

Optimal transport theory

Student : Toby van Gastelen

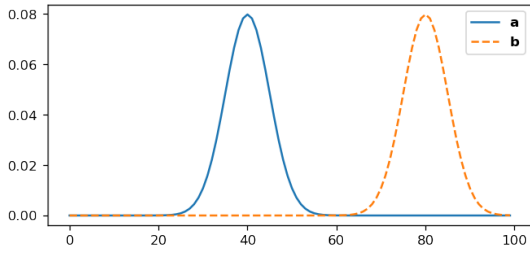
Supervisors : Dr. Saad Yalouz - Pr. Dr. Lucas Visscher
Vrije Universiteit Amsterdam

I. INTRODUCTION

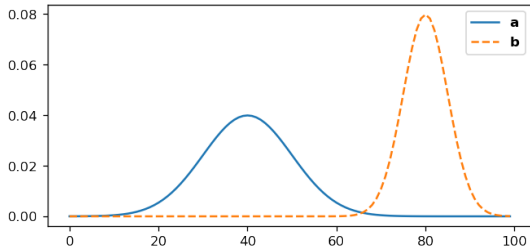
When we want to compare two vectors of the same length, e.g. \mathbf{a} and \mathbf{b} , we usually resort to the l^2 -norm of the difference:

$$\Delta = \|\mathbf{a} - \mathbf{b}\|_2, \quad (1)$$

also known as the Euclidian distance. Larger Δ would in this case imply a greater dissimilarity between the vectors. We will now demonstrate the problem with this approach. Let us first look at two normalized Gaussians with the same standard deviation σ , but different means μ and look at the l^2 -norm of the difference between the vectors containing the function evaluations for $\mathbf{x} = (0, 1, \dots, 99)^T$, $a(x_i) = a_i$ and $b(x_i) = b_i$ (FIG. 1a). We properly normalize each of the vectors to represent a discrete probability distribution on $[0, 99]$. Let us



(a) Example one: Two vectors \mathbf{a} and \mathbf{b} containing the Gaussian function evaluations for $\mu = 40$ and $\sigma = 5$, and $\mu = 80$ and $\sigma = 5$, respectively. In this case we have $\Delta = 0.3359$ and $W_2 = 40.00$.



(b) Example two: Two vectors \mathbf{a} and \mathbf{b} containing the Gaussian function evaluations for $\mu = 40$ and $\sigma = 10$, and $\mu = 80$ and $\sigma = 5$, respectively. In this case we have $\Delta = 0.2907$ and $W_2 = 40.31$.

Figure 1: Comparison between two Gaussians using Δ and W_2 .

also do the same thing, after increasing the σ of \mathbf{a} from 5 to 10 (FIG. 1b). According to the l^2 -norms of both examples, we find that the Gaussians with different σ (example two) are more similar as compared to the Gaussians with the same σ

(example one). This is caused by the overlap between \mathbf{a} and \mathbf{b} , due to the increased σ of \mathbf{b} . Let us now take a look what would happen if we calculate the Wasserstein distance W_2 for both examples (FIG. 1). Again, a larger value would imply a greater dissimilarity between the vectors. For this we used the **Python optimal transport**¹ (POT) library¹. Now we find that the difference between the two vectors is actually smaller in the first example, as compared to the second example, due to the similar shapes of the evaluated functions. We also know that when comparing two vectors containing the same shapes at different locations, W_2 equals the absolute distance between both of the centers of mass. Looking at example one this seems to be the case. Another nice property of the Wasserstein distance is that it allows us unambiguously compare vectors of different dimension. The procedure can easily be extended to matrices and higher dimensional objects.

II. WASSERSTEIN DISTANCE

Let us now dive a little bit deeper into this metric using mostly the theory presented in the book **Computational optimal transport**². We first define the Wasserstein distance as

$$W_p(\mathbf{a}, \mathbf{b}) = \left\{ \min_{\gamma} \sum_{ij} \gamma_{ij} M_{ij}(p) \right\}^{1/p}, \quad (2)$$

for which $a(\mathbf{x}_i) = a_i$ and $b(\mathbf{y}_i) = b_i$. W_p can be thought of as a measure to compare the resemblance of two shapes. Here we define the cost matrix \mathbf{M} as

$$M_{ij,p} = \|\mathbf{x}_i - \mathbf{y}_j\|_2^p \rightarrow \mathbf{M}_p. \quad (3)$$

This matrix encodes the penalty of moving mass from \mathbf{x}_i to \mathbf{y}_j . In the case of matrices $a(\mathbf{x}_i)$ and $b(\mathbf{y}_i)$ map the positions $\mathbf{x}_i/\mathbf{y}_i = (\text{row}, \text{column})$ to the value of the matrix entry, i.e. $a([x_{i1}, x_{i2}]) = a_i$ and $b([y_{i1}, y_{i2}]) = b_i$. For \mathbf{M} and γ we have

$$\dim(\mathbf{M}) = \dim(\gamma) = \dim(\mathbf{a}) \times \dim(\mathbf{b}), \quad (4)$$

i.e. $\mathbf{M}_{\dim(\mathbf{a}) \times \dim(\mathbf{b})}$ and $\gamma_{\dim(\mathbf{a}) \times \dim(\mathbf{b})}$. For this procedure to work both \mathbf{a} and \mathbf{b} need to abide by

$$\begin{aligned} \sum_i a_i &= 1, \\ \sum_i b_i &= 1 \\ a_i, b_i &\geq 0 \quad \forall i, \end{aligned} \quad (5)$$

such that \mathbf{a} and \mathbf{b} can be thought of as probability densities. The following constraints are put on γ :

$$\begin{aligned} \sum_{ij} \gamma_{ij} &= 1, \\ \gamma_{ij} &\geq 0, \forall i, j, \\ \sum_i \gamma_{ij} &= b_j, \\ \sum_j \gamma_{ij} &= a_i, \end{aligned} \quad (6)$$

which we will refer to as $\gamma \in \mu(\mathbf{a}, \mathbf{b})$. Now γ represents a probability density such that

$$P(\mathbf{x}_i, \mathbf{y}_j) = \gamma_{ij} \geq 0. \quad (7)$$

In reality γ encodes the permutations that need to be carried out to transform one vector into the other, which is then weighted by \mathbf{M}_p . Regarding the Wasserstein distance the following relations hold:

$$\begin{aligned} W_p(\mathbf{a}, \mathbf{b}) &= 0 \text{ iff } \mathbf{a} = \mathbf{b}, \\ W_p(\mathbf{a}, \mathbf{b}) &= W_p(\mathbf{b}, \mathbf{a}), \\ W_p(\mathbf{a}, \mathbf{b}) &\leq W_p(\mathbf{a}, \mathbf{c}) + W_p(\mathbf{c}, \mathbf{b}), \forall \mathbf{a}, \mathbf{b}, \mathbf{c} \end{aligned} \quad (8)$$

for $p > 1$, which makes that W_p classifies as a distance in the space of probability densities. In determining W_p the objective is to find the minimizer γ . For this we apply numerical optimization techniques. Algorithms to obtain W_p scale at best to order $O(n^2)$, where n is the dimension of the input vectors \mathbf{a} and \mathbf{b} . To make numerical optimization more feasible for large input data we introduce the addition of an entropic term $\Omega(\gamma)$ weighted by constant $\lambda > 0$:

$$\begin{aligned} W_p^\lambda(\mathbf{a}, \mathbf{b}) &= \{\min_{\gamma} \sum_{ij} (\gamma_{ij} M_{ij,p}) + \lambda \sum_{ij} \gamma_{ij} (\log(\gamma_{ij}) - 1)\}^{1/p} \\ &= \{\min_{\gamma} \sum_{ij} (\gamma_{ij} M_{ij,p}) + \lambda \Omega(\gamma)\}^{1/p}, \end{aligned} \quad (9)$$

resulting in

$$\gamma^*(\lambda) = \arg \min_{\gamma} \{\sum_{ij} (\gamma_{ij} M_{p,ij}) + \lambda \Omega(\gamma)\}. \quad (10)$$

This can in turn be optimized using the Sinkhorn algorithm. This algorithm scales to $O(n/\lambda^2)$. The algorithm is treated in more detail in SEC. III. For $p > 1$ the solution of γ (or approximate solution γ^*) is indepent of chosen p , i.e. the solution is always the same for $p > 1$. This is why there is no need to write the exponent anymore in EQ. 10. Increasing the parameter λ results in faster convergence to the solution but also results in a larger error, as compared to $\lambda = 0$. After numerically obtaining the solution for γ (or γ^*) we obtain W_p (or W_p^λ) by doing

$$W_p = (\sum_{ij} \gamma_{ij} M_{p,ij})^{1/p} = \text{Tr}(\mathbf{M}_p^T \gamma)^{1/p}, \quad (11)$$

omitting the entropic part.

III. SINKHORN ALGORITHM

Next, we will look at the Sinkhorn algorithm used to solve the minimization problem described in EQ. 10. To solve this problem we define the Lagrangian as

$$\begin{aligned} \mathcal{L}(\gamma, \mathbf{v}, \mathbf{u}) &= \sum_{ij} [\gamma_{ij} M_{ij,p} + \lambda \gamma_{ij} (\log(\gamma_{ij}) - 1)] + \\ &\quad \begin{cases} -\sum_i v_i [\sum_j (\gamma_{ij}) - a_i] - \sum_j u_j [\sum_i (\gamma_{ij}) - b_j] & \text{if } \gamma \in \mu(\mathbf{a}, \mathbf{b}) \\ \infty & \text{if } \gamma \notin \mu(\mathbf{a}, \mathbf{b}) \end{cases} \end{aligned} \quad (12)$$

where \mathbf{v} and \mathbf{u} contain the Lagrange multipliers that ensure that the constraints in EQ. 6 are satisfied.³ The $+\infty$ condition is added to ensure $\mathcal{L}(\gamma, \mathbf{v}, \mathbf{u})$ remains strictly convex in γ_{ij} for $\gamma_{ij} \in \mathbb{R}$. To remind ourselves: a function $f(x)$ is strictly convex⁴ for $x \in \mathbb{R}$ if

$$\forall x_1 \neq x_2 \in \mathbb{R}, \forall t \in (0, 1) : f(tx_1 + (1-t)x_2) < tf(x_1) + (1-t)f(x_2). \quad (13)$$

Implications of this are that that

$$\frac{d^2 f(x)}{dx^2} > 0, \forall x \in \mathbb{R} \quad (14)$$

and $f(x)$ contains at most one minimum, which in turn is then the global minimum. We can then find this global minimum by solving $\frac{df(x)}{dx} = 0$, which has now only one solution. We can check that EQ. 14 holds for $\mathcal{L}(\gamma, \mathbf{v}, \mathbf{u})$ by doing

$$\frac{\partial^2 \mathcal{L}(\gamma, \mathbf{v}, \mathbf{u})}{\partial \gamma_{ij}^2} = \frac{\lambda}{\gamma_{ij}} > 0, \forall \gamma_{ij} \in \mu(\mathbf{a}, \mathbf{b}). \quad (15)$$

Now we can obtain γ^* by doing

$$\frac{\partial \mathcal{L}(\gamma, \mathbf{v}, \mathbf{u})}{\partial \gamma_{ij}} = M_{ij,p} + \lambda \log \gamma_{ij} - v_i - u_j = 0. \quad (16)$$

The solution to this equation is

$$\gamma_{ij}^* = \exp\left(\frac{v_i}{\lambda}\right) \exp\left(\frac{u_j}{\lambda}\right) \exp\left(-\frac{M_{ij,p}}{\lambda}\right) = f_i g_j K_{ij}, \quad (17)$$

where the K_{ij} resembles the Gibbs kernel as used in statistical mechanics.⁵ λ Can be thought of as a temperature parameter similar to $k_b T$. Making λ larger increasingly desensitizes the system to unoptimal permutations present in γ . To make the notation more simple we briefly change to integral notation. We assume $\gamma^*(x, y)$ has the following form

$$\gamma^*(x, y) = f(x)g(y)K(x, y). \quad (18)$$

Now, we can rewrite the constraints as

$$\begin{aligned} \int \gamma^*(x, y) dy &= a(x), \\ \int \gamma^*(x, y) dx &= b(y). \end{aligned} \quad (19)$$

Note here that the first two constraints are now also satisfied, because of the constraints we placed on \mathbf{a} and \mathbf{b} (EQ. 5). Now we write

$$\begin{aligned} g(y) &= \frac{b(y)}{\int f(x)K(x, y)dx}, \\ f(x) &= \frac{a(x)}{\int g(y)K(x, y)dy}, \end{aligned} \quad (20)$$

which can be transformed into an iterative scheme

$$\begin{aligned} g^{n+1}(y) &= \frac{b(y)}{\int f^n(x)K(x, y)dx}, \\ f^{n+1}(x) &= \frac{a(x)}{\int g^{n+1}(y)K(x, y)dy}. \end{aligned} \quad (21)$$

When going back to the matrix/vector notation we can now write the complete algorithm, which is presented in the next section.

A. The algorithm

- Initialization:

$$\begin{aligned} \mathbf{g}^0 &= (1, 1, \dots, 1)^T, \text{ where } \dim(\mathbf{g}) = \dim(\mathbf{b}) \\ \mathbf{f}^0 &= \mathbf{a}^0 \\ \gamma_{ij}^{*0} &= f_i^0 K_{ij} g_j^0 \leftrightarrow \boldsymbol{\gamma}^{*0} = \text{diag}(\mathbf{f}^0) \mathbf{K} \text{diag}(\mathbf{g}^0) \\ \text{iterations} &= 0 \end{aligned}$$

- While $\text{iterations} < \text{max-iterations}$ and $\|\mathbf{g} \odot (\mathbf{K}^T \mathbf{f}) - \mathbf{b}\|_2 > \text{threshold}$:

$$\begin{aligned} g_i^{n+1} &= \frac{b_i}{(\mathbf{K}^T \mathbf{f}^n)_i} \\ f_i^{n+1} &= \frac{a_i}{(\mathbf{K} \mathbf{g}^{n+1})_i} \\ \gamma_{ij}^{*n+1} &= f_i^{n+1} K_{ij} g_j^{n+1} \leftrightarrow \boldsymbol{\gamma}^{*n+1} = \text{diag}(\mathbf{f}^{n+1}) \mathbf{K} \text{diag}(\mathbf{g}^{n+1}) \\ \text{iterations} &\rightarrow \text{iterations} + 1, \end{aligned}$$

where \odot represents the Hadamard product. One should be careful when changing the order of updating \mathbf{g} and \mathbf{f} , due to how the convergence condition is defined.

¹ Rémi Flamary and Nicolas Courty. Python optimal transport. <https://pot.readthedocs.io/en/stable/>, 2019. Accessed: 28-01-2020.

² Gabriel Peyré and Marco Cuturi. *Foundations and Trends in Machine Learning*, 11:355–607, 2019.

³ Paul Dawkins. Paul’s online notes, section 3-5: Lagrange multi-

pliers, 2020. Accessed: 29-01-2020.

⁴ Markus Grasmair. Basic properties of convex functions. Accessed: 29-01-2020.

⁵ A. Bovier, F. Dunlop, A. van Enter, F. den Hollander, and J. Dalibard. *Mathematical statistical physics*. Elsevier, 2006.