

## A Standard results

Before detailing our proofs, we first recall some well-known results regarding the Kullback-Leibler divergence and the SoftMin operator defined in (12).

### A.1 The Kullback-Leibler divergence

**First properties.** For any pair of Radon measures  $\alpha, \beta \in \mathcal{M}^+(\mathcal{X})$  on the compact metric set  $(\mathcal{X}, d)$ , the Kullback-Leibler divergence is defined through

$$\text{KL}(\alpha, \beta) \stackrel{\text{def.}}{=} \begin{cases} \langle \alpha, \log \frac{d\alpha}{d\beta} - 1 \rangle + \langle \beta, 1 \rangle & \text{if } \alpha \ll \beta \\ +\infty & \text{otherwise.} \end{cases}$$

It can be rewritten as an  $f$ -divergence associated to

$$\psi : x \in \mathbb{R}_{\geq 0} \mapsto x \log(x) - x + 1 \in \mathbb{R}_{\geq 0},$$

with  $0 \cdot \log(0) = 0$ , as

$$\text{KL}(\alpha, \beta) = \begin{cases} \langle \beta, \psi(\frac{d\alpha}{d\beta}) \rangle & \text{if } \alpha \ll \beta \\ +\infty & \text{otherwise.} \end{cases} \quad (28)$$

Since  $\psi$  is a strictly convex function with a unique global minimum at  $\psi(1) = 0$ , we thus get that  $\text{KL}(\alpha, \beta) \geq 0$  with equality iff.  $\alpha = \beta$ .

**Dual formulation.** The convex conjugate of  $\psi$  is defined for  $u \in \mathbb{R}$  by

$$\begin{aligned} \psi^*(u) &\stackrel{\text{def.}}{=} \sup_{x>0} (xu - \psi(x)) \\ &= e^u - 1, \end{aligned}$$

$$\text{and we have } \psi(x) + \psi^*(u) \geq xu \quad (29)$$

for all  $(x, u) \in \mathbb{R}_{\geq 0} \times \mathbb{R}$ , with equality if  $x > 0$  and  $u = \log(x)$ . This allows us to rewrite the Kullback-Leibler divergence as the solution of a dual concave problem:

**Proposition 5** (Dual formulation of KL). *Under the assumptions above,*

$$\text{KL}(\alpha, \beta) = \sup_{h \in \mathcal{F}_b(\mathcal{X}, \mathbb{R})} \langle \alpha, h \rangle - \langle \beta, e^h - 1 \rangle \quad (30)$$

where  $\mathcal{F}_b(\mathcal{X}, \mathbb{R})$  is the space of bounded measurable functions from  $\mathcal{X}$  to  $\mathbb{R}$ .

*Proof. Lower bound on the sup.* If  $\alpha$  is not absolutely continuous with respect to  $\beta$ , there exists a Borel set  $A$  such that  $\alpha(A) > 0$  and  $\beta(A) = 0$ . Consequently, for  $h = \lambda \mathbf{1}_A$ ,

$$\langle \alpha, h \rangle - \langle \beta, e^h - 1 \rangle = \lambda \alpha(A) \xrightarrow{\lambda \rightarrow +\infty} +\infty.$$

Otherwise, if  $\alpha \ll \beta$ , we define  $h_* = \log \frac{d\alpha}{d\beta}$  and see that

$$\langle \alpha, h_* \rangle - \langle \beta, e^{h_*} - 1 \rangle = \text{KL}(\alpha, \beta).$$

If  $h_n = \log(\frac{d\alpha}{d\beta}) \mathbf{1}_{1/n \leq d\alpha/d\beta \leq n} \in \mathcal{F}_b(\mathcal{X}, \mathbb{R})$ , the monotone and dominated convergence theorems then allow us to show that

$$\langle \alpha, h_n \rangle - \langle \beta, e^{h_n} - 1 \rangle \xrightarrow{n \rightarrow +\infty} \text{KL}(\alpha, \beta).$$

*Upper bound on the sup.* If  $h \in \mathcal{F}_b(\mathcal{X}, \mathbb{R})$  and  $\alpha \ll \beta$ , combining (28) and (29) allow us to show that

$$\begin{aligned} \text{KL}(\alpha, \beta) - \langle \alpha, h \rangle + \langle \beta, e^h - 1 \rangle \\ = \langle \beta, \psi(\frac{d\alpha}{d\beta}) + \psi^*(h) - h \frac{d\alpha}{d\beta} \rangle \geq 0. \end{aligned}$$

The optimal value of  $\langle \alpha, h \rangle - \langle \beta, e^h - 1 \rangle$  is bounded above and below by  $\text{KL}(\alpha, \beta)$ : we get (30).  $\square$

Since  $\langle \alpha, h \rangle - \langle \beta, e^h - 1 \rangle$  is a convex function of  $(\alpha, \beta)$ , taking the supremum over test functions  $h \in \mathcal{F}_b(\mathcal{X}, \mathbb{R})$  defines a *convex* divergence:

**Proposition 6.** *The KL divergence is a (jointly) convex function on  $\mathcal{M}^+(\mathcal{X}) \times \mathcal{M}^+(\mathcal{X})$ .*

Going further, the density of continuous functions in the space of bounded measurable functions allows us to restrict the optimization domain:

**Proposition 7.** *Under the same assumptions,*

$$\text{KL}(\alpha, \beta) = \sup_{h \in \mathcal{C}(\mathcal{X}, \mathbb{R})} \langle \alpha, h \rangle - \langle \beta, e^h - 1 \rangle \quad (31)$$

where  $\mathcal{C}(\mathcal{X}, \mathbb{R})$  is the space of (bounded) continuous functions on the compact set  $\mathcal{X}$ .

*Proof.* Let  $h = \sum_{i \in I} h_i \mathbf{1}_{A_i}$  be a simple Borel function on  $\mathcal{X}$ , and let us choose some error margin  $\delta > 0$ . Since  $\alpha$  and  $\beta$  are Radon measures, for any  $i$  in the finite set of indices  $I$ , there exists a compact set  $K_i$  and an open set  $V_i$  such that  $K_i \subset A_i \subset V_i$  and

$$\sum_{i \in I} \max[\alpha(V_i \setminus K_i), \beta(V_i \setminus K_i)] \leq \delta.$$

Moreover, for any  $i \in I$ , there exists a continuous function  $\varphi_i$  such that  $\mathbf{1}_{K_i} \leq \varphi_i \leq \mathbf{1}_{V_i}$ . The continuous function  $g = \sum_{i \in I} h_i \varphi_i$  is then such that

$$|\langle \alpha, g - h \rangle| \leq \|h\|_\infty \delta \quad \text{and} \quad |\langle \beta, e^g - e^h \rangle| \leq \|e^h\|_\infty \delta$$

so that

$$\begin{aligned} |(\langle \alpha, h \rangle - \langle \beta, e^h - 1 \rangle) - (\langle \alpha, g \rangle - \langle \beta, e^g - 1 \rangle)| \\ \leq (\|h\|_\infty + \|e^h\|_\infty) \delta. \end{aligned}$$

As we let our simple function approach any measurable function in  $\mathcal{F}_b(\mathcal{X}, \mathbb{R})$ , choosing  $\delta$  arbitrarily small, we then get (31) through (30).  $\square$

We can then show that the Kullback-Leibler divergence is weakly lower semi-continuous:

**Proposition 8.** *If  $\alpha_n \rightharpoonup \alpha$  and  $\beta_n \rightharpoonup \beta$  are weakly converging sequences in  $\mathcal{M}^+(\mathcal{X})$ , we get*

$$\liminf_{n \rightarrow +\infty} \text{KL}(\alpha_n, \beta_n) \geq \text{KL}(\alpha, \beta).$$

*Proof.* According to (31), the KL divergence is defined as a pointwise supremum of weakly continuous applications

$$\varphi_h : (\alpha, \beta) \mapsto \langle \alpha, h \rangle - \langle \beta, e^h - 1 \rangle,$$

for  $h \in \mathcal{C}(\mathcal{X}, \mathbb{R})$ . It is thus lower semi-continuous for the convergence in law.  $\square$

## A.2 SoftMin Operator

**Proposition 9** (The SoftMin interpolates between a minimum and a sum). *Under the assumptions of the definition (12), we get that*

$$\begin{aligned} \min_{x \sim \alpha} \varphi(x) &\xrightarrow{\varepsilon \rightarrow 0} \min_{x \in \text{Supp}(\alpha)} \varphi(x) \\ &\xrightarrow{\varepsilon \rightarrow +\infty} \langle \alpha, \varphi \rangle. \end{aligned}$$

If  $\varphi$  and  $\psi$  are two continuous functions in  $\mathcal{C}(\mathcal{X})$  such that  $\varphi \leq \psi$ ,

$$\min_{x \sim \alpha} \varphi(x) \leq \min_{x \sim \alpha} \psi(x). \quad (32)$$

Finally, if  $K \in \mathbb{R}$  is constant with respect to  $x$ , we have that

$$\min_{x \sim \alpha} [K + \varphi(x)] = K + \min_{x \sim \alpha} [\varphi(x)]. \quad (33)$$

**Proposition 10** (The SoftMin operator is continuous). *Let  $(\alpha_n)$  be a sequence of probability measures converging weakly towards  $\alpha$ , and  $(\varphi_n)$  be a sequence of continuous functions that converges uniformly towards  $\varphi$ . Then, for  $\varepsilon > 0$ , the SoftMin of the values of  $\varphi_n$  on  $\alpha_n$  converges towards the SoftMin of the values of  $\varphi$  on  $\alpha$ , i.e.*

$$\begin{aligned} (\alpha_n \rightharpoonup \alpha, \varphi_n \xrightarrow{\|\cdot\|_\infty} \varphi) \\ \implies \min_{x \sim \alpha_n} \varphi_n(x) \rightarrow \min_{x \sim \alpha} \varphi(x). \end{aligned}$$

## B Proofs

### B.1 Dual Potentials

We first state some important properties of solutions  $(f, g)$  to the dual problem (8). Please note that these results hold under the assumption that  $(\mathcal{X}, d)$  is a compact metric space, endowed with a *ground cost* function  $C : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$  that is  $\kappa$ -Lipschitz with respect to both of its input variables.

The existence of an optimal pair  $(f, g)$  of potentials that reaches the maximal value of the dual objective is proved using the contractance of the Sinkhorn map  $T$ , defined in (11), for the Hilbert projective metric (Franklin and Lorenz, 1989).

While optimal potentials are only defined  $(\alpha, \beta)$ -a.e., as highlighted in Proposition 1, they are extended to the whole domain  $\mathcal{X}$  by imposing, similarly to the classical theory of OT (Santambrogio, 2015, Remark 1.13), that they satisfy

$$f = T(\beta, g) \quad \text{and} \quad g = T(\alpha, f), \quad (34)$$

with  $T$  defined in (11). We thus assume in the following that this condition holds. The following propositions studies the uniqueness and the smoothness (with respect to the spacial position and with respect to the input measures) of these functions  $(f, g)$  defined on the whole space.

**Proposition 11** (Uniqueness of the dual potentials up to an additive constant). *Let  $(f_0, g_0)$  and  $(f_1, g_1)$  be two optimal pairs of dual potentials for a problem  $\text{OT}_\varepsilon(\alpha, \beta)$  that satisfy (34). Then, there exists a constant  $K \in \mathbb{R}$  such that*

$$f_0 = f_1 + K \quad \text{and} \quad g_0 = g_1 - K. \quad (35)$$

*Proof.* For  $t \in [0, 1]$ , let us define  $f_t = f_0 + t(f_1 - f_0)$ ,  $g_t = g_0 + t(g_1 - g_0)$  and

$$\begin{aligned} \varphi(t) &= \langle \alpha, f_t \rangle + \langle \beta, g_t \rangle \\ &\quad - \varepsilon \langle \alpha \otimes \beta, \exp\left(\frac{1}{\varepsilon}(f_t \oplus g_t - C)\right) - 1 \rangle, \end{aligned}$$

the value of the dual objective between the two optimal pairs. As  $\varphi$  is a concave function bounded above by  $\varphi(0) = \varphi(1) = \text{OT}_\varepsilon(\alpha, \beta)$ , it is *constant* with respect to  $t$ . Hence, for all  $t$  in  $[0, 1]$ ,

$$\begin{aligned} 0 &= \varphi''(t) \\ &= -\frac{1}{\varepsilon} \langle \alpha \otimes \beta, e^{(f_t \oplus g_t - C)/\varepsilon} ((f_1 - f_0) \oplus (g_1 - g_0))^2 \rangle. \end{aligned}$$

This is only possible if,  $\alpha \otimes \beta$ -a.e. in  $(x, y)$ ,

$$(f_1(x) - f_0(x) + g_1(y) - g_0(y))^2 = 0,$$

i.e. there exists a constant  $K \in \mathbb{R}$  such that

$$\begin{aligned} f_1(x) - f_0(x) &= +K \quad \alpha\text{-a.e.} \\ g_1(y) - g_0(y) &= -K \quad \beta\text{-a.e.} \end{aligned}$$

As we extend the potentials through (34), the SoftMin operator commutes with the addition of  $K$  (33) and lets our result hold on the whole feature space.  $\square$

**Proposition 12** (Lipschitz property). *The optimal potentials  $(f, g)$  of the dual problem (8) are both  $\kappa$ -Lipschitz functions on the feature space  $(\mathcal{X}, d)$ .*

*Proof.* According to (34),  $f$  is a SoftMin combination of  $\kappa$ -Lipschitz functions of the variable  $x$ ; using the algebraic properties of the SoftMin operator detailed in (32-33), one can thus show that  $f$  is a  $\kappa$ -Lipschitz function on the feature space. The same argument holds for  $g$ .  $\square$

**Proposition 13** (The dual potentials vary continuously with the input measures). *Let  $\alpha_n \rightharpoonup \alpha$  and  $\beta_n \rightharpoonup \beta$  be weakly converging sequences of measures in  $\mathcal{M}_1^+(\mathcal{X})$ . Given some arbitrary anchor point  $x_o \in \mathcal{X}$ , let us denote by  $(f_n, g_n)$  the (unique) sequence of optimal potentials for  $\text{OT}_\varepsilon(\alpha_n, \beta_n)$  such that  $f_n(x_o) = 0$ .*

*Then,  $f_n$  and  $g_n$  converge uniformly towards the unique pair of optimal potentials  $(f, g)$  for  $\text{OT}_\varepsilon(\alpha, \beta)$  such that  $f(x_o) = 0$ . Up to the value at the anchor point  $x_o$ , we thus have that*

$$(\alpha_n \rightharpoonup \alpha, \beta_n \rightharpoonup \beta) \implies (f_n \xrightarrow{\|\cdot\|_\infty} f, g_n \xrightarrow{\|\cdot\|_\infty} g).$$

*Proof.* For all  $n$  in  $\mathbb{N}$ , the potentials  $f_n$  and  $g_n$  are  $\kappa$ -Lipschitz functions on the compact, bounded set  $\mathcal{X}$ . As  $f_n(x_o)$  is set to zero, we can bound  $|f_n|$  on  $\mathcal{X}$  by  $\kappa$  times the diameter of  $\mathcal{X}$ ; combining this with (34), we can then produce a uniform bound on both  $f_n$  and  $g_n$ : there exists a constant  $M \in \mathbb{R}$  such that

$$\forall n \in \mathbb{N}, \forall x \in \mathcal{X}, -M \leq f_n(x), g_n(x) \leq +M.$$

Being equicontinuous and uniformly bounded on the compact set  $\mathcal{X}$ , the sequence  $(f_n, g_n)_n$  satisfies the hypotheses of the Ascoli-Arzelà theorem: there exists a subsequence  $(f_{n_k}, g_{n_k})_k$  that converges uniformly towards a pair  $(f, g)$  of continuous functions.  $k$  tend to infinity, we see that  $f(x_o) = 0$  and, using the continuity of the SoftMin operator (Proposition 10) on the optimality equations (10), we show that  $(f, g)$  is an optimal pair for  $\text{OT}_\varepsilon(\alpha, \beta)$ .

Now, according to Proposition 11, such a limit pair of optimal potentials  $(f, g)$  is *unique*.  $(f_n, g_n)_n$  is thus a *compact* sequence with a *single* possible adherence value: it has to converge, uniformly, towards  $(f, g)$ .  $\square$

## B.2 Proof of Proposition 2

The proof is mainly inspired from (Santambrogio, 2015, Proposition 7.17). Let us consider  $\alpha, \delta\alpha, \beta, \delta\beta$  and times  $t$  in a neighborhood of 0, as in the statement above. We define  $\alpha_t = \alpha + t\delta\alpha, \beta_t = \beta + t\delta\beta$  and the variation ratio  $\Delta_t$  given by

$$\Delta_t \stackrel{\text{def.}}{=} \frac{\text{OT}_\varepsilon(\alpha_t, \beta_t) - \text{OT}_\varepsilon(\alpha, \beta)}{t}.$$

Using the very definition of  $\text{OT}_\varepsilon$  and the continuity property of Proposition 13, we now provide lower and upper bounds on  $\Delta_t$  as  $t$  goes to 0.

**Weak\* continuity.** As written in (13),  $\text{OT}_\varepsilon(\alpha, \beta)$  can be computed through a straightforward, *continuous* expression that does not depend on the value of the optimal dual potentials  $(f, g)$  at the anchor point  $x_o$ :

$$\text{OT}_\varepsilon(\alpha, \beta) = \langle \alpha, f \rangle + \langle \beta, g \rangle.$$

Combining this equation with Proposition 13 (that guarantees the *uniform* convergence of potentials for weakly converging sequences of probability measures) allows us to conclude.

**Lower bound.** First, let us remark that  $(f, g)$  is a *suboptimal* pair of dual potentials for  $\text{OT}_\varepsilon(\alpha_t, \beta_t)$ . Hence,

$$\begin{aligned} \text{OT}_\varepsilon(\alpha_t, \beta_t) &\geq \langle \alpha_t, f \rangle + \langle \beta_t, g \rangle \\ &\quad - \varepsilon \langle \alpha_t \otimes \beta_t, \exp\left(\frac{1}{\varepsilon}(f \oplus g - C)\right) - 1 \rangle \end{aligned}$$

and thus, since

$$\begin{aligned} \text{OT}_\varepsilon(\alpha, \beta) &= \langle \alpha, f \rangle + \langle \beta, g \rangle \\ &\quad - \varepsilon \langle \alpha \otimes \beta, \exp\left(\frac{1}{\varepsilon}(f \oplus g - C)\right) - 1 \rangle, \end{aligned}$$

one has

$$\begin{aligned} \Delta_t &\geq \langle \delta\alpha, f \rangle + \langle \delta\beta, g \rangle \\ &\quad - \varepsilon \langle \delta\alpha \otimes \beta + \alpha \otimes \delta\beta, \exp\left(\frac{1}{\varepsilon}(f \oplus g - C)\right) \rangle + o(1) \\ &\geq \langle \delta\alpha, f - \varepsilon \rangle + \langle \delta\beta, g - \varepsilon \rangle + o(1), \end{aligned}$$

since  $g$  and  $f$  satisfy the optimality equations (10).

**Upper bound.** Conversely, let us denote by  $(g_t, f_t)$  the optimal pair of potentials for  $\text{OT}_\varepsilon(\alpha_t, \beta_t)$  satisfying  $g_t(x_o) = 0$  for some arbitrary anchor point  $x_o \in \mathcal{X}$ . As  $(f_t, g_t)$  are suboptimal potentials for  $\text{OT}_\varepsilon(\alpha, \beta)$ , we get that

$$\begin{aligned} \text{OT}_\varepsilon(\alpha, \beta) &\geq \langle \alpha, f_t \rangle + \langle \beta, g_t \rangle \\ &\quad - \varepsilon \langle \alpha \otimes \beta, \exp\left(\frac{1}{\varepsilon}(f_t \oplus g_t - C)\right) - 1 \rangle \end{aligned}$$

and thus, since

$$\begin{aligned} \text{OT}_\varepsilon(\alpha_t, \beta_t) &= \langle \alpha_t, f_t \rangle + \langle \beta_t, g_t \rangle \\ &\quad - \varepsilon \langle \alpha_t \otimes \beta_t, \exp\left(\frac{1}{\varepsilon}(f_t \oplus g_t - C)\right) - 1 \rangle, \end{aligned}$$

$$\begin{aligned} \Delta_t &\leq \langle \delta\alpha, f_t \rangle + \langle \delta\beta, g_t \rangle \\ &\quad - \varepsilon \langle \delta\alpha \otimes \beta_t + \alpha_t \otimes \delta\beta, \exp\left(\frac{1}{\varepsilon}(f_t \oplus g_t - C)\right) \rangle + o(1) \\ &\leq \langle \delta\alpha, f_t - \varepsilon \rangle + \langle \delta\beta, g_t - \varepsilon \rangle + o(1). \end{aligned}$$

**Conclusion.** Now, let us remark that as  $t$  goes to 0

$$\alpha + t\delta\alpha \rightharpoonup \alpha \quad \text{and} \quad \beta + t\delta\beta \rightharpoonup \beta.$$

Thanks to Proposition 13, we thus know that  $f_t$  and  $g_t$  converge uniformly towards  $f$  and  $g$ . Combining the lower and upper bound, we get

$$\begin{aligned} \Delta_t &\xrightarrow{t \rightarrow 0} \langle \delta\alpha, f - \varepsilon \rangle + \langle \delta\beta, g - \varepsilon \rangle = \langle \delta\alpha, f \rangle + \langle \delta\beta, g \rangle, \end{aligned}$$

since  $\delta\alpha$  and  $\delta\beta$  both have an overall mass that sums up to zero.

### B.3 Proof of Proposition 3

The definition of  $\text{OT}_\varepsilon(\alpha, \alpha)$  is that

$$\text{OT}_\varepsilon(\alpha, \alpha) = \max_{(f,g) \in \mathcal{C}(\mathcal{X})^2} \langle \alpha, f + g \rangle - \varepsilon \langle \alpha \otimes \alpha, e^{(f \oplus g - C)/\varepsilon} - 1 \rangle.$$

**Reduction of the problem.** Thanks to the symmetry of this concave problem with respect to the variables  $f$  and  $g$ , we know that there exists a pair  $(f, g = f)$  of optimal potentials on the diagonal, and

$$\text{OT}_\varepsilon(\alpha, \alpha) = \max_{f \in \mathcal{C}(\mathcal{X})} 2\langle \alpha, f \rangle - \varepsilon \langle \alpha \otimes \alpha, e^{(f \oplus f - C)/\varepsilon} - 1 \rangle.$$

Thanks to the density of continuous functions in the set of simple measurable functions, just as in the proof of Proposition 7, we show that this maximization can be done in the full set of measurable functions  $\mathcal{F}_b(\mathcal{X}, \mathbb{R})$ :

$$\begin{aligned} \text{OT}_\varepsilon(\alpha, \alpha) &= \max_{f \in \mathcal{F}_b(\mathcal{X}, \mathbb{R})} 2\langle \alpha, f \rangle \\ &\quad - \varepsilon \langle \alpha \otimes \alpha, e^{(f \oplus f - C)/\varepsilon} - 1 \rangle \\ &= \max_{f \in \mathcal{F}_b(\mathcal{X}, \mathbb{R})} 2\langle \alpha, f \rangle \\ &\quad - \varepsilon \langle \exp(f/\varepsilon)\alpha, k_\varepsilon \star \exp(f/\varepsilon)\alpha \rangle + \varepsilon, \end{aligned}$$

where  $\star$  denotes the smoothing (convolution) operator defined through

$$[k \star \mu](x) = \int_{\mathcal{X}} k(x, y) d\mu(y)$$

for  $k \in \mathcal{C}(\mathcal{X} \times \mathcal{X})$  and  $\mu \in \mathcal{M}^+(\mathcal{X})$ .

**Optimizing on measures.** Through a change of variables

$$\mu = \exp(f/\varepsilon)\alpha \quad \text{i.e.} \quad f = \varepsilon \log \frac{d\mu}{d\alpha},$$

keeping in mind that  $\alpha$  is a probability measure, we then get that

$$\begin{aligned} \text{OT}_\varepsilon(\alpha, \alpha) &= \varepsilon \max_{\mu \in \mathcal{M}^+(\mathcal{X}), \alpha \ll \mu \ll \alpha} 2\langle \alpha, \log \frac{d\mu}{d\alpha} \rangle \\ &\quad - \langle \mu, k_\varepsilon \star \mu \rangle + 1 \\ -\frac{1}{2}\text{OT}_\varepsilon(\alpha, \alpha) &= \varepsilon \min_{\mu \in \mathcal{M}^+(\mathcal{X}), \alpha \ll \mu \ll \alpha} \langle \alpha, \log \frac{d\mu}{d\alpha} \rangle \\ &\quad + \frac{1}{2}\langle \mu, k_\varepsilon \star \mu \rangle - \frac{1}{2}, \end{aligned}$$

where we optimize on positive measures  $\mu \in \mathcal{M}^+(\mathcal{X})$  such that  $\alpha \ll \mu$  and  $\mu \ll \alpha$ .

**Expansion of the problem.** As  $k_\varepsilon(x, y) = \exp(-C(x, y)/\varepsilon)$  is positive for all  $x$  and  $y$  in  $\mathcal{X}$ , we can remove the  $\mu \ll \alpha$  constraint from the optimization problem:

$$\begin{aligned} -\frac{1}{2}\text{OT}_\varepsilon(\alpha, \alpha) &= \varepsilon \min_{\mu \in \mathcal{M}^+(\mathcal{X}), \alpha \ll \mu} \langle \alpha, \log \frac{d\mu}{d\alpha} \rangle \\ &\quad + \frac{1}{2}\langle \mu, k_\varepsilon \star \mu \rangle - \frac{1}{2}. \end{aligned}$$

Indeed, restricting a positive measure  $\mu$  to the support of  $\alpha$  lowers the right-hand term  $\langle \mu, k_\varepsilon \star \mu \rangle$  without having any influence on the density of  $\alpha$  with respect to  $\mu$ . Finally, let us remark that the  $\alpha \ll \mu$  constraint is already encoded in the  $\log \frac{d\mu}{d\alpha}$  operator, which blows up to infinity if  $\alpha$  has no density with respect to  $\mu$ ; in all in all, we thus have:

$$\begin{aligned} F_\varepsilon(\alpha) &= -\frac{1}{2}\text{OT}_\varepsilon(\alpha, \alpha) \\ &= \varepsilon \min_{\mu \in \mathcal{M}^+(\mathcal{X})} \langle \alpha, \log \frac{d\mu}{d\alpha} \rangle + \frac{1}{2}\langle \mu, k_\varepsilon \star \mu \rangle - \frac{1}{2}, \end{aligned}$$

which is the desired result.

**Existence of the optimal measure  $\mu$ .** In the expression above, the existence of an optimal  $\mu$  is given as a consequence of the well-known fact from OT theory that optimal dual potentials  $f$  and  $g$  exist, so that the dual OT problem (8) is a max and not a mere supremum. Nevertheless, since this property of  $F_\varepsilon$  is key to the metrization of the convergence in law by Sinkhorn divergences, let us endow it with a direct, alternate proof:

**Proposition 14.** *For any  $\alpha \in \mathcal{M}_1^+(\mathcal{X})$ , assuming that  $\mathcal{X}$  is compact, there exists a unique  $\mu_\alpha \in \mathcal{M}^+(\mathcal{X})$  such that*

$$F_\varepsilon(\alpha) = \varepsilon \left[ \langle \alpha, \log \frac{d\mu_\alpha}{d\alpha} \rangle + \frac{1}{2}\langle \mu_\alpha, k_\varepsilon \star \mu_\alpha \rangle - \frac{1}{2} \right].$$

Moreover,  $\alpha \ll \mu_\alpha \ll \alpha$ .

*Proof.* Notice that for  $(\alpha, \mu) \in \mathcal{M}_1^+(\mathcal{X}) \times \mathcal{M}^+(\mathcal{X})$ ,

$$\begin{aligned} E_\varepsilon(\alpha, \mu) &\stackrel{\text{def.}}{=} \langle \alpha, \log \frac{d\mu}{d\alpha} \rangle + \frac{1}{2}\langle \mu, k_\varepsilon \star \mu \rangle \\ &= \text{KL}(\alpha, \mu) + \langle \alpha - \mu, 1 \rangle + \frac{1}{2}\|\mu\|_{k_\varepsilon}^2 - \frac{1}{2}. \end{aligned}$$

Since  $C$  is bounded on the compact set  $\mathcal{X} \times \mathcal{X}$  and  $\alpha$  is a probability measure, we can already say that

$$\frac{1}{\varepsilon}F_\varepsilon(\alpha) \leq E_\varepsilon(\alpha, \alpha) - \frac{1}{2} = \frac{1}{2}\langle \alpha \otimes \alpha, e^{-C/\varepsilon} \rangle - \frac{1}{2} < +\infty.$$

**Upper bound on the mass of  $\mu$ .** Since  $\mathcal{X} \times \mathcal{X}$  is compact and  $k_\varepsilon(x, y) > 0$ , there exists  $\eta > 0$  such that  $k(x, y) > \eta$  for all  $x$  and  $y$  in  $\mathcal{X}$ . We thus get

$$\|\mu\|_{k_\varepsilon}^2 \geq \langle \mu, 1 \rangle^2 \eta$$

and show that

$$\begin{aligned} E_\varepsilon(\alpha, \mu) &\geq \langle \alpha - \mu, 1 \rangle + \frac{1}{2} \|\mu\|_{k_\varepsilon}^2 - \frac{1}{2} \\ &\geq \langle \mu, 1 \rangle (\langle \mu, 1 \rangle \eta - 1) - \frac{1}{2}. \end{aligned}$$

As we build a minimizing sequence  $(\mu_n)$  for  $F_\varepsilon(\alpha)$ , we can thus assume that  $\langle \mu_n, 1 \rangle$  is uniformly bounded by some constant  $M > 0$ .

**Weak continuity.** Crucially, the Banach-Alaoglu theorem asserts that

$$\{\mu \in \mathcal{M}^+(\mathcal{X}) \mid \langle \mu, 1 \rangle \leq M\}$$

is weakly compact; we can thus extract a weakly converging subsequence  $\mu_{n_k} \rightharpoonup \mu_\infty$  from the minimizing sequence  $(\mu_n)$ . Using Proposition 8 and the fact that  $k_\varepsilon$  is continuous on  $\mathcal{X} \times \mathcal{X}$ , we show that  $\mu \mapsto E_\varepsilon(\alpha, \mu)$  is a weakly lower semi-continuous function:  $\mu_\infty = \mu_\alpha$  realizes the minimum of  $E_\varepsilon$  and we get our existence result.

**Uniqueness.** We assumed that our kernel  $k_\varepsilon$  is *positive universal*. The squared norm  $\mu \mapsto \|\mu\|_{k_\varepsilon}^2$  is thus a strictly convex functional and using Proposition 6, we can show that  $\mu \mapsto E_\varepsilon(\alpha, \mu)$  is *strictly* convex. This ensures that  $\mu_\alpha$  is uniquely defined.  $\square$

#### B.4 Proof of Proposition 4

Let us take a pair of measures  $\alpha_0 \neq \alpha_1$  in  $\mathcal{M}_1^+(\mathcal{X})$ , and  $t \in (0, 1)$ ; according to Proposition 14, there exists a pair of measures  $\mu_0, \mu_1$  in  $\mathcal{M}^+(\mathcal{X})$  such that

$$\begin{aligned} (1-t)F_\varepsilon(\alpha_0) + tF_\varepsilon(\alpha_1) &= \varepsilon(1-t)E_\varepsilon(\alpha_0, \mu_0) + \varepsilon tE_\varepsilon(\alpha_1, \mu_1) \\ &> \varepsilon E_\varepsilon((1-t)\alpha_0 + t\alpha_1, (1-t)\mu_0 + t\mu_1) \\ &\geq F_\varepsilon((1-t)\alpha_0 + t\alpha_1), \end{aligned}$$

which is enough to conclude. To show the strict inequality, let us remark that

$$\begin{aligned} (1-t)E_\varepsilon(\alpha_0, \mu_0) + tE_\varepsilon(\alpha_1, \mu_1) \\ = E_\varepsilon((1-t)\alpha_0 + t\alpha_1, (1-t)\mu_0 + t\mu_1) \end{aligned}$$

would imply that  $\mu_0 = \mu_1$ , since  $\mu \mapsto \|\mu\|_{k_\varepsilon}^2$  is strictly convex. As  $\alpha \mapsto \text{KL}(\alpha, \beta)$  is strictly convex on the set of measures  $\alpha$  that are absolutely continuous with respect to  $\beta$ , we would then have  $\alpha_0 = \alpha_1$  and a contradiction with our first hypothesis.

#### B.5 Proof of the Metrization of the Convergence in Law

The regularized OT cost is weakly continuous, and the uniform convergence for dual potentials ensures that

$H_\varepsilon$  and  $S_\varepsilon$  are both continuous too. Paired with (6), this property guarantees the convergence towards 0 of the Hausdorff and Sinkhorn divergences, as soon as  $\alpha_n \rightharpoonup \alpha$ .

Conversely, let us assume that  $S_\varepsilon(\alpha_n, \alpha) \rightarrow 0$  (resp.  $H_\varepsilon(\alpha_n, \alpha)$ ). Any weak limit  $\alpha_{n_\infty}$  of a subsequence  $(\alpha_{n_k})_k$  is equal to  $\alpha$ : since our divergence is weakly continuous, we have  $S_\varepsilon(\alpha_{n_\infty}, \alpha) = 0$  (resp.  $H_\varepsilon(\alpha_{n_\infty}, \alpha)$ ), and positive definiteness holds through (6).

In the meantime, since  $\mathcal{X}$  is compact, the set of probability Radon measures  $\mathcal{M}_1^+(\mathcal{X})$  is sequentially compact for the weak- $\star$  topology.  $\alpha_n$  is thus a compact sequence with a unique adherence value: it converges, towards  $\alpha$ .