

470 **A ADDITIONAL PROOFS**

471 **A.1 Proof of Lemma III (strongly convex case)**

472 *Proof.* The major part of the proof is adapted from [Muzellec et al \[2021\]](#), Lemma 3.1]. We denote
 473 $T = \nabla f$ and T_0 an OT map between μ and ν such that there exists a convex potential f_0 verifying
 474 $\nabla f_0 = T_0$. By definition of the Fenchel-Legendre transform, we have for any convex function f
 475 and point x

$$f(x) + f^*(\nabla f(x)) = x^\top \nabla f(x), \quad (2)$$

476 Integrating this relation over μ for the optimal potential f_0 yields

$$\langle f_0, \mu \rangle + \langle f_0^* \circ T_0, \mu \rangle = \int x^\top T_0(x) d\mu(x). \quad (3)$$

477 Using the property $T_\#(\mu) = \nu$, we obtain that the r.h.s. is equal to J_0 . We use this same property to
 478 re-write $J(f) = \int f(x) + f^*(T_0(x)) d\mu(x)$. Finally, using the Legendre identity stated above, we
 479 have $J(f) = \int f^*(T_0(x)) - f^*(\nabla f(x)) + x^\top \nabla f(x) d\mu(x)$, which leads to

$$J(f) - J_0 = \int f^*(T_0(x)) - f^*(\nabla f(x)) - (T_0(x) - \nabla f(x))^\top x d\mu(x). \quad (4)$$

480 Recalling that $\partial f^*(\nabla f(x)) = x$, where f^* is a subgradient of f^* , we identify in the integrand a
 481 Bregman divergence $D_{f^*}(T_0(x), \nabla f(x))$ where for a convex function h the Bregman divergence
 482 $D_h(y, x) = h(y) - h(x) - \partial h(x)^\top (y - x)$. When f is assumed γ -strongly convex, f^* is $\frac{1}{\gamma}$ -smooth
 483 and $D_{f^*}(T_0(x), \nabla f(x))$ is upper-bounded by $\frac{1}{2\gamma} \|T(x) - T_0(x)\|^2$ which yields

$$J(f) - J_0 \leq \frac{1}{2\gamma} \int \|T(x) - T_0(x)\|^2 d\mu(x). \quad (5)$$

484 Conversely, when f is assumed M -smooth, f^* is $\frac{1}{2M}$ -strongly convex and $D_{f^*}(T_0(x), \nabla f(x))$ is
 485 lower-bounded by $\frac{1}{2M} \|T(x) - T_0(x)\|^2$ which yields

$$J(f) - J_0 \geq \frac{1}{2M} \int \|T(x) - T_0(x)\|^2 d\mu(x). \quad (6)$$

486 □

487 **A.2 Proof of Prop. III**

488 *Proof.* Define the potential $g_0(x) = |x| + \frac{x^2}{2}$ and for $0 \leq \lambda \leq \frac{1}{2}$, define the translated potential
 489 $g_\lambda = g_0(\cdot - \lambda)$. Let us start by computing the Legendre transform of g_0 . The Legendre transform
 490 of g_0 is defined for all y as

$$g_0^*(y) = \sup_{x \in \mathbb{R}} xy - g_0(x). \quad (7)$$

491 Since $g_0(x) \geq |x|$ for all x , then if $y \in [-1, 1]$ then $g_0^*(y) = 0$. If $y > 1$, since g_0 is pair and positive,
 492 the maximum is attained on \mathbb{R}^+ . Denoting $\phi(x, y) = xy - g_0(x)$, we have that $\phi(\cdot, y)$ increases
 493 between $[0, y - 1]$ and decreases between $[y - 1, +\infty[$ hence the maximum is attained in $x = y - 1$
 494 which yields $g_0^*(y) = \frac{(y-1)^2}{2}$. Conversely, if $y < -1$ we have $g_0^*(y) = \frac{(y+1)^2}{2}$. From this result, we
 495 can compute g_λ^* in virtue of the relation $g_\lambda^*(y) = g^*(y) + \lambda y$.

496 Let us now compute the semi-dual $J(g_\lambda)$. The first term is given by

$$\int_{-\frac{1}{2}}^{\frac{1}{2}} g_\lambda(x) dx = \int_{-\frac{1}{2}}^\lambda g(x - \lambda) dx + \int_\lambda^{\frac{1}{2}} g(x - \lambda) dx \quad (8)$$

$$= \int_{-\frac{1}{2}}^\lambda \frac{(x - \lambda)^2}{2} + (\lambda - x) dx + \int_\lambda^{\frac{1}{2}} \frac{(x - \lambda)^2}{2} + (x - \lambda) dx \quad (9)$$

$$= \left[\frac{(\cdot - \lambda)^3}{6} \right]_{-\frac{1}{2}}^{\frac{1}{2}} + \left[-\frac{(\cdot - \lambda)^2}{2} \right]_{-\frac{1}{2}}^\lambda + \left[-\frac{(\lambda - \cdot)^2}{2} \right]_\lambda^{\frac{1}{2}} \quad (10)$$

$$= \frac{(\frac{1}{2} - \lambda)^3 + (\frac{1}{2} + \lambda)^3}{6} + \frac{(\frac{1}{2} + \lambda)^2 + (\frac{1}{2} - \lambda)^2}{2} \quad (11)$$

$$= \frac{1}{4} + \frac{\frac{1}{4} + 3\lambda^2}{6} + \lambda^2. \quad (12)$$

497 The second term is given by

$$\int_{\mathbb{R}} g_\lambda^*(y) d(\nabla g_0)(\mu)(y) = \int_{-\frac{1}{2}}^{\frac{1}{2}} g_\lambda^*(\nabla g_0(y)) dy \quad (13)$$

$$= \int_{-\frac{1}{2}}^0 g_0^*(y - 1 - \lambda) dy + \int_0^{\frac{1}{2}} g_0^*(y + 1 - \lambda) dy \quad (14)$$

$$= \int_{-\frac{1}{2}}^0 \frac{(u - \lambda)^2}{2} du + \int_\lambda^{\frac{1}{2}} \frac{(u - \lambda)^2}{2} du \quad (15)$$

$$= \left[\frac{(\cdot - \lambda)^3}{6} \right]_{-\frac{1}{2}}^0 + \left[\frac{(u - \lambda)^3}{6} \right]_\lambda^{\frac{1}{2}} \quad (16)$$

$$= \frac{(\frac{1}{2} + \lambda)^3 + (\frac{1}{2} - \lambda)^3 - \lambda^3}{6} \quad (17)$$

$$= \frac{\frac{1}{4} + 3\lambda^2 - \lambda^3}{6}. \quad (18)$$

498 Hence the semi-dual is given by $J(g_\lambda) = \frac{1}{3} + 2\lambda^2 - \lambda^3$. Finally, let us compute the error e_μ

$$e_\mu(g_\lambda) = \int_{-\frac{1}{2}}^0 (\nabla g_0(x) - \nabla g_\lambda(x))^2 dx + \int_0^\lambda (\nabla g_0(x) - \nabla g_\lambda(x))^2 dx + \int_\lambda^{\frac{1}{2}} (\nabla g_0(x) - \nabla g_\lambda(x))^2 dx \quad (19)$$

$$= \frac{\lambda^2}{2} + \lambda(2 + \lambda)^2 + \lambda^2(\frac{1}{2} - \lambda) \quad (20)$$

$$= 4\lambda + 5\lambda^2. \quad (21)$$

499

□

500 A.3 Proof of Prop. [2](#)

501 *Proof.* We begin with splitting $\hat{J}(f_{i_0}) - J_0$ in non-stochastic and stochastic terms

$$\hat{J}(f_{i_0}) - J_0 = J(f_{i_0}) - J_0 + \hat{J}(f_{i_0}) - J(f_{i_0}), \quad (22)$$

502 where we denoted \hat{J} the empirical semi-dual $J_{\hat{\mu}, \hat{\nu}}$. Using Lemma [3](#), we get the lower bound

$$\hat{J}(f_{i_0}) - J_0 \geq \frac{1}{2M} e_\mu(f_{i_0}) + \hat{J}(f_{i_0}) - J(f_{i_0}). \quad (23)$$

503 By construction, f_{i_0} verifies for all $1 \leq i \leq p$

$$\begin{aligned} \hat{J}(f_{i_0}) - J_0 &\leq \hat{J}(f_i) - J_0 \\ &= J(f_i) - J_0 + \hat{J}(f_i) - J(f_i). \end{aligned}$$

504 Picking $i = i_1$ and using Lemma [11](#), we obtain

$$\hat{J}(f_{i_0}) - J_0 \leq \frac{1}{2\gamma} e_\mu(f_{i_1}) + \hat{J}(f_{i_1}) - J(f_{i_1}). \quad (24)$$

505 Equations [\(23\)](#) and [\(24\)](#) give

$$e_\mu(f_{i_0}) \leq \frac{M}{\gamma} e_\mu(f_{i_1}) + 2M(\hat{J}(f_{i_1}) - J(f_{i_1})) \quad (25)$$

$$+ 2M(J(f_{i_0}) - \hat{J}(f_{i_0})). \quad (26)$$

506 The Hoeffding lemma gives for all $t > 0$

$$\mathbb{P}(\langle f_i, \hat{\mu} - \mu \rangle \geq t) \leq \exp\left(-\frac{2nt^2}{\|f_i\|_{osc,X}^2}\right). \quad (27)$$

507 We place ourselves on the event

$$A = (\langle f_{i_1}, \hat{\mu} - \mu \rangle \geq t) \cup (\langle f_{i_1}^*, \hat{\nu} - \nu \rangle \geq t) \cup (\langle f_{i_0}, \mu - \hat{\mu} \rangle \geq t) \cup (\langle f_{i_0}^*, \nu - \hat{\nu} \rangle \geq t). \quad (28)$$

508 We want to set $\mathbb{P}(A) \leq \delta$. By triangle inequality, we get the upper-bound

$$\mathbb{P}(A) \leq 4 \exp\left(-\frac{2nt^2}{C^2}\right), \quad (29)$$

509 where $C = \max(C_{i_0}, C_{i_1})$ and C_i defined as $C_i = \max(\|f_i\|_{X,o}, \|f_i^*\|_{Y,o})$. Hence setting, $t =$

510 $C\sqrt{\frac{\ln(4/\delta)}{2n}}$, we have with probability at least $1 - \delta$

$$e_\mu(f_{i_0}) \leq \frac{M}{\gamma} e_\mu(f_{i_1}) + 8MC\sqrt{\frac{\ln(4/\delta)}{2n}}. \quad (30)$$

511

□

512 **A.4 Proof of Prop. [3](#)**

513 *Proof.* Take $\mu \sim [0, 1]$ and $f_0 \equiv 0$ and let us compute the error for $g = M\frac{x^2}{2}$.

$$e_\mu(g) = \int_0^1 (Mx)^2 dx \quad (31)$$

$$= \frac{M^2}{3}. \quad (32)$$

514 Conversely, defining $h_\epsilon = \gamma\frac{x^2}{2} + (\epsilon + \alpha_{M,\gamma})x$ with

$$\alpha_{M,\gamma} = \frac{\gamma}{2} \left[\sqrt{1 + \frac{4(M-\gamma)}{3\gamma}} - 1 \right], \quad (33)$$

515 the error $e_\mu(h_\epsilon)$ is given by

$$e_\mu(h) = \int_0^1 (\gamma x + (\epsilon + \alpha_{M,\gamma}))^2 dx \quad (34)$$

$$= \frac{\gamma^2}{3} + \gamma(\alpha_{M,\gamma} + \epsilon) + (\alpha_{M,\gamma} + \epsilon)^2 \quad (35)$$

$$= \frac{\gamma^2}{3} + \gamma\epsilon + \gamma\alpha_{M,\gamma} + \epsilon^2 + 2\epsilon\alpha_{M,\gamma} + \alpha_{M,\gamma}^2 \quad (36)$$

$$= \frac{\gamma^2}{3} + \gamma\epsilon + \frac{\gamma^2}{2} \sqrt{1 + \frac{4(M-\gamma)}{3\gamma}} - \frac{\gamma^2}{2} + \epsilon^2 + 2\epsilon\alpha_{M,\gamma} + \frac{\gamma^2}{4} \left[2 + \frac{4(M-\gamma)}{3\gamma} - 2\sqrt{1 + \frac{4(M-\gamma)}{3\gamma}} \right] \quad (37)$$

$$= \frac{\gamma^2}{3} + \gamma\epsilon + \epsilon^2 + 2\epsilon\alpha_{M,\gamma} + \frac{\gamma(M-\gamma)}{3} \quad (38)$$

$$= \frac{M\gamma}{3} + \gamma\epsilon + \epsilon^2 - \gamma\epsilon + \gamma\epsilon \sqrt{1 + \frac{4(M-\gamma)}{3\gamma}} \quad (39)$$

$$= \frac{M\gamma}{3} \left[1 + \frac{3\epsilon^2}{M\gamma} + \frac{3\epsilon}{M} \sqrt{1 + \frac{4(M-\gamma)}{3\gamma}} \right]. \quad (40)$$

516 In particular, we obtain $\frac{e_\mu(g)}{e_\mu(h_\epsilon)} = \frac{M}{\gamma} \times \frac{1}{1 + \frac{3\epsilon}{M} \left[\frac{\epsilon}{\gamma} + \sqrt{1 + \frac{4(M-\gamma)}{3\gamma}} \right]} \xrightarrow{\epsilon \rightarrow 0} \frac{M}{\gamma}$. Now, let us compute the

517 semi-duals $J(g)$ and $J(h_\epsilon)$. Since $f_0 \equiv 0$, we have $\nu = \delta_0$ a Dirac mass in 0. Hence we simply
518 need to compute the Legendre transform of g and h_ϵ in 0

$$g^*(0) = \sup_x -M \frac{x^2}{2} \quad (41)$$

$$= 0, \quad (42)$$

519 and

$$h_\epsilon^*(0) = \sup_x -\gamma \frac{x^2}{2} - (\epsilon + \alpha_{M,\gamma})x \quad (43)$$

$$= \frac{(\epsilon + \alpha_{M,\gamma})^2}{2\gamma}. \quad (44)$$

520 Hence we obtain

$$J(g) = \int_0^1 M \frac{x^2}{2} dx \quad (45)$$

$$= \frac{M}{6}, \quad (46)$$

521 and

$$J(h_\epsilon) = \int_0^1 \gamma \frac{x^2}{2} + (\epsilon + \alpha_{M,\gamma})x dx + \frac{(\epsilon + \alpha_{M,\gamma})^2}{2\gamma} \quad (47)$$

$$= \frac{\gamma}{6} + \frac{\epsilon + \alpha_{M,\gamma}}{2} + \frac{(\epsilon + \alpha_{M,\gamma})^2}{2\gamma} \quad (48)$$

$$= \frac{\gamma}{6} + \frac{\epsilon + \alpha_{M,\gamma}}{2} + \frac{\epsilon^2 + 2\epsilon\alpha_{M,\gamma} + \alpha_{M,\gamma}^2}{2\gamma} \quad (49)$$

$$= \frac{1}{2\gamma} \left(\frac{\gamma^2}{3} + \gamma(\epsilon + \alpha_{M,\gamma}) + \epsilon^2 + 2\epsilon\alpha_{M,\gamma} + \alpha_{M,\gamma}^2 \right). \quad (50)$$

522 We recognize between brackets the same expression as $e_\mu(h_\epsilon)$ in Equation (52) hence we obtain

$$J(h_\epsilon) = \frac{M}{6} \left[1 + \frac{3\epsilon^2}{M\gamma} + \frac{3\epsilon}{M} \sqrt{1 + \frac{4(M-\gamma)}{3\gamma}} \right]. \quad (51)$$

523

□

524 **A.5 Proof of Prop. 4**

525 *Proof.* Applying Lemma 11, we have for all $\delta > 0$ the inequality $J(Q_\delta(f)) - J_0 \leq \frac{1}{2\delta} e_\mu(Q_\delta(f))$.
 526 The right hand side is decomposed in $\langle Q_\delta(f), \mu \rangle + \langle Q_\delta(f)^*, \nu \rangle$. The first term is simply $\langle f, \mu \rangle +$
 527 $\delta \langle q, \mu \rangle$. For the second-term we use the standard result $Q_\delta(f)^* = M_\delta(f)$ the Moreau-Yosida
 528 transform of f reading $M_\tau(f) = \inf_y f(y) + \frac{q(x-y)}{\tau}$. If f^* is L-Lipschitz on the support of ν , we
 529 have the lower bound $M_\delta(f) \geq f - \frac{L^2\delta}{2}$. Hence we recover

$$J(f) - J_0 - \frac{L^2\delta}{2} + \delta \langle q, \mu \rangle \leq \frac{1}{2\delta} e_\mu(Q_\delta(f)). \quad (52)$$

530 The term $e_\mu(Q_\delta(f))$ is upper-bounded by $2(e_\mu(f) + 2\delta^2 \langle q, \mu \rangle)$ which gives $J(f) - J_0 \leq \frac{e_\mu(f)}{\delta} +$
 531 $\delta \langle q, \mu \rangle + \frac{L^2\delta}{2}$. Optimizing on δ leads to

$$J(f) - J_0 \leq 2\sqrt{e_\mu(f) \left(\frac{L^2}{2} + \langle q, \mu \rangle \right)}. \quad (53)$$

532 □

533 **A.6 Proof of Prop. 6**

534 *Proof.* Recall that the Fenchel-Legendre of a standard Log-Sum-Exp function $\text{LSE}(x) =$
 535 $\log(\sum_{i=1}^n e^{x_i})$ is given by

$$\text{LSE}^*(y) = \sum_{i=1}^n y_i \log(y_i) + \iota(y \in \mathcal{S}_n) \quad (54)$$

$$= -\text{Ent}(y) + \iota(y \in \mathcal{S}_n), \quad (55)$$

536 where \mathcal{S}_n is the probability simplex. More generally, defining $\text{LSE}_b(x) = \log(\sum_{i=1}^n e^{x_i + b_i})$, using
 537 the fact that $f^*(\cdot + \tau) = f^*(\cdot) - \tau^\top$, we have

$$(\text{LSE}_b)^*(y) = -\text{Ent}(y) - b^\top y + \iota(y \in \mathcal{S}_n). \quad (56)$$

538 At the optimum, for empirical measures $\hat{\mu} = \frac{1}{m} \sum_{i=1}^m \delta_{x_i}$, $\hat{\nu} = \frac{1}{n} \sum_{i=1}^n \delta_{y_i}$ the empirical Sinkhorn
 539 Kantorovitch potentials $(\hat{\phi}_\varepsilon, \hat{\psi}_\varepsilon)$ are linked as

$$\hat{\phi}_\varepsilon(x) = -\varepsilon \log \left(\frac{1}{n} \sum_{i=1}^n e^{2\frac{\hat{\psi}_\varepsilon(y_i) - \|x - y_i\|^2}{2\varepsilon}} \right), \quad (57)$$

540 hence the Sinkhorn Brenier potential f_ε can be written as

$$f_\varepsilon(x) = \varepsilon \text{LSE}_{b_\varepsilon}(C_\varepsilon x), \quad (58)$$

541 where we defined

$$\begin{cases} C_\varepsilon = (\frac{y_i}{\varepsilon})_{1 \leq i \leq n} \in \mathbb{R}^{n \times d} \\ b_{\varepsilon, n} = (\frac{2\hat{\psi}_\varepsilon(y_i) - \|y_i\|^2}{2\varepsilon} - \log(n))_{1 \leq i \leq n} \in \mathbb{R}^n \end{cases}. \quad (59)$$

542 Now recall that

- 543 • $(\varepsilon f(\cdot))^* = \varepsilon f^*(\frac{\cdot}{\varepsilon})$.
- 544 • $\forall z, (f(A))^*(z) = \inf_{Ay=z} f^*(y)$.

545 Hence we can deduce

$$f_\varepsilon^*(y) = \varepsilon \inf_{C\Delta=y} -\text{Ent}(\Delta) - \Delta^\top b_{\varepsilon, n} + \iota(\Delta \in \mathcal{S}_n),$$

546 where $C \in \mathbb{R}^{n \times d}$ is the matrix of the samples (y_i) . In particular if f_ε^* is evaluated outside the convex
 547 hull of $\hat{\nu}$, it is infinite. Since ν has continuous density, there almost surely exists (y_0, r) , $r > 0$ such
 548 that $B(y_0, r) \subset \text{Supp}(\nu)$ and $B(y_0, r) \cap \text{Conv}(\hat{\nu}) = \emptyset$ where $\text{Conv}(\hat{\nu})$ is the convex hull of the
 549 samples $\hat{\nu}$. In particular, almost surely

$$\langle f_\varepsilon^*, \nu \rangle = +\infty. \quad (60)$$

550 □

551 **A.7 Proof of Prop. 4**

552 The proof is largely inspired from an article on the online blog of Francis Bach⁴.

553 Since the 2-self-concordance is scaling invariant, we shall simply prove that $f(x) = \text{LSE}_b(C)$ is
 554 $(2, D(C))$ self-concordant with $b \in \mathbb{R}_+^n$, $C \in \mathbb{R}^{n \times d}$ the matrix whose rows are centers $(c_i)_{1 \leq i \leq n}$
 555 and $D(C) = \max_{i,j} \|c_i - c_j\|$.

556 *Proof.* Defining the (non-normalized) distribution $\zeta = \frac{1}{n} \sum_{i=1}^n b_i \delta_{c_i}$, we can remark that f is the
 557 normalizing factor of the conditional exponential distribution

$$h(c|x) \propto e^{c^\top x} d\zeta(c) \quad (61)$$

$$= e^{c^\top x - f(x)} d\zeta(c). \quad (62)$$

558 The gradient of f is given by

$$\nabla f(x) = \frac{\int c e^{c^\top x} d\zeta(c)}{\int e^{c^\top x} d\zeta(c)} \quad (63)$$

$$= \mathbb{E}_h(c), \quad (64)$$

559 and using the results of Pistone and Wynn [1999], we have for higher order derivatives

$$\nabla^p f(x) = \mathbb{E}_h(\otimes_{j=1}^p (c - \nabla f(x))), \quad (65)$$

560 where for a vector $v \in \mathbb{R}^d$, $\otimes_{j=1}^p v$ is a tensor V_p in \mathbb{R}^{d^p} whose entries are $(v_{i_1} \times \dots \times v_{i_p})$. In
 561 particular, applying the formula for $p = 3$ and denoting $H = (c - \nabla f(x)) \otimes (c - \nabla f(x))$

$$\nabla^3 f(x) = \mathbb{E}_h[(c - \nabla f(x)) \otimes H]. \quad (66)$$

562 Using the linearity of the expectation, we have

$$|(\nabla^3 f(x)[v]u)^\top u| = |\mathbb{E}_h[(c - \nabla f(x))^\top v \times (Hu)^\top u]| \quad (67)$$

$$\leq \mathbb{E}_h[|(c - \nabla f(x))^\top v| \times |(Hu)^\top u|]. \quad (68)$$

563 Since $\nabla f(x) \in \text{Conv}(C)$, we have in particular that $\|c - \nabla f(x)\| \leq D(C)$. Furthermore since H
 564 is a positive matrix, we obtain the following upper-bound

$$|(\nabla^3 f(x)[v]u)^\top u| \leq D(C) \|v\| \mathbb{E}_h[(Hu)^\top u] \quad (69)$$

$$\leq D(C) \|v\| (\nabla^2 f(x)u)^\top u. \quad (70)$$

565 □

566 **A.8 Proof of Prop. 5**

567 *Proof.* The Sinkhorn Brenier empirical potentials are of the form $f_\varepsilon = \varepsilon \text{LSE}_{b_{\varepsilon,n}}(C_{\varepsilon,\cdot})$ where C_ε
 568 and $b_{\varepsilon,n}$ are defined in (59). Using the formulas from the previous proof, we simply have to bound
 569 $H_{c,x} = (c - \nabla f(x)) \otimes (c - \nabla f(x))$

$$u^\top H_{c,x} u = (u^\top (c - \nabla f(x)))^2 \quad (71)$$

$$\leq \|u\|_2^2 \|c - \nabla f(x)\|_2^2. \quad (72)$$

570 Since $\nabla f(x)$ is in the convex hull of $\frac{\hat{v}}{\varepsilon}$ and $c \in \text{Supp}(\frac{\hat{v}}{\varepsilon})$, we deduce that $\|H_{c,x}\|_{op} \leq \frac{D^2(\hat{v})}{\varepsilon^2}$, where
 571 $\|\cdot\|_{op}$ is the spectral norm. In particular $\|\nabla^2 f(x)\|_{op} \leq \frac{D^2(\hat{v})}{\varepsilon}$. □

572 **B MISCELLANEOUS**

573 **B.1 DA experiment**

574 We present here the results in the Domain Adaptation experiment where the source terms are (D) and
 575 (W) respectively. The results are displayed on Table 5: again, the best accuracy for the downstream
 576 classification task is not correlated with the minimization of the semi-dual, in particular the best OT
 577 maps are not suited for label transfer.

⁴<https://francisbach.com/self-concordant-analysis-for-logistic-regression/>

	ICNN		Sinkhorn		SSNB	
	acc(f_{i_1})	acc(f_{i_0})	acc(f_{i_1})	acc(f_{i_0})	acc(f_{i_1})	acc(f_{i_0})
D/A	0.5	0.47 (2/48)	0.91	0.78 (5/5)	0.91	0.84 (11/11)
D/C	0.54	0.43 (2/48)	0.83	0.74 (5/5)	0.83	0.75 (9/11)
D/W	0.52	0.28 (11/48)	0.96	0.85 (4/5)	0.99	0.95 (8/11)
W/A	0.48	0.25 (17/48)	0.89	0.78 (4/5)	0.87	0.77 (11/11)
W/C	0.4	0.2 (13/48)	0.77	0.73 (4/5)	0.78	0.74 (10/11)
W/D	0.62	0.51 (2/48)	0.95	0.9 (3/5)	1.0	1.0 (1/11)

Table 5: Potential Selection for Domain-Adaptation. The column $\text{acc}(f_{i_1})$ corresponds to the best (highest) accuracy and $\text{acc}(f_{i_0})$ corresponds to the accuracy of the potential selected with the Brenier criterion. On this Table, the potentials are ranked with respect to the accuracy; the closer to one, the better the classification. In bold, the highest accuracy after being calibrated with the semi-dual.

578 B.2 SSNB algorithm

579 For $l < L$, the SSNB model is defined as

$$\inf_{f \in \mathcal{F}_{l,L}} W_2^2((\nabla f)_\#(\mu), \nu), \quad (73)$$

580 where $\mathcal{F}_{l,L}$ is the set of l -strongly convex, L -smooth functions. For empirical potentials $\hat{\mu} =$
581 $\frac{1}{n} \sum_{i=1}^n \delta_{x_i}$ and $\hat{\nu} = \frac{1}{m} \sum_{i=1}^m \delta_{y_i}$, the authors propose to solve the non-convex problem (73) in
582 an alternate fashion: for a fixed $f \in \mathcal{F}_{l,L}$, they estimate the transport coupling $(P_{ij}) \in \mathbb{R}^{n \times m}$ from
583 $(\nabla f)_\#(\hat{\mu})$ to $\hat{\nu}$ by solving the associated linear program (or an entropic approximation) and then,
584 once the coupling is fixed, they estimate f (pointwise on $\hat{\mu}$) by solving

$$\min_{(z_1, \dots, z_n) \in \mathbb{R}^n \times d, u \in \mathbb{R}^n} \sum_{ij} P_{ij} \|z_i - y_j\|_2^2$$

$$\text{subject to } u_i \geq u_j + z_j^\top (x_i - x_j) + \frac{1}{2(1-l/L)} \left(\frac{1}{L} \|z_i - z_j\|^2 + \frac{1}{l} \|x_i - x_j\|_2^2 - \frac{2l}{L} (z_j - z_i)^\top (x_i - x_j) \right), \quad (74)$$

585 where $z_i = \nabla f(x_i)$ and $u_i = f(x_i)$. The problem above is a convex Quadratically Constrained
586 Quadratic Problem and can be numerically solved with CVXPY for instance. However, when such
587 an option is chosen the $n(n-1)$ constraints must be computed at each iterations which induces a
588 large overhead. Instead, we reformulate this problem as a standard linear conic problem of the form
589 $Ax - b \in \mathcal{K}$, with \mathcal{K} a fixed cone to be compiled only once.

590 **From QCQP to SOCP** First we show how to reformulate a (convex) QCQP without equality
591 constraints into an SOCP. The standard formulation of a QCQP is

$$\inf_x \frac{1}{2} x^\top Q_0 x + c_0^\top x$$

$$\text{s. t. } \frac{1}{2} x^\top Q_i x + c_i^\top x + r_i \leq 0, \quad i = 1, \dots, p. \quad (75)$$

592 Introducing the slack variables $(t_0, t_1, \dots, t_p) = \frac{1}{2}(x^\top Q_0 x, x^\top Q_1 x, \dots, x^\top Q_p x)$, we re-write the
593 problem as

$$\inf_{x,t} t_0 + c_0^\top x$$

$$\text{s. t. } t_i + c_i^\top x + r_i = 0, \quad i = 1, \dots, p$$

$$t_i \geq \frac{1}{2} x^\top Q_i x, \quad i = 0, \dots, p. \quad (76)$$

594 Decomposing Q_i as $Q_i = F_i^\top F_i$ with F_i having p rows, the constraint $t_i = \frac{1}{2} x^\top Q_i x$ becomes
595 $(1, t_i, F_i x) \in \mathcal{Q}_r^{d+2}$, where \mathcal{Q}_r^{d+2} is the rotated $(d+2)$ -dimensional Lorentz cone defined as

$$\mathcal{Q}_r^{d+2} = \{(x_1, x_2, \dots, x_{d+2}) \text{ s.t. } 2x_1x_2 \geq \sum_{k=1}^d x_{i+2}^2\}. \quad (77)$$

596 We obtain a MOSEK-friendly formulation of the QCQP as

$$\begin{aligned} & \inf_{x,t} t_0 + c_0^\top x \\ & \text{s. t. } t_i + c_i^\top x + r_i = 0, \quad i = 1, \dots, p \\ & (1, t_i, F_i x) \in \mathcal{Q}_r^{d+2}, \quad i = 0, \dots, p, \end{aligned} \quad (78)$$

597 which has the form $Ax - b \in \mathcal{K}$ where \mathcal{K} is a fixed product of Lorentz cone whose number and
598 dimensions solely depend on n and d in the case of SSNB. Hence we can compile \mathcal{K} only once for
599 fixed (n, d) , which allows us to considerably reduce the overhead.

600 **Decomposition of Q_{ij}** In the SSNB model the symmetric positive matrices $Q_{ij} \in \mathcal{S}_{n(d+1)}^+(\mathbb{R})$ are
601 defined up to a common scaling parameter as

$$\begin{cases} q_{kl} = 1 \text{ if } k = l \in \{di, \dots, (d+1)i\} \cup \{dj, \dots, (d+1)j\} \\ q_{kl} = -1 \text{ if } l = k + dj, k \in \{di, \dots, (d+1)i\} \\ q_{kl} = -1 \text{ if } k = l + dj, l \in \{di, \dots, (d+1)i\}. \end{cases} \quad (79)$$

602 The matrix Q_{ij} is factorized as $F_{ij}^\top F_{ij}$ with $F_{ij} \in \mathbb{R}^{d \times n(d+1)}$ defined as

$$\begin{cases} f_{kl} = 1 \text{ if } l = k + di, k \in \{1, \dots, d\} \\ f_{kl} = -1 \text{ if } l = k + dj, k \in \{1, \dots, d\}. \end{cases} \quad (80)$$

603 B.3 Models parameters

604 **ICNN** We used a 3-layers ICNN with softplus activations. The number of hidden neurons was
605 chosen in $\{64, 128, 256\}$, the soft convexity penalty for the potential g and the matching mo-
606 ment/variance penalty were both chosen in $\{0, 0.001, 0.01, 0.1\}$. As recommended by the authors,
607 the batch size was set to 60, the number of epochs was set to 60, the number of inner iterations
608 to approximate the conjugate was set to 25 and the learning rate is initially set to 1e-4 and is then
609 divided by 2 every 2-epochs.

610 To compute the semi-dual, we regularized the potential f by adding $\frac{\delta}{2}\|x\|^2$ with $\delta = 1e-3$. The
611 numerical optimization was done with SciPy with a stopping condition set to 0.001 ; for a lower
612 stopping criterion, the minimization would not converge.

613 **Sinkhorn** The temperature ε was chosen in $\{0.5, 0.1, 0.05, 0.01, 0.005\}$. We stopped the training
614 when the optimality conditions are almost met

$$\begin{cases} \langle |\phi_\varepsilon(\cdot) + \varepsilon \log(\int_y e^{\frac{\psi_\varepsilon(y) - c(\cdot, y)}{\varepsilon}} d\hat{\nu}(y))|, \hat{\mu} \rangle \leq 1e-5 \\ \langle |\psi_\varepsilon(\cdot) + \varepsilon \log(\int_x e^{\frac{\phi_\varepsilon(y) - c(x, \cdot)}{\varepsilon}} d\hat{\mu}(x))|, \hat{\nu} \rangle \leq 1e-5. \end{cases} \quad (81)$$

615 The resulting Sinkhorn Brenier potential \hat{f}_ε is regularized with $\frac{\delta}{2}\|x\|^2$, $\delta = 0.001$. When the semi-
616 dual is computed on a point y_i , the stopping criterion is given by

$$\|\nabla \hat{f}_\varepsilon(z_t) - y_i\| \leq 1e-5, \quad (82)$$

617 where z_t is the current point of the optimization at time step t .

618 **SSNB** The strong convexity parameter l is chosen in $\{0.2, 0.5, 0.7, 0.9\}$ and the smoothness pa-
619 rameter L is chosen in $\{0.2, 0.5, 0.7, 0.9, 1.2\}$ with $l < L$. The number of iterations in the alternate
620 minimization is set to 10. The conjugate is computed with a first order scheme with learning rate $\frac{1}{2L}$
621 and is stopped with the same criterion as above.

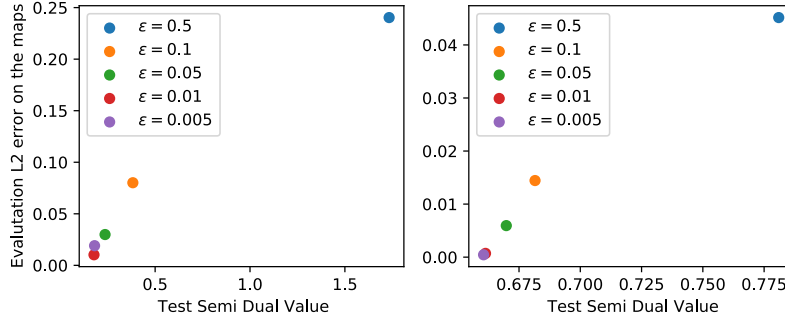


Figure 4: Empirical Semi-Dual against Quadratic Error on the Quadratic and Log-Sum-Exp experiments for the Sinkhorn model, $n = 10000$ and $d = 8$.

622 **B.4 Additional Experiment Sinkhorn**

623 We run 10 times the Quadratic and Log-Sum-Exp experiments with the Sinkhorn model but on
 624 $n = 10000$ points for the training of the model, the semi-dual and the computation of the error. The
 625 results are reported on Figure [B.4](#). Just as for SSNB, the semi-dual can accurately rank the potentials
 626 according to their error $e_\mu(f_i) = \int \|\nabla f_i(x) - T_0(x)\|_2^2 d\mu(x)$ where T_0 is the ground truth OT map.