

# Lecture 6: Divergences between Probability Measures

Lénaïc Chizat

March 10, 2021

The material of today's lecture is partly adapted from Gabriel Peyré's lecture notes.

## 1 Motivating problem: density fitting

In statistics, imaging, or machine learning, one of the most fundamental problems is to compare a probability distribution  $\nu \in \mathcal{P}(\mathbb{R}^d)$  arising from measurements to a model, namely a parameterized family of distributions  $\{\mu_\theta, \theta \in \Theta\}$  where typically  $\Theta \subset \mathbb{R}^p$ . A suitable parameter can be obtained by minimizing

$$\min_{\theta \in \Theta} F(\theta) := D(\mu_\theta, \nu)$$

where  $D : \mathcal{P}(\mathbb{R}^d) \times \mathcal{P}(\mathbb{R}^d) \rightarrow [0, +\infty]$  is a *divergence*<sup>1</sup>, i.e. a quantity that quantifies the discrepancy between  $\mu$  and  $\nu$ .

**Example 1.1.** One can choose  $D(\mu, \nu) = W_p^p(\mu, \nu)$ . When  $\nu$  is an empirical measure and with  $p = 2$ , this is called the *Minimum Kantorovich estimator*. A drawback of this discrepancy is that it is computationally expensive, compared to other discrepancies that we will see in this lecture.

**Example 1.2** (Maximum likelihood). Let  $x_1, \dots, x_n \in \mathbb{R}^d$  be independent samples from  $\nu$ . When  $\mu_\theta$  has a density  $\rho_\theta$  with respect to a reference measure  $\sigma$  (e.g. the Lebesgue measure), the maximum likelihood estimator (MLE) is obtained by solving

$$\min_{\theta \in \Theta} -\frac{1}{n} \sum_{i=1}^n \log(\rho_\theta(x_i)).$$

This corresponds to using an empirical counterpart of the Kullback-Leibler loss since this converges towards  $-\int_X \log(\rho_\theta(x)) d\nu(x) = \text{KL}(\mu_\theta, \nu) - \int \log(d\nu/d\sigma) d\nu$  (one can show this equality provided  $\mu_\theta \ll \nu$  and  $\nu \ll \sigma$  and all the terms are finite).

The MLE is a statistically optimal estimation procedure in certain cases, but fails:

- when there is no natural reference measure  $\sigma$ ;
- when the density  $\rho_\theta$  is difficult to compute;
- the resulting objective  $F$  is too complicated to minimize.

---

<sup>1</sup>In (applied) mathematics, *divergence* generally refers to a nonnegative quantity  $D(a, b)$  that satisfies  $D(a, b) = 0$  if  $a = b$  (such as a distance, a squared distance, the Kullback-Leibler divergence,...)

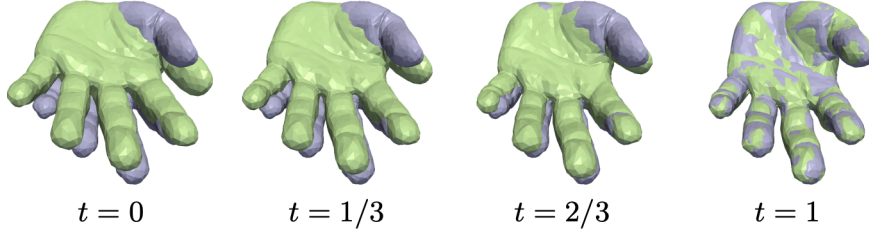


Figure 1: Gradient descent algorithm for density fitting. Here  $(h_\theta)_\#$  is a parameterized set of diffeomorphisms,  $(h_{\theta_t})_\#\zeta$  is in green and the target  $\nu$  is in purple. Image from [3].

**Generative models.** A typical set-up where all these problems appear is for so-called generative models, where the parametric measure is written as a push-forward of a fixed reference measure  $\zeta \in \mathcal{P}(Z)$

$$\mu_\theta = (h_\theta)_\#\zeta \quad \text{where} \quad h_\theta : Z \rightarrow \mathbb{R}^d.$$

This leads to the objective function  $F(\theta) = D((h_\theta)_\#\zeta, \nu)$ .

The typical approach to tackle such problems numerically, is the gradient descent algorithm: initialize  $\theta_0 \in \Theta$  and define for  $t \geq 0$ ,

$$\theta_{t+1} = \theta_t - \eta \nabla F(\theta_t)$$

where  $\eta > 0$  is a step-size (with potentially a projection on  $\Theta$  if it is not a vector space). See Figure 1 for an example. Since  $F$  is non-convex, there is no guarantee that  $F(\theta_t)$  converges to the minimum. In practice, the algorithm behaves better when the divergence is “geometrically faithful” (such as  $W_2^2$ ). Let us give a formula for the gradient under strong regularity assumptions (which could be relaxed). Here  $E : \mu \mapsto D(\mu, \nu)$ .

**Proposition 1.3** (Chain rule for generative models). *Assume that  $E : \mathcal{P}_2(\mathbb{R}^d) \rightarrow \mathbb{R}$  is such that for all  $\mu \in \mathcal{P}_2(\mathbb{R}^d)$ , there exists a function  $E'(\mu) \in \mathcal{C}^1(\mathbb{R}^d)$  with  $\nabla E'(\mu)$  Lipschitz, and such that for all  $\nu \in \mathcal{P}_2(\mathbb{R}^d)$ ,*

$$E(\nu) - E(\mu) = \int_{\mathbb{R}^d} E'(\mu) d(\nu - \mu) + o(W_2(\mu, \nu)).$$

*Assume moreover that  $h : \mathbb{R}^p \rightarrow L^2(\zeta; \mathbb{R}^d)$  is (Fréchet) differentiable, with partial derivatives at  $\theta$  denoted by  $\partial_i h_\theta \in L^2(\zeta; \mathbb{R}^d)$ . Then  $F : \theta \mapsto E((h_\theta)_\#\zeta)$  is (Fréchet) differentiable with gradient, for  $i = 1, \dots, p$ ,*

$$[\nabla F(\theta)]_i = \int_Z \nabla E'((h_\theta)_\#\zeta)(h_\theta(z))^\top \partial_i h_\theta(z) d\zeta(z).$$

*Proof.* We first study  $G : f \mapsto E(f_\#\zeta)$  and show that  $G$  is Fréchet differentiable with differential  $DG(f)(\delta f) = \int \nabla E'(f_\#\zeta)(f(z))^\top \delta f(z) d\zeta(z)$ . Then the conclusion follows by the usual chain rule for Fréchet differentials.

For  $f, \delta f \in L^2(\zeta, \mathbb{R}^d)$ , we have that  $W_2^2(f_\#\zeta, (f + \delta f)_\#\zeta) \leq \|\delta f\|_{L^2(\zeta)}^2$  by taking  $(f, f + \delta f)_\#\zeta$  as an admissible transport plan. Thus, by our assumption on  $E$ ,

$$\begin{aligned} E((f + \delta f)_\#\zeta) - E(f_\#\zeta) &= \int_Z [E'(f_\#\zeta)(f(z) + \delta f(z)) - E'(f_\#\zeta)(f(z))] d\zeta(z) + o(\|\delta f\|) \\ &= \int_Z \nabla E'(f_\#\zeta)(f(z))^\top \delta f(z) d\zeta(z) + O(\text{Lip}(\nabla E'(f_\#\zeta)) \|\delta f\|^2) + o(\|\delta f\|) \end{aligned}$$

by a Taylor expansion of  $E'(f_{\#}\zeta)$  in the integral. This shows the formula for  $DG(f)$  and concludes the proof.  $\square$

**Example 1.4.** Show that if  $W : \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}$  is symmetric and differentiable with a Lipschitz gradient, then  $E(\mu) := \int W(x, y) d\mu(x) d\mu(y)$  satisfies the assumptions of Prop 1.3 with  $E'(\mu)(x) = \int W(x, y) d\mu(y)$ .

**Rest of this lecture** We will now introduce various notions of divergences : Csiszár-divergences, dual norms, MMD and Sinkhorn divergences. Much can be said about each of them, but we will focus on discussing (i) the divergence property and (ii) the weak continuity. In the rest of this lecture, we assume that  $X$  is a compact metric space.

## 2 Csiszár divergences

Maybe the most classical way to compare two probability measures are the total variation norm and the Kullback-Leibler divergence. They belong to the family of *Csiszár divergences* – also known as *f-divergences* – which consist in comparing the relative densities to 1. They are simple to compute between discrete distributions ( $O(n)$  operations for distributions with  $n$  atoms) but are not weakly continuous.

**Definition 2.1** (*f-divergence*). Let  $f : \mathbb{R} \rightarrow \mathbb{R} \cup \{+\infty\}$  be a convex function. For any  $\mu, \nu \in \mathcal{P}(X)$ , let  $\mu = \frac{d\mu}{d\nu}\nu + \mu^\perp$  be the Lebesgue decomposition of  $\mu$  with respect to  $\nu$ . The divergence is defined by

$$D_f(\mu, \nu) := \int_X f\left(\frac{d\mu}{d\nu}\right) d\nu + f'_\infty(1) \cdot \mu^\perp(X).$$

where  $f'_\infty(x) = \lim_{t \rightarrow \infty} f(tx)/t \in \mathbb{R} \cup \{\infty\}$  is the asymptotic speed of growth of  $f$  in the direction  $x$ .

If  $f'_\infty(1) = \infty$  then  $f$  grows faster than any linear function ( $f$  is said *superlinear*).

**Proposition 2.2.** Let  $f$  be convex such that  $\min f = 0$  and  $\arg \min f = \{1\}$ . Then  $D_f(\mu, \nu) \geq 0$  with equality if and only if  $\mu = \nu$ .

*Proof.* If  $\mu = \nu$ , then  $d\mu/d\nu = 1 \in L^1(\nu)$  and  $\mu^\perp = 0$ , thus  $D_f(\mu, \nu) = \int f(1) d\nu = 0$ . Conversely, if  $D_f(\mu, \nu) = 0$  then  $\mu^\perp = 0$  (because  $f'_\infty(1) \geq f(2) - f(1) > 0$ ) and  $d\mu/d\nu = 1 \in L^1(\nu)$  so  $\mu = \nu$ .  $\square$

**Example 2.3** (Kullback-Leibler divergence/relative entropy). This is the Csiszár divergence associated to the function

$$f(s) = \begin{cases} s \log(s) - s + 1 & \text{if } s > 0 \\ 1 & \text{if } s = 0 \\ +\infty & \text{if } s < 0 \end{cases}$$

which is convex, lsc, with unique minimum  $f(1) = 0$ . If  $\mu \ll \nu$  then

$$D_f(\mu, \nu) = \int_X \left( \frac{d\mu}{d\nu} \log\left(\frac{d\mu}{d\nu}\right) - \frac{d\mu}{d\nu} + 1 \right) d\nu = \int_X \log\left(\frac{d\mu}{d\nu}\right) d\mu = \text{KL}(\mu, \nu)$$

and  $D_f(\mu, \nu) = +\infty$  otherwise since  $f'_\infty(1) = +\infty$ .

**Example 2.4** (Total variation). This is the Csiszár divergence associated to

$$f(s) = \begin{cases} |s - 1| & \text{if } s \geq 0 \\ +\infty & \text{otherwise} \end{cases}$$

We have  $f'_\infty(1) = 1$  thus

$$D_f(\mu, \nu) = \int_X \left( \left| \frac{d\mu}{d\nu} - 1 \right| d\nu + d\mu^\perp \right) \stackrel{(*)}{=} \int_X d|\mu - \nu| = |\mu - \nu|(X)$$

where  $(*)$  comes from the fact that  $(\mu - \nu)_+ = \max\{0, d\mu/d\nu - 1\}\nu + \mu^\perp$  and  $(\mu - \nu)_- = \max\{0, 1 - d\mu/d\nu\}\nu$ . Beware that in the probability literature there is sometimes a factor  $1/2$  in front of the definition (so that it takes values in  $[0, 1]$  for probability distributions).

In the context of generative models, a drawback of  $f$ -divergences is that they are not weakly continuous : for instance  $D_f(\delta_x, \delta_y) = f'_\infty(1) \cdot 1_{x=y}$  is not continuous at  $x = y$ . We have however weak lower-semicontinuity.

**Proposition 2.5.** *If  $f : \mathbb{R} \rightarrow \mathbb{R} \cup \{+\infty\}$  is convex, lsc, and not identically  $+\infty$ , then  $D_f(\mu, \nu)$  is (jointly) convex and weakly lower-semicontinuous and one has*

$$D_f(\mu, \nu) = \sup_{\varphi, \psi \in \mathcal{C}(X)} \int \varphi d\mu + \int \psi d\nu \quad \text{s.t.} \quad \varphi(x) + f^*(\psi(x)) \leq 0, \forall x \in X$$

where  $f^* : s \mapsto \sup_{u \in \mathbb{R}} us - \cdot f(u)$  is the convex conjugate of  $f$ .

*Proof idea.* This is a special case of a more general property concerning *integral functionals of perspective functions*. Let  $\psi_f : \mathbb{R}_+ \times \mathbb{R} \rightarrow \mathbb{R} \cup \{+\infty\}$  be the perspective of  $f$ , defined as

$$\psi_f(t, x) = \begin{cases} t \cdot f(x/t) & \text{if } t > 0 \\ f'_\infty(x) & \text{if } t = 0 \\ +\infty & \text{if } t < 0 \end{cases}$$

By direct computation, it can be seen  $\psi_f$  is the convex conjugate of the (convex and lsc) function of  $(s, y) \in \mathbb{R}^2$  that is worth 0 if  $s + f^*(y) \leq 0$  and  $+\infty$  otherwise (exercise). Thus  $\psi_f$  is convex and lsc. For any  $\sigma \in \mathcal{P}(X)$  such that  $\mu, \nu \ll \sigma$ , it holds

$$D_f(\mu, \nu) = \int \psi_f\left(\frac{d\nu}{d\sigma}, \frac{d\mu}{d\sigma}\right) d\sigma = \int \left( \sup_{\psi + f^*(\varphi) \leq 0} \psi \frac{d\nu}{d\sigma} + \varphi \frac{d\mu}{d\sigma} \right) d\sigma \stackrel{(*)}{=} \sup_{\substack{\varphi, \psi \in \mathcal{C}(X) \\ \psi + f^* \circ \varphi \leq 0}} \int \varphi d\mu + \int \psi d\nu$$

To exchange  $\int$  and  $\sup$ , we have used a (non-trivial) interversion theorem [5, Thm. 6]. The convexity and weak lsc follow directly by this dual representation. See [1, Sec. 2.6] for a direct proof of lower-semicontinuity.  $\square$

**Example 2.6** (Total variation). For the total variation, we have by direct computation  $f^*(s) = \max\{-1, s\}$  for  $s \leq 1$  and  $f^*(s) = +\infty$  for  $s > 1$ , so we recover the usual dual characterization

$$|\mu - \nu|(X) = \sup_{\substack{\varphi \in \mathcal{C}(X) \\ \varphi \leq 1}} \int \varphi d\mu - \int \max\{-1, \varphi\} d\nu = \sup_{\substack{\varphi \in \mathcal{C}(X) \\ \|\varphi\|_\infty \leq 1}} \int \varphi d(\mu - \nu).$$

### 3 Integral probability metrics (dual norms)

#### 3.1 General case

For a symmetric set  $B$  of measurable functions from  $X$  to  $\mathbb{R}$  and  $\alpha \in \mathcal{M}(X)$  a signed measure, let

$$\|\alpha\|_B := \sup_{f \in B} \int_X f(x) d\alpha(x) \quad (3.1)$$

The divergences associated to such dual norms, obtained with  $\alpha = \mu - \nu$ , for  $\mu, \nu \in \mathcal{P}(X)$

$$D_B(\mu, \nu) := \|\mu - \nu\|_B = \sup_{f \in B} \int f(x) d(\mu(x) - \nu(x))$$

are often called “integral probability metrics”, see [6] (or also “maximum mean discrepancy” but the latter is sometimes reserved to the special case discussed later).

**Proposition 3.1.** *If  $B$  is symmetric, bounded in sup-norm and contains 0, then  $\|\cdot\|_B$  is a seminorm on  $\mathcal{M}(X)$  (it is nonnegative, positively homogeneous and subadditive).*

**Example 3.2** (Total variation). It is recovered with  $B = \{f \in \mathcal{C}(X) ; \|f\|_\infty \leq 1\}$ .

**Example 3.3** (Wasserstein-1). It is the integral probability metric induced by the set of 1-Lipschitz functions  $B = \{f \in \mathcal{C}(X) ; \text{Lip}(f) \leq 1\}$ .

**Example 3.4** (Flat norm and the Dudley metric). If the set  $B$  is bounded in  $\|\cdot\|_\infty$ , then  $\|\cdot\|_B$  is a norm on the whole space  $\mathcal{M}(X)$  of signed measures. This is not the case for  $\|\cdot\|_{W_1}$ , which is only finite for  $\alpha$  such that  $\int_X d\alpha = 0$ . This can be alleviated by imposing a bound on the value of the potential  $f$ , in order to define for instance the *flat norm*

$$B = \{f ; \text{Lip}(f) \leq 1 \quad \text{and} \quad \|f\|_\infty \leq 1\}$$

It is similar to the *Dudley* metric, which uses

$$B = \{f ; \text{Lip}(f) + \|f\|_\infty \leq 1\}.$$

The following proposition shows that to metrize the weak convergence, the set  $B$  should not be too large nor too small.

**Proposition 3.5.** *Let  $(\alpha_k)_k$  be a bounded (for total variation  $\|\cdot\|_{TV}$ ) sequence in  $\mathcal{M}(X)$ .*

- (i) *If  $\mathcal{C}(X) \subset \overline{\text{span}(B)}^{\|\cdot\|_\infty}$ , i.e. if the span of  $B$  is dense in the set of continuous functions endowed with the sup-norm), then  $\|\alpha_k - \alpha\|_B \rightarrow 0$  implies  $\alpha_k \rightarrow \alpha$*
- (ii) *If  $B \subset \mathcal{C}(X)$  is compact (i.e. if it is closed, uniformly continuous and bounded) then  $\alpha_k \rightarrow \alpha$  implies  $\|\alpha_k - \alpha\|_B \rightarrow 0$ .*

*Proof.* (i) If  $\|\alpha_k - \alpha\|_B \rightarrow 0$  then by duality, for any  $f \in B$ , since  $|\langle f, \alpha_k - \alpha \rangle| \leq \|\alpha_k - \alpha\|_B$  then  $\langle f, \alpha_k \rangle \rightarrow \langle f, \alpha \rangle$ . By linearity, this extends to  $\text{span}(B)$  and then to  $\overline{\text{span}(B)}^{\|\cdot\|_\infty}$  since  $|\langle f, \alpha_k \rangle - \langle f', \alpha_k \rangle| \leq \|f - f'\|_\infty \sup_k \|\alpha_k\|_{TV}$ .

(ii) We assume that  $\alpha_k \rightarrow \alpha$  and we consider a subsequence  $\alpha_{n_k}$  such that

$$\|\alpha_{n_k} - \alpha\|_B \rightarrow \limsup_k \|\alpha_k - \alpha\|_B.$$

Since  $B$  is compact, the maximum appearing in the definition of  $\|\alpha_{n_k} - \alpha\|_B$  is reached, so there exists some  $f \in B$  such that  $\langle \alpha_{n_k} - \alpha, f_{n_k} \rangle = \|\alpha_{n_k} - \alpha\|_B$ . By compactness, we can again extract a subsequence  $f_{n_k}$  (not relabelled for simplicity) that converges to some  $f \in B \subset \mathcal{C}(X)$ . One has

$$\|\alpha_{n_k} - \alpha\|_B = \langle \alpha_{n_k} - \alpha, f \rangle + \langle \alpha_{n_k}, f_{n_k} - f \rangle - \langle \alpha, f_{n_k} - f \rangle \rightarrow 0$$

because  $\alpha_{n_k} - \alpha \rightarrow 0$  and  $\|f_{n_k} - f\|_\infty \rightarrow 0$ .  $\square$

Observe that this proof is a direct generalization of our proof that  $W_1$  metrizes the weak topology on a compact  $X$  in Lecture 4.

### 3.2 Kernel Maximum Mean Discrepancies

We now describe an important class of integral probability metrics.

**Definition 3.6** (Positive definite kernel). A symmetric function  $k : X \times X \rightarrow \mathbb{R}$  is said to be positive definite (p.d.) if for any  $n \geq 1$ , for any family  $x_1, \dots, x_n \in X$  the matrix  $(k(x_i, x_j))_{i,j}$  is positive semi-definite<sup>2</sup>, i.e. for all  $r \in \mathbb{R}^n$ ,

$$\sum_{i,j=1}^n r_i r_j k(x_i, x_j) \geq 0. \quad (3.2)$$

The kernel is said to be conditionally positive definite if Eq. (3.2) holds for all zero mean vectors  $r$ , i.e. such that  $\sum_i r_i = 0$ .

**Definition 3.7** (MMD). Given a continuous and positive semi-definite kernel  $k : X \times X \rightarrow \mathbb{R}$ , we define for  $\alpha \in \mathcal{M}(X)$  (finite signed Borel measure)

$$\|\alpha\|_k^2 = \iint_X k(x, y) d\alpha(x) d\alpha(y).$$

The squared MMD between  $\mu, \nu \in \mathcal{P}(X)$  is then

$$\|\mu - \nu\|_k^2 = \iint k d\mu \otimes \mu + \iint k d\nu \otimes \nu - 2 \iint k d\mu \otimes \nu.$$

This definition as a squared quantity makes sense thanks to the following result.

**Proposition 3.8.** *If  $k \in \mathcal{C}(X^2)$  is conditionally p.d.,  $\iint k d\alpha \otimes \alpha \geq 0$  if  $\int d\alpha = 0$ .*

*Proof.* Let  $\alpha_n \in \mathcal{M}(X)$  be measures with finite support and zero mass such that  $\alpha_n \rightarrow \alpha$ . Since  $\alpha_n \otimes \alpha_n \rightarrow \alpha \otimes \alpha$  (Lem.2.2 from Lect. 5), we have  $0 \leq \iint k d\alpha_n \otimes \alpha_n \rightarrow \iint k d\alpha \otimes \alpha$ .  $\square$

**MMD as a dual norm.** To develop a finer understanding of MMD, one needs to develop the theory of *Reproducible Kernel Hilbert Spaces* (RKHS). Since this is beyond the scope of this course, we limit ourselves to showing a link with dual norms via the following result.

**Theorem 3.9** (Aronzsjajn).  *$k$  is a p.d. kernel on the set  $X$  if and only if there exists a Hilbert space  $\mathcal{H}$  and a mapping  $\Phi : X \rightarrow \mathcal{H}$  such that  $k(x, x') = \langle \Phi(x), \Phi(x') \rangle_{\mathcal{H}}$ .*

<sup>2</sup>It would be more consistent to call such a kernel “positive semi-definite” but we are using the convention from the literature [7].

As a consequence the MMD consists in embedding  $\mathcal{M}(X)$  into a Hilbert space  $\mathcal{H}$  via the *kernel mean embedding*  $\mu \mapsto \int \Phi d\mu$ , since  $\|\mu - \nu\|_k = \|\int \Phi d\mu - \int \Phi d\nu\|_{\mathcal{H}}$ . We also have

$$\|\alpha\|_k = \sup_{\|h\|_{\mathcal{H}} \leq 1} \left\langle h, \int \Phi d\alpha \right\rangle = \sup_{f \in B} \int f d\alpha$$

where  $B = \{x \mapsto \langle h, \Phi(x) \rangle ; \|h\|_{\mathcal{H}} \leq 1\}$  so it is an integral probability metric.

**Metriizing weak convergence.** It can be shown that if  $k$  is *universal* (i.e. (i) of Prop. 3.5) holds), continuous and conditionally positive definite then  $\|\cdot\|_k$  metrizes weak convergence in  $\mathcal{P}(X)$ , see e.g. [?]. Examples of such kernels on  $\mathbb{R}^d$  are:

- the Gaussian kernel  $k(x, y) = e^{-\frac{\|y-x\|_2^2}{2\sigma^2}}$  with  $\sigma > 0$ ;
- the distance kernel  $k(x, y) = -\text{dist}(x, y)$  (its MMD is called the “Energy distance”);

**Discrete case.** In the special case of discrete measures  $\mu = \sum_{i=1}^m a_i \delta_{x_i}$  and  $\nu = \sum_{j=1}^n b_j \delta_{y_j}$  then we have

$$\|\mu - \nu\|_k^2 = \sum_{i,i'} a_i a_{i'} k(x_i, x_{i'}) + \sum_{j,j'} b_j b_{j'} k(y_j, y_{j'}) - 2 \sum_{i,j} a_i b_j k(x_i, y_j).$$

This requires  $O((m+n)^2)$  operations to compute.

## 4 Sinkhorn divergences

Recall that  $W_p^p$  are often good choices of divergence, but are computationally expensive (another downside, not discussed in this course, is that they are difficult to estimate from random samples, in particular in high dimension [8]). Next, we build divergences from entropy regularized optimal transport.

### 4.1 Reminders on entropy regularized OT

With  $c \in \mathcal{C}(X^2)$ , the definition of entropy regularized optimal transport is

$$T_{c,\lambda}(\mu, \nu) := \min_{\gamma \in \Pi(\mu, \nu)} \int c(x, y) d\gamma(x, y) + \lambda \text{KL}(\gamma; \mu \otimes \nu)$$

Note that this definition differs from the one in Lecture 3 by a constant  $\lambda$  (because the quantity  $\mathcal{H}$  from Lecture 3 differs from KL by 1). We recall the following facts from Lecture 3:

- (Duality)

$$T_{c,\lambda}(\mu, \nu) = \sup_{\varphi, \psi \in \mathcal{C}(X)} \int \varphi(x) d\mu(x) + \int \psi(y) d\nu(y) + \lambda \left( 1 - \iint e^{(\varphi(x) + \psi(y) - c(x,y))/\lambda} d\mu(x) d\nu(y) \right)$$

- (Optimality conditions) There exists maximizers  $(\varphi_\lambda, \psi_\lambda)$  and a unique minimizer  $\gamma_\lambda$  linked by the optimality condition

$$d\gamma_\lambda(x, y) = e^{(\varphi_\lambda(x) + \psi_\lambda(y) - c(x,y))/\lambda} d\mu(x) d\nu(y)$$

It follows in particular that

$$T_{c,\lambda}(\mu, \nu) = \int \varphi_\lambda d\mu + \int \psi_\lambda d\nu.$$

## 4.2 Is $T_{c,\lambda}$ a suitable divergence?

**Proposition 4.1** (Interpolation properties). *For  $\mu, \nu \in \mathcal{P}(X)$  and  $c \in \mathcal{C}(X \times X)$ , it holds*

$$T_{c,\lambda}(\mu, \nu) \rightarrow \begin{cases} T_c(\mu, \nu) := T_{c,0}(\mu, \nu) & \text{as } \lambda \rightarrow 0 \\ \int c(x, y) d\mu(x) d\nu(y) & \text{as } \lambda \rightarrow \infty \end{cases}$$

Moreover, denoting  $\gamma_\lambda$  the unique minimizer for  $T_{c,\lambda}$ , it holds  $\gamma_\lambda \rightarrow \mu \otimes \nu$  as  $\lambda \rightarrow \infty$ .

*Proof.* (i) We first study the limit  $\lambda \rightarrow 0$ . Since  $\text{KL} \geq 0$ , it holds  $T_{c,\lambda} \geq T_c$ . To prove the reverse inequality as  $\lambda \rightarrow 0$ , let  $\gamma_0 \in \Pi(\mu, \nu)$  be optimal for  $T_c(\mu, \nu)$ . For any  $\epsilon > 0$ , there exists  $\gamma_\epsilon \in \Pi(\mu, \nu)$  such that  $|\int c d\gamma - \int c d\gamma_\epsilon| \leq \epsilon$  and  $\text{KL}(\gamma_\epsilon, \mu \otimes \nu) < +\infty$ . One possible way to build  $\gamma_\epsilon$  is to take  $(Q_i)$  a partition of  $X$  into sets of diameter less than  $\text{Lip}(c)/(2\epsilon)$  and to take

$$\gamma_\epsilon = \sum_{i,j} \frac{\gamma_0(Q_i \times Q_j)}{\mu(Q_i)\nu(Q_j)} d(\mu|_{Q_i} \otimes \nu|_{Q_j})$$

(Exercise: show that  $\gamma_\epsilon$  indeed satisfies our requirements). It follows that

$$T_{c,\lambda} \leq \epsilon + \lambda \text{KL}(\gamma_\epsilon, \mu \otimes \nu) \xrightarrow{\lambda \rightarrow 0} \epsilon.$$

As  $\epsilon > 0$  was arbitrary, this shows that  $T_{c,\lambda}(\mu, \nu) \xrightarrow{\lambda \rightarrow 0} T_c(\mu, \nu)$ .

(ii) Now we study the limit  $\lambda \rightarrow \infty$ . First, it is clear that  $T_{c,\lambda}(\mu, \nu) \leq \int c d\mu \otimes \nu$  since  $\mu \otimes \nu \in \Pi(\mu, \nu)$ . Let  $(\lambda_k)_k$  be a positive sequence that diverges to  $+\infty$ , and let  $\gamma_k$  be the corresponding sequence of (unique) minimizers for  $T_{c,\lambda_k}$ . By optimality, we have  $\int c d\gamma_k + \lambda_k \text{KL}(\gamma_k, \mu \otimes \nu) \leq \int c d\mu \otimes \nu$  and thus

$$\text{KL}(\gamma_k, \mu \otimes \nu) \leq \frac{1}{\lambda_k} \left( \int c d\mu \otimes \nu - \int c d\gamma_k \right) \rightarrow 0.$$

Moreover, by compactness of  $\Pi(\mu, \nu)$  we can extract a converging subsequence  $\gamma_{n_k} \rightarrow \gamma_\infty$ . Since  $\text{KL}$  is weakly lower-semicontinuous (Prop. 2.5), it holds

$$\text{KL}(\gamma_\infty, \mu \otimes \nu) \leq \liminf_{k \rightarrow \infty} \text{KL}(\gamma_{n_k}, \mu \otimes \nu) = 0$$

Hence  $\gamma_\infty = \mu \otimes \nu$  hence the conclusion.  $\square$

One could imagine replacing Wasserstein by its entropy regularized version, but the previous result shows that when  $\lambda$  is large,  $T_{c,\lambda}$  behaves like an inner product rather than like a divergence. In particular, even for standard costs  $c = \text{dist}(x, y)^p$ ,  $\mu \mapsto T_c^\lambda(\mu, \nu)$  is in general not minimized at  $\mu = \nu$ .

**Corollary 4.2.** *Let  $\nu \in \mathcal{P}(X)$  be such that  $\arg \min_{y \in X} \int c(x, y) d\nu(y)$  is a singleton, denoted  $x^*$  and let  $\mu_\lambda \in \arg \min_{\mu \in \mathcal{P}(X)} T_{c,\lambda}(\mu, \nu)$ . Then as  $\lambda \rightarrow \infty$ , one has  $\mu_\lambda \rightarrow \delta_{x^*}$ .*

*Proof.* By assumption, the continuous function  $f$  defined by  $f(x) = \int c(x, y) d\nu(y)$  admits a unique minimizer  $x^* \in X$ . It follows that for any  $r > 0$ , there exists  $\epsilon > 0$  such that  $|f(x) - f(x^*)| \leq \epsilon$  implies  $\text{dist}(x, x^*) \leq r$  (to see this, observe that  $g(r) := \min_{x \in X, \text{dist}(x, x^*) \geq r} f(x)$  satisfies  $g(r) > f(x^*)$  and  $\lim_{r \rightarrow 0^+} g(r) = f(x^*)$ ).

Now consider an increasing unbounded sequence  $(\lambda_k)_k$  and let  $(\mu_k)_k$  be the corresponding set of minimizers. Taking  $\delta_{x^*}$  as a competitor, the optimality of  $\mu_k$  together with Prop. 4.1 implies that for any  $\epsilon > 0$ , there exists  $k_0$  such that for  $k \geq k_0$ ,

$$\int f d\mu_k = \int c d\mu_k \otimes \nu \leq \int f(x^*) + \epsilon$$



Thus for any  $r > 0$ , there exists  $k'_0$  such that for  $k \geq k'_0$ ,  $\mu_k(B_r(x^*)) \geq 1 - r$  from which we deduce that  $\mu_k$  converges weakly to  $\delta_{x^*}$ .  $\square$

For instance, when  $c(x, y) = \frac{1}{2}\|y - x\|_2^2$ , then  $\mu_\lambda$  converges to a Dirac mass located at the mean  $\int x d\nu(x)$  of  $\nu$ .

### 4.3 Debiased quantity: the Sinkhorn divergence

Thinking of  $-T_{c,\lambda}$  as an inner product suggests to define

$$S_{c,\lambda}(\mu, \nu) := T_{c,\lambda}(\mu, \nu) - \frac{1}{2}T_{c,\lambda}(\mu, \mu) - \frac{1}{2}T_{c,\lambda}(\nu, \nu)$$

From a computational aspect, the debiasing terms add an essentially negligible cost because the Sinkhorn iterations for those problems are well-conditioned. We can already see that these correction terms allow to correct the asymptotic behavior when  $\lambda$  is large.

**Proposition 4.3** (Interpolation properties). *For  $\mu, \nu \in \mathcal{P}(X)$  and  $c \in \mathcal{C}(X \times X)$  it holds*

$$S_{c,\lambda}(\mu, \nu) \longrightarrow \begin{cases} T_c(\mu, \nu) & \text{as } \lambda \rightarrow 0 \\ \frac{1}{2}\|\mu - \nu\|_{-c}^2 & \text{as } \lambda \rightarrow \infty \end{cases}$$

where  $\|\cdot\|_{-c}$  is the MMD associated to the kernel  $-c$ .

*Proof.* This is immediate from Proposition 4.1 which deals with  $T_{c,\lambda}$ . In particular, when  $\lambda \rightarrow \infty$ , it shows that

$$S_{c,\lambda}(\mu, \nu) \rightarrow \int c d\mu \otimes \nu - \frac{1}{2} \int c d\mu \otimes \mu - \frac{1}{2} \int c d\nu \otimes \nu$$

which is precisely the definition of  $\frac{1}{2}\|\mu - \nu\|_{-c}^2$ .  $\square$

One can show that under regularity assumptions over  $\mu$  and  $\nu$  and for the cost  $c(x, y) = \|y - x\|_2^2$  on  $\mathbb{R}^d$  that  $|S_{c,\lambda} - T_c| = O(\lambda^2)$ , see [2]. Finally, under assumptions on the cost, we can show that  $\mu \mapsto S_{c,\lambda}(\mu, \nu)$  is minimized at  $\nu$ .

**Proposition 4.4** (Positive semi-definiteness). *If  $k(x, y) = e^{-c(x,y)/\lambda}$  is a positive semi-definite kernel, then  $S_{c,\lambda}(\mu, \nu) \geq 0$  with equality if  $\mu = \nu$ .*

*Proof.* In the following, we fix  $\lambda > 0$  and denote  $(\varphi, \varphi)$  a solution for  $T_{c,\lambda}(\mu, \mu)$  (there exists a solution of this form by symmetry and concavity of the dual problem) and  $(\psi, \psi)$  a solution for  $T_{c,\lambda}(\nu, \nu)$ . Using the suboptimal function  $(\varphi, \psi)$  in the dual maximization problem, one obtains

$$T_{c,\lambda}(\mu, \nu) \geq \int \varphi d\mu + \int \psi d\nu + \lambda \left( 1 - \int \int e^{(\varphi \oplus \psi - c)/\lambda} d\mu \otimes \nu \right)$$

Using the fact that  $T_{c,\lambda}(\mu, \mu) = 2 \int \varphi d\mu$  and  $T_{c,\lambda}(\nu, \nu) = 2 \int \psi d\nu$ , this implies

$$\frac{1}{\lambda} S_{c,\lambda}(\mu, \nu) \geq 1 - \int \int e^{(\varphi \oplus \psi - c)/\lambda} d\mu \otimes \nu = 1 - \langle \tilde{\mu}, \tilde{\nu} \rangle_k$$

with  $\tilde{\mu} = e^\varphi / \lambda \mu$ ,  $\tilde{\nu} = e^\psi / \lambda \nu$  and  $\langle \cdot, \cdot \rangle_k$  is the positive semi-definite inner product associated with the kernel  $k := e^{-c/\lambda}$ . By the optimality condition satisfied by  $\varphi$ , it holds

$$\|\tilde{\mu}\|_k^2 = \int \int e^{(\varphi(x) + \varphi(y) - c(x,y))/\lambda} d\nu(x) d\nu(y) = 1$$

and similarly  $\|\tilde{\nu}\|_k^2 = 1$ . So by Cauchy-Schwartz inequality, one has  $\langle \tilde{\mu}, \tilde{\nu} \rangle_k \leq 1$  and finally  $S_{c,\lambda}(\mu, \nu) \geq 0$ .  $\square$

In case  $e^{-c/\lambda}$  is a conditionally positive definite and universal kernel, we can further show that  $S_{c,\lambda}(\mu_n, \mu) \rightarrow 0$  if and only if  $\mu_n \rightarrow \mu$ , see [4].

## 5 Practical session

The last practical session (taken from Gabriel Peyré’s numerical tours website) is at the following link: [https://nbviewer.jupyter.org/github/gpeyre/numerical-tours/blob/master/python/optimaltransp\\_6\\_entropic\\_adv.ipynb](https://nbviewer.jupyter.org/github/gpeyre/numerical-tours/blob/master/python/optimaltransp_6_entropic_adv.ipynb). Only exercises 4 and 5 are not already solved, but take time to understand each block of code. For those who do not wish to use Jupyter, you may of course paste the code in a different terminal. Note that there is a sign mistake before cell 11: one should read +KL instead of −KL and  $P$  is the solution to entropy regularized optimal transport. For those who wish to go further on this topic, you can check out Jean Feydy’s great tutorial [http://www.math.ens.fr/~feidy/Teaching/DataScience/gradient\\_flows.html](http://www.math.ens.fr/~feidy/Teaching/DataScience/gradient_flows.html).

## References

- [1] Luigi Ambrosio, Nicola Fusco, and Diego Pallara, *Functions of bounded variation and free discontinuity problems*, Courier Corporation, 2000.
- [2] Lénaïc Chizat, Pierre Roussillon, Flavien Léger, François-Xavier Vialard, and Gabriel Peyré, *Faster wasserstein distance estimation with the sinkhorn divergence*, Advances in Neural Information Processing Systems (H. Larochelle, M. Ranzato, R. Hadsell, M. F. Balcan, and H. Lin, eds.), vol. 33, Curran Associates, Inc., 2020, pp. 2257–2269.
- [3] Jean Feydy, Benjamin Charlier, François-Xavier Vialard, and Gabriel Peyré, *Optimal transport for diffeomorphic registration*, International Conference on Medical Image Computing and Computer-Assisted Intervention, Springer, 2017, pp. 291–299.
- [4] Jean Feydy, Thibault Séjourné, François-Xavier Vialard, Shun-ichi Amari, Alain Trounev, and Gabriel Peyré, *Interpolating between optimal transport and MMD using sinkhorn divergences*, The 22nd International Conference on Artificial Intelligence and Statistics, PMLR, 2019, pp. 2681–2690.
- [5] Ralph Rockafellar, *Integrals which are convex functionals. ii*, Pacific Journal of Mathematics **39** (1971), no. 2, 439–469.
- [6] Bharath K. Sriperumbudur, Kenji Fukumizu, Arthur Gretton, Bernhard Schölkopf, and Gert R.G. Lanckriet, *On integral probability metrics,  $\varphi$ -divergences and binary classification*, arXiv preprint arXiv:0901.2698 (2009).
- [7] Bharath K Sriperumbudur, Kenji Fukumizu, and Gert RG Lanckriet, *Universality, characteristic kernels and rkhs embedding of measures.*, Journal of Machine Learning Research **12** (2011), no. 7.
- [8] Jonathan Weed, Francis Bach, et al., *Sharp asymptotic and finite-sample rates of convergence of empirical measures in wasserstein distance*, Bernoulli **25** (2019), no. 4A, 2620–2648.