathbb{R}^d) \\ P_i \in \mathcal{M}_{d_i, d}(\mathbb{R})}} \sum_{i=1}^p \lambda_i W_2^2(\nu_i, P_i \# \gamma) \quad (\text{BGWB})$$

This generalisation consists in adding a degree of freedom: the linear maps P_i , which we had considered fixed or random beforehand. We called this new problem "Blind" since the optimisation has to determine itself a possible correspondence between a coupling measure and the input measures.

5.2 Theoretical Properties of BGWB

5.2.1 Non-Convexity

First note that the (BGWB) energy,

$$J(Y, (\pi_i), (P_i)) = \sum_{i=1}^p \lambda_i \sum_{(k,l) \in \llbracket 1, K_i \rrbracket \times \llbracket 1, L \rrbracket} \|x_k^{(i)} - P_i y_l\|_2^2 \pi_{k,l}^{(i)} \quad (\text{J})$$

is separately convex (even quadratic or linear) in Y, p_i and P_i .

In order to gauge the computational difficulty of the problem, we shall use the closed form for the P_i (21) and study the partial problem with γ fixed:

$$\min_{P \in \mathcal{M}_{h,d}(\mathbb{R})} W_2^2(\nu, P \# \gamma) = \min_{P \in \mathcal{M}_{h,d}(\mathbb{R})} \min_{\pi \in \Pi(a,b)} M \cdot \pi = \min_{\pi \in \Pi(a,b)} \min_{P \in \mathcal{M}_{h,d}(\mathbb{R})} M \cdot \pi, \quad (20)$$

where $M_{k,l} = \|x_k - P y_l\|_2^2$.

Then for $\pi \in \Pi(a,b)$, by (21), $P^* := \operatorname{argmin}_{P \in \mathcal{M}_{h,d}(\mathbb{R})} = X^T \pi D$, with $D := Y \left(\sum_{l=1}^L b_l y_l y_l^T \right)^{-1} \in \mathcal{M}_{L,d}(\mathbb{R})$.

The associated OT matrix is $M_{k,l} = \|x_k - P^* y_l\|_2^2 = \|x_k - X^T \pi D y_l\|_2^2$, and (20) reads:

$$\min_{\pi \in \Pi(a,b)} \sum_{(k,l) \in \llbracket 1, K_i \rrbracket \times \llbracket 1, L \rrbracket} \|x_k - X^T \pi D y_l\|_2^2 \pi_{k,l},$$

which is a cubic problem with linear constraints, thus (potentially highly) non-convex in general. Further note that this numerically difficult problem is rendered even more complex by the subsequent optimisation in the positions Y .

A classical way of solving this cubic problem would be to alternate optimisation on π and a surrogate variable π' :

$$\min_{\pi, \pi' \in \Pi(a,b)} \sum_{(k,l) \in \llbracket 1, K_i \rrbracket \times \llbracket 1, L \rrbracket} \|x_k - X^T \pi' D y_l\|_2^2 \pi_{k,l},$$

yielding a linear program in π and a quadratic program in π' .

5.2.2 Invariants

In broad terms, the (BGWB) energy is invariant to any invertible linear transformation in the coupling space.

For any invertible $R \in \mathcal{M}_d(\mathbb{R})$, we have $J(YR^T, (\pi_i), (P_i R^{-1})) = J(Y, (\pi_i), (P_i))$. Essentially, the energy value remains the same if the points y_l and the maps P_i are transformed reciprocally.

This is proven by $P_i R^{-1} R y_l = P_i y_l$, or in a Wasserstein space viewpoint by the equally simple remark: $W_2^2(\nu, P\#\gamma) = W_2^2(\nu, (PR^T)\#(R^{-1}\#\gamma))$.

It is up to debate whether breaking this invariant would be beneficial in practice. One way of going about this could be to impose P_1 to be an augmented identity matrix, and only optimise on P_2, \dots, P_p .

5.3 Extending the Gradient Descent Solvers

In order to adapt our Gradient Descent-based methods for (GWB) to (BGWB) (Algorithm 4 and Algorithm 5), we only need to compute the gradients in the linear maps P_i .

5.3.1 Gradient in P

Using the same notations as §3.2.1:

We use mixed convention, where for $(\alpha, \beta) \in \llbracket 1, d_i \rrbracket \times \llbracket 1, d \rrbracket$, $\frac{\partial J}{\partial M}$ and $\frac{\partial M}{\partial P_{\alpha, \beta}}$ have the shape (K, L) .

Using the chain rule³ we can compute $\frac{\partial J}{\partial P_{\alpha, \beta}} = \text{Tr} \left(\begin{bmatrix} \frac{\partial J}{\partial M} \\ \frac{\partial M}{\partial P_{\alpha, \beta}} \end{bmatrix}^T \right)$, where:

- We have computed $\frac{\partial J}{\partial M} = \pi^*$
- Recall $M_{k,l} = \|x_k - P y_l\|_2^2$. We compute $\frac{\partial M}{\partial P_{\alpha, \beta}}$ by computing⁴ $\frac{\partial M_{k,l}}{\partial P} = -2x_k y_l^T + 2P y_l y_l^T$

As discussed earlier with the theoretical properties of (BGWB), we normalise the linear maps P_i : at each GD iteration, we project each row of each P_i onto the L^2 sphere \mathbb{S}^d , which we will see as projecting onto the set of $d_i \times d$ matrices with normalised rows $R_{d_i, d}$.

Similarly to the positions y_l in the (GWB) GD discussion, we apply of Shapiro's Theorem (3.2.1) on the primal problem, allowing us to plug-in the optimal value of y_l which yields a sub-gradient in P .

5.3.2 GD for BGWB Algorithm

We present our GD solver for the (BGWB) problem. Note that technically, the method is Projected Alternated Gradient Descent.

³See *The Matrix Cookbook* [16], §2.8.1 (126)

⁴See *The Matrix Cookbook* [16], §2.4.1 (70) and §2.4.2 (77)

Algorithm 8: (BGWB) resolution with Gradient Descent

Data: Input measure points $(X_i)_{i \in \llbracket 1, p \rrbracket} \in \prod_{i=1}^p \mathcal{M}_{K_i, d_i}(\mathbb{R})$ and weights $(a_i)_{i \in \llbracket 1, p \rrbracket} \in \prod_{i=1}^p \Sigma_{K_i}$
 Number of barycentre points L , barycentric coefficients $\lambda \in \Sigma_p$, precision ε ,
 iterations N , learning rate η , l.r. decay ρ .

Result: Barycentre positions $Y \in \mathcal{M}_{L, d}(\mathbb{R})$, barycentre weights $b \in \Sigma_L$,
 and linear maps $(P_i)_{i \in \llbracket 1, p \rrbracket} \in \prod_{i=1}^p \mathbb{R}^{d_i}$.

```

1 Initialisation: Draw  $Y \in \mathcal{M}_{L, d}(\mathbb{R})$ ,  $b \in \Sigma_L$ ,  $(P_i)_{i \in \llbracket 1, p \rrbracket} \in \prod_{i=1}^p \mathbb{R}^{d_i}$  and let  $J_0 := +\infty$ ;
2 for  $t \in \llbracket 1, N \rrbracket$  do
3   for  $i \in \llbracket 1, p \rrbracket$  do
4     Compute  $J^{(i)} = \min_{\pi_i \in \Pi(a_i, b)} M_i \cdot \pi_i$  where  $M_{k, l}^{(i)} = \|x_k^{(i)} - P_i y_l\|_2^2$ ;
5   end
6   Compute the loss  $J_t = \sum_{i=1}^p \lambda_i J^{(i)}$  and its gradients w.r.t.  $Y, b, P_i$ ;
7   for  $i \in \llbracket 1, p \rrbracket$  do
8     Step the linear map  $P_i$ :  $P_i \leftarrow \Pi_{\mathbb{R}^{d_i, d}} \left( P_i - \rho^t \eta \frac{\partial J_t}{\partial P_i} \right)$ ;
9   end
10  Step the positions  $Y$ :  $Y \leftarrow Y - \rho^t \eta \frac{\partial J_t}{\partial Y}$ ;
11  Step the weights  $b$ :  $b \leftarrow \Pi_{\Sigma_L} \left( b - \rho^t \eta \frac{\partial J_t}{\partial b} \right)$ ;
12  if  $J_{t-1} - J_t < \varepsilon$  then
13    Declare convergence and terminate.
14  end
15 end

```

5.3.3 SGD for BGWB Algorithm

Below is a stochastic variant of the GD solver for BGWB (Algorithm 8), which is technically Stochastic Projected Alternated Gradient Descent.

Algorithm 9: (BGWB) resolution with Stochastic Gradient Descent

Data: Input measure points $(X_i)_{i \in \llbracket 1, p \rrbracket} \in \prod_{i=1}^p \mathcal{M}_{K_i, d_i}(\mathbb{R})$ and weights $(a_i)_{i \in \llbracket 1, p \rrbracket} \in \prod_{i=1}^p \Sigma_{K_i}$
 Number of barycentre points L , barycentric coefficients $\lambda \in \Sigma_p$, precision ε ,
 iterations N , learning rate η and l.r. decay ρ .

Result: Barycentre positions $Y \in \mathcal{M}_{L, d}(\mathbb{R})$, barycentre weights $b \in \Sigma_L$,

and linear maps $(P_i)_{i \in \llbracket 1, p \rrbracket} \in \prod_{i=1}^p \mathbb{R}^{d_i}$.

- 1 **Initialisation:** Draw $Y \in \mathcal{M}_{L, d}(\mathbb{R})$, $b \in \Sigma_L$, $(P_i)_{i \in \llbracket 1, p \rrbracket} \in \prod_{i=1}^p \mathbb{R}^{d_i}$ and let $J_0 := +\infty$;
- 2 **for** $t \in \llbracket 1, N \rrbracket$ **do**
- 3 **for** $_ \in \llbracket 1, p \rrbracket$ **do**
- 4 Draw $i \sim \mathbb{D}$;
- 5 Compute $J^{(i)} = \min_{\pi_i \in \Pi(a_i, b)} M_i \cdot \pi_i$ where $M_{k, l}^{(i)} = \|x_k^{(i)} - P_i y_l\|_2^2$;
- 6 **end**
- 7 Compute the loss $J_t = \sum_{i=1}^p J^{(i)}$ and its gradients w.r.t. $Y, b, (P_i)_{i \in \llbracket 1, p \rrbracket}$;
- 8 **for** $i \in \llbracket 1, p \rrbracket$ **do**
- 9 Step the linear map $P_i : P_i \leftarrow \Pi_{\mathbb{R}^{d_i, d}} \left(P_i - \rho^t \eta \frac{\partial J_t}{\partial P_i} \right)$;
- 10 **end**
- 11 Step the positions $Y : Y \leftarrow Y - \rho^t \eta \frac{\partial J_t}{\partial Y}$;
- 12 Step the weights $b : b \leftarrow \Pi_{\Sigma_L} \left(b - \rho^t \eta \frac{\partial J_t}{\partial b} \right)$;
- 13 **if** $J_{t-1} - J_t < \varepsilon$ **then**
- 14 Declare convergence and terminate.
- 15 **end**
- 16 **end**

5.3.4 Visual Experiments

First, notice (see [Figure 20](#)) that with the parameters b, P let constant, the GD solver on the (GWB) problem finds a similar solution to `free_support_barycenter` on (GWB'), since the two problems are equivalent:

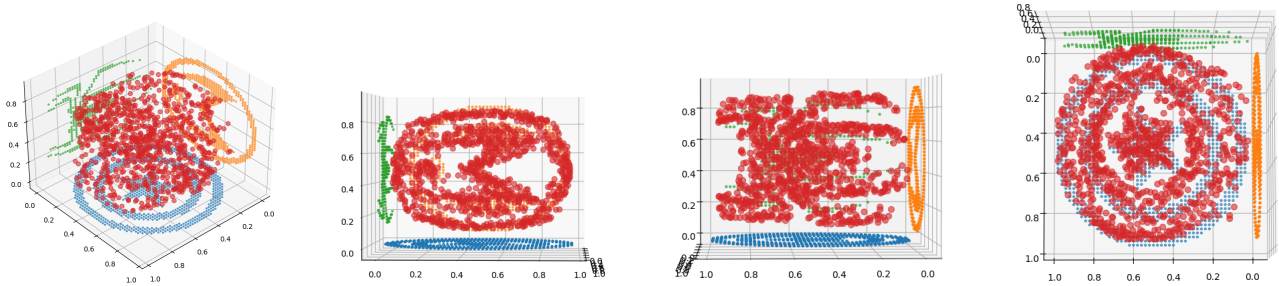


Figure 20: (GWB) resolution using GD on a toy dataset (fixed maps (P_i)).

And below ([Figure 21](#)) is the result of the GD solution for the (BGWB) problem (thus also optimising the linear maps P_i):

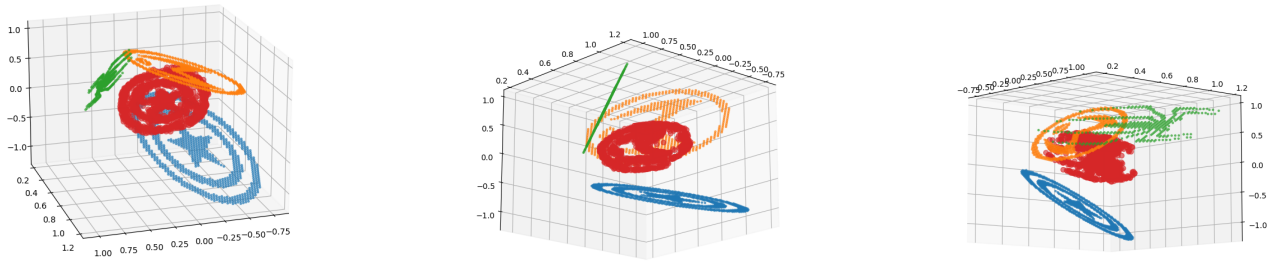


Figure 21: (BGWB) resolution using GD on a toy dataset.

5.4 BCD resolution

5.4.1 Closed form in P

Using the same notations as [§3.3](#):

The structure of the problem makes it sufficient to find a closed form solution of the following problem:

$$\operatorname{argmin}_{P \in \mathcal{M}_{d',d}(\mathbb{R})} \sum_{(k,l) \in \llbracket 1,K \rrbracket \times \llbracket 1,L \rrbracket} \|x_k - Py_l\|_2^2 \pi_{k,l}.$$

Note that the energy (denoted as J') above is convex and quadratic in P . We compute:

$$\nabla J'(P) = 2 \sum_{k,l} \pi_{k,l} (Py_l y_l^T - x_k y_l^T), \text{ yielding an optimality condition } PB = C$$

$$\text{where } B := \sum_{k=1}^K \sum_{l=1}^L \pi_{k,l} y_l y_l^T = \sum_{l=1}^L b_l y_l y_l^T \text{ and } C := \sum_{k=1}^K \sum_{l=1}^L \pi_{k,l} x_k y_l^T.$$

Assuming the invertibility of B we have a unique solution:

$$P^* = \left(\sum_{k=1}^K \sum_{l=1}^L \pi_{k,l} x_k y_l^T \right) \left(\sum_{l=1}^L b_l y_l y_l^T \right)^{-1}. \quad (21)$$

Since the $\pi_{k,l} \geq 0$, B is symmetric positive semi-definite. Then B is invertible iff $\text{Span}(y_l)_{l \in \llbracket 1, L \rrbracket} = \mathbb{R}^d$, which is a weak assumption in practice.

Note that unlike GD, we cannot enforce normalisation on the P_i apart from their initialisation.

For the matrix form, we use $C = X^T \pi Y$ and $B = Y^T b Y$.

5.4.2 BCD for BGWB Algorithm

Algorithm 10: (BGWB) resolution with Block-Coordinate Descent

Data: Input measure points $(X_i)_{i \in \llbracket 1, p \rrbracket} \in \prod_{i=1}^p \mathcal{M}_{K_i, d_i}(\mathbb{R})$ and weights $(a_i)_{i \in \llbracket 1, p \rrbracket} \in \prod_{i=1}^p \Sigma_{K_i}$
 Barycentre weights $b \in \Sigma_L$, barycentric coefficients $\lambda \in \Sigma_p$, precision ε , iterations N .

Result: Barycentre positions $Y \in \mathcal{M}_{L, d}(\mathbb{R})$, and linear maps $(P_i)_{i \in \llbracket 1, p \rrbracket} \in \prod_{i=1}^p \mathbb{R}^{d_i}$.

```

1 Initialisation: Draw  $Y \in \mathcal{M}_{L, d}(\mathbb{R})$ ,  $(P_i)_{i \in \llbracket 1, p \rrbracket} \in \prod_{i=1}^p \mathbb{R}^{d_i}$  and let  $J_0 := +\infty$ ;
2 for  $t \in \llbracket 1, N \rrbracket$  do
3   for  $i \in \llbracket 1, p \rrbracket$  do
4     Compute the OT distance matrix  $M_{k,l}^{(i)} = \|x_k^{(i)} - P_i y_l\|_2^2$ ;
5     Compute the OT map  $\pi_i$  by solving  $\min_{\pi_i \in \Pi(a_i, b)} M_i \cdot \pi_i$ ;
6   end
7   for  $l \in \llbracket 1, L \rrbracket$  do
8     Update  $y_l$ : compute  $y_l = \left( \sum_{i=1}^p \lambda_i b_l P_i^T P_i \right)^{-1} \left( \sum_{i=1}^p \lambda_i (\pi_{\cdot, l}^{(i)})^T X_i P_i \right)$ ;
9   end
10  for  $i \in \llbracket 1, p \rrbracket$  do
11    Update  $P_i$ : compute  $P_i = \left( X^T \pi_i Y \right) \left( Y^T b Y \right)^{-1}$ ;
12  end
13  Compute the energy  $J_t = \sum_{i=1}^p \lambda_i M_i \cdot \pi_i$ ;
14  if  $J_{t-1} - J_t < \varepsilon$  then
15    Declare convergence and terminate.
16  end
17 end

```

6 Perspectives and Conclusion

6.1 Perspectives

There are two main axes of progression for this project:

- First of all, the work on the reconstruction problem and its ties to the Sliced Wasserstein distance still has more potential. Substantial work has been done in order to attempt to compute the law of cell movement (i.e., starting from a configuration \mathbf{m} , what is the probability of the next configuration being \mathbf{m}' ?). This question is still unanswered due to random matrix theory technical obstacles, but computing the aforementioned law would allow a significantly better understanding of local optima. A crucial question to answer is the existence of local optima for the Sliced Wasserstein distance (with p projections and at the limits $p \rightarrow +\infty$, which our Reconstruction Problem analysis is starting to shed light on).
- The second perspective is the continuation of the Blind version of GWB: several questions remain. What constraints should we impose upon the P_i , and is it beneficial to drop the invariants? Can we derive a dual and a better understanding of the solutions? An application that we considered and couldn't put together in time was Domain Adaptation: the BGWB problem could be a way of using non-labelled data, by learning a heterogeneous mapping that preserves labels.

6.2 Conclusion

In this internship we put together effective solvers for the (Blind) Generalised Wasserstein Barycentre problem. These numerical solutions are available and well documented on our repository, and are on their way to the open source module POT.

On the theoretical side, we developed an extensive understanding of the set of solutions of a reconstruction problem, which is a restriction of the GWB problem. Furthermore, these insights have led to a unique approach of studying the local optima of Sliced Wasserstein distances, using Random Matrix Theory.

During my stay at the laboratoire MAP5 I had the pleasure of working with Julie Delon and Rémy Flamary who had the kindness of sharing valuable advice on my work and my career as a future PhD student. Beyond the abundant academic help that they have generously provided, they allowed me to grow further as a human being and an aspiring researcher. At MAP5 I also met with the doctoral students and the other permanent researchers, and I am overjoyed to have been welcomed so warmly into what will be my home for the next three years. I am eager to continue this work and much more during my PhD at MAP5, under the guidance of Julie Delon, Agnès Desolneux and Rémy Flamary.

References

- [1] Isabelle Abraham et al. “Tomographic reconstruction from a few views: a multi-marginal optimal transport approach”. In: *Applied Mathematics & Optimization* 75.1 (2017), pp. 55–73.
- [2] Radosław Adamczak et al. *Sharp bounds on the rate of convergence of the empirical covariance matrix*. 2010. DOI: [10.48550/ARXIV.1012.0294](https://doi.org/10.48550/ARXIV.1012.0294). URL: <https://arxiv.org/abs/1012.0294>.
- [3] Martial Agueh and Guillaume Carlier. “Barycenters in the Wasserstein Space”. In: *SIAM Journal on Mathematical Analysis* 43.2 (2011), pp. 904–924. DOI: [10.1137/100805741](https://doi.org/10.1137/100805741). eprint: <https://doi.org/10.1137/100805741>. URL: <https://doi.org/10.1137/100805741>.
- [4] Jason M. Altschuler and Enric Boix-Adserà. “Wasserstein Barycenters Are NP-Hard to Compute”. In: *SIAM Journal on Mathematics of Data Science* 4.1 (2022), pp. 179–203. DOI: [10.1137/21M1390062](https://doi.org/10.1137/21M1390062). eprint: <https://doi.org/10.1137/21M1390062>. URL: <https://doi.org/10.1137/21M1390062>.
- [5] Pierre Bernhard and Alain Rapaport. “On a theorem of Danskin with an application to a theorem of Von Neumann-Sion”. In: *Nonlinear Analysis: Theory, Methods & Applications* 24.8 (1995), pp. 1163–1181. ISSN: 0362-546X. DOI: [https://doi.org/10.1016/0362-546X\(94\)00186-L](https://doi.org/10.1016/0362-546X(94)00186-L). URL: <https://www.sciencedirect.com/science/article/pii/0362546X9400186L>.
- [6] J. Frédéric Bonnans and Alexander Shapiro. “Optimization Problems with Perturbations: A Guided Tour”. In: *SIAM Rev.* 40.2 (1998), pp. 228–264. DOI: [10.1137/S0036144596302644](https://doi.org/10.1137/S0036144596302644). URL: <https://doi.org/10.1137/S0036144596302644>.
- [7] Nicolas Bonneel et al. “Sliced and radon wasserstein barycenters of measures”. In: *Journal of Mathematical Imaging and Vision* 51.1 (2015), pp. 22–45.
- [8] Guillaume Carlier and Ivar Ekeland. “Matching for teams”. In: *Economic Theory* 42 (Feb. 2010), pp. 397–418. DOI: [10.1007/s00199-008-0415-z](https://doi.org/10.1007/s00199-008-0415-z).
- [9] Marco Cuturi and Arnaud Doucet. “Fast Computation of Wasserstein Barycenters”. In: (2013). DOI: [10.48550/ARXIV.1310.4375](https://doi.org/10.48550/ARXIV.1310.4375). URL: <https://arxiv.org/abs/1310.4375>.
- [10] John M Danskin. “The theory of max-min, with applications”. In: *SIAM Journal on Applied Mathematics* 14.4 (1966), pp. 641–664.
- [11] Julie Delon, Nathaël Gozlan, and Alexandre Saint-Dizier. *Generalized Wasserstein barycenters between probability measures living on different subspaces*. 2021. DOI: [10.48550/ARXIV.2105.09755](https://doi.org/10.48550/ARXIV.2105.09755). URL: <https://arxiv.org/abs/2105.09755>.
- [12] Ishan Deshpande, Ziyu Zhang, and Alexander Schwing. “Generative Modeling Using the Sliced Wasserstein Distance”. In: June 2018, pp. 3483–3491. DOI: [10.1109/CVPR.2018.00367](https://doi.org/10.1109/CVPR.2018.00367).
- [13] Aingeru Fernandez-Bertolin, Philippe Jaming, and Karlheinz Grochenig. “Determining point distributions from their projections”. In: July 2017, pp. 164–168. DOI: [10.1109/SAMPTA.2017.8024381](https://doi.org/10.1109/SAMPTA.2017.8024381).
- [14] Rémi Flamary et al. “POT: Python Optimal Transport”. In: *Journal of Machine Learning Research* 22.78 (2021), pp. 1–8. URL: <http://jmlr.org/papers/v22/20-451.html>.
- [15] Soheil Kolouri et al. “Sliced Wasserstein Auto-Encoders”. In: *International Conference on Learning Representations*. 2019. URL: <https://openreview.net/forum?id=H1xaJn05FQ>.
- [16] K. B. Petersen and M. S. Pedersen. *The Matrix Cookbook*. Version 20081110. Oct. 2008. URL: <http://www2.imm.dtu.dk/pubdb/p.php?3274>.
- [17] Gabriel Peyré and Marco Cuturi. “Computational Optimal Transport”. In: (2018). DOI: [10.48550/ARXIV.1803.00567](https://doi.org/10.48550/ARXIV.1803.00567). URL: <https://arxiv.org/abs/1803.00567>.
- [18] Filippo Santambrogio. “Optimal transport for applied mathematicians”. In: *Birkäuser, NY* 55.58-63 (2015), p. 94.