



# On the Gradient Formula for learning Generative Models with Regularized Optimal Transport Costs

Antoine Houdard, Arthur Leclaire, Nicolas Papadakis, Julien Rabin

## ► To cite this version:

Antoine Houdard, Arthur Leclaire, Nicolas Papadakis, Julien Rabin. On the Gradient Formula for learning Generative Models with Regularized Optimal Transport Costs. Transactions on Machine Learning Research Journal, 2023. hal-03740368

**HAL Id: hal-03740368**

**<https://hal.science/hal-03740368>**

Submitted on 29 Jul 2022

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

---

# On the Gradient Formula for learning Generative Models with Regularized Optimal Transport Costs

Antoine Houdard  
Arthur Leclaire  
Nicolas Papadakis

*Univ. Bordeaux, Bordeaux INP, CNRS, IMB, UMR 5251,F-33400 Talence, France*

*Antoine.Houdard@math.u-bordeaux.fr  
Arthur.Leclaire@math.u-bordeaux.fr  
Nicolas.Papadakis@math.u-bordeaux.fr*

Julien Rabin

*Normandie Univ., ENSICAEN, CNRS, GREYC*

*Julien.Rabin@unicaen.fr*

## Abstract

The use of optimal transport costs for learning generative models has become popular with Wasserstein Generative Adversarial Networks (WGANs). Training a WGAN requires the computation of the differentiation of the optimal transport cost with respect to the parameters of the generative model. In this work, we provide sufficient conditions for the existence of a gradient formula in two different frameworks: the case of semi-discrete optimal transport (i.e. with a discrete target distribution) and the case of regularized optimal transport (i.e. with an entropic penalty). Both cases are based on the dual formulation of the transport cost, and the gradient formula involves a solution of the dual problem. The learning problem is addressed with an alternate algorithm, whose behavior is examined for the problem of MNIST digits generation. In particular, we analyze the impact of entropic regularization both on visual results and convergence speed.

## 1 Introduction

Generative modeling is at the heart of various problems in data science, either to approximate the data distribution in order to draw new samples, or to interpolate the data points. Beyond the purpose of image synthesis or editing, adopting such a generative model can also be used to reconstruct or restore corrupted data (Bora et al., 2017; Hand & Joshi, 2019; Hyder & Asif, 2020; Heckel & Soltanolkotabi, 2020; Menon et al., 2020; Shamshad & Ahmed, 2020; Damara et al., 2021; Leong, 2021) or to propose a geometric structure for the data that may reveal some interpretable dimensions (Radford et al., 2015; Shen et al., 2020). Supervised or not, learning generative models on large datasets thus opens new perspectives on the resolution of inverse problems.

Given the empirical distribution  $\nu$  of the data supported on a compact set  $\mathcal{Y} \subset \mathbb{R}^d$ , estimating a generative network consists in minimizing

$$\theta \mapsto \mathcal{L}(\mu_\theta, \nu) \tag{1}$$

where  $\mu_\theta$  is the distribution of generated samples (parameterized by a  $\theta$  in an open subset  $\Theta \subset \mathbb{R}^q$ ) with support included in a compact set  $\mathcal{X} \subset \mathbb{R}^d$ , and where  $\mathcal{L}$  is a loss function between probability distributions. The distribution  $\mu_\theta$  is often considered to be the law of a random variable  $g_\theta(Z)$  where  $g_\theta$  is a neural network and  $Z$  a random variable, and samples of the model can then be obtained by passing new realizations of  $Z$  through the network  $g_\theta$ . These distributions are often built upon features computed from samples and data points (such as the latent space of a variational auto-encoder (Kingma & Welling, 2014)) which may be integrated in the loss function (1). In this context, we face the long-standing problem of quantifying the discrepancy between probability distributions in a relevant and efficient manner.

## 1.1 Wasserstein Generative models

In the seminal work of Goodfellow *et al.* (Goodfellow et al., 2014) on adversarial training of generative networks, the considered loss is related (in a dual sense) to the Jensen-Shannon divergence between feature distributions. The major innovation of such a framework is that these features are simultaneously learnt from the dataset by training a binary classification network which competes against the generative network to discriminate between data point and generated samples. Arjovsky et al. (2017) later remarked that the Jensen-Shannon divergence has major flaws that directly impact the learning of GANs such as the convergence and robustness, and then proposed to use optimal transport (OT) costs instead, leading to new generative models called Wasserstein GANs (WGANs).

**Wasserstein distance** The OT cost between probability distributions  $\mu_\theta$  and  $\nu$  is defined by

$$W(\mu_\theta, \nu) = \inf_{\pi \in \Pi(\mu_\theta, \nu)} \int_{\mathcal{X} \times \mathcal{Y}} c(x, y) d\pi(x, y), \quad (2)$$

where  $\Pi(\mu_\theta, \nu)$  is the set of probability distributions on  $\mathcal{X} \times \mathcal{Y}$  having marginals  $\mu_\theta$  and  $\nu$ , while the ground cost  $c : \mathcal{X} \times \mathcal{Y} \rightarrow \mathbb{R}$  is a continuous function ( $c(x, y)$  represents the elementary cost between locations  $x \in \mathcal{X}$  and  $y \in \mathcal{Y}$ ). A simple and popular choice is the Euclidean distance between points to the power  $p \geq 1$ , *i.e.*  $c(x, y) = \|x - y\|^p = \left(\sum_{i=1}^d x_i^2\right)^{p/2}$ , for which  $W(\cdot, \cdot)^{\frac{1}{p}}$  defines the well-known  $p$ -Wasserstein distance. Another possible, but more complex choice, is to define the cost function as a metric in feature space.

As we will recall later, such an OT cost (2) admits a dual formulation (Santambrogio, 2015)

$$W(\mu_\theta, \nu) = \sup_{(\varphi, \psi) \in \mathcal{K}_c} \int \varphi d\mu_\theta + \int \psi d\nu, \quad (3)$$

where  $(\varphi, \psi)$  is a couple of dual variables that belongs to the set

$$\mathcal{K}_c = \{(\varphi, \psi) \in \mathcal{C}(\mathcal{X}) \times \mathcal{C}(\mathcal{Y}), \text{ subject to } \varphi(x) + \psi(y) \leq c(x, y) \quad \mu \otimes \nu \text{ a.e.} \} \quad (4)$$

where  $\mathcal{C}(\mathcal{X})$  indicates the set of real continuous functions on  $\mathcal{X}$  and  $\otimes$  the product of two measures. Optimizing one of the dual variables in (3) amounts to taking the  $c$ -transform

$$\psi^c(x) = \min_{y \in \mathcal{Y}} c(x, y) - \psi(y), \quad (5)$$

thus leading to another expression of the OT cost, often called “semi-dual formulation”:

$$W(\mu_\theta, \nu) = \sup_{\psi \in \mathcal{C}(\mathcal{Y})} \int \psi^c d\mu_\theta + \int \psi d\nu. \quad (6)$$

Solving the constrained dual problem (3) is a difficult (and possibly infinite-dimensional) optimization problem. One possibility to get an unconstrained optimization problem on  $(\varphi, \psi)$  is to rely on the entropic regularization of OT (Chizat, 2017; Peyré & Cuturi, 2019). The entropy-regularized OT admits a similar semi-dual formulation (6) except that a smoothed version of the minimum (called softmin) is used in the computation of the  $c$ -transform (5). In the discrete setting, the regularized OT cost can be computed efficiently with the Sinkhorn algorithm (Peyré & Cuturi, 2019), that exhibits geometric convergence. When one of the distribution is continuous (as in equation 1), one has to rely on stochastic algorithms (Genevay et al., 2016), which are provably convergent, but considerably slower. First attempts of using entropy-regularized OT for generative modeling have been made in (Genevay et al., 2018; Seguy et al., 2018) and we pursue here this investigation.

**Learning with Wasserstein losses** Once chosen the loss function  $\mathcal{L}$ , the minimization problem (1) can be solved with a gradient-based algorithm, for example a stochastic gradient descent or the ADAM algorithm. Hence the main topic of this paper is the computation of the gradients of (1) with respect to  $\theta$  in the case

where  $\mathcal{L}$  is an OT cost of the form (2), with or without entropic regularization. As we will prove later, the gradient of (1) is directly linked to the dual variable introduced in the dual formulation (3). Indeed, we will give conditions ensuring that the gradient at a point  $\theta_0$  can be expressed as

$$\nabla_{\theta}(W(\mu_{\theta}, \nu))|_{\theta=\theta_0} = \nabla_{\theta} \left( \int \varphi_* d\mu_{\theta} \right)|_{\theta=\theta_0} \quad (7)$$

where  $(\varphi_*, \psi_*)$  is an optimal dual variable for  $W(\mu_{\theta_0}, \nu)$ , *i.e.* a solution of the dual problem (3).

Such formula was proved in (Arjovsky et al., 2017) for the 1-Wasserstein cost, with the hypothesis that both sides of the equality exist. This proof was adapted by the authors of (Sanjabi et al., 2018) to the case of regularized Wasserstein costs. Both these proofs are based on some version of the so-called “envelope theorem” (also called Danskin’s theorem in the context of convex optimization), which allows to differentiate under the maximum. This theorem requires some regularity assumptions that should be carefully checked. As we will see in Section 2.5, in the discrete setting, there exist some irregular cases where these assumptions are not sufficient to make formula (7) licit.

To sum up, the main goal of this paper is to provide a new set of hypotheses that validates (7) and to show how these results apply to generative models parameterized by neural networks.

## 1.2 Related works

The computation of the gradient in equation 7 involves the optimal dual variable  $\varphi_*$ . When learning a Wasserstein generative model, the performance of the dual solver used for estimating the OT cost (3) is therefore a key point. Many attempts have already been made in the literature, with several ways to parameterize the problem (3).

The method of Arjovsky et al. (2017), being based on the 1-Wasserstein distance, only requires one dual variable that is constrained to be 1-Lipschitz. In practice, this dual variable is parameterized by a neural network, and the Lipschitz constraint is enforced by weight clipping (WGAN-WC). On a similar formulation, Gulrajani et al. (2017) suggest to impose Lipschitzness by including a gradient penalty in the dual loss (WGAN-GP).

In contrast, Seguy et al. (2018) consider regularized OT costs with a generic cost function. This leads to an unconstrained dual problem, but with two dual variables. In practice, the authors choose to parameterize both dual variables with neural networks. Closely related, the work by Sanjabi et al. (2018) is based on the same formulation of WGAN training with regularized OT and also relies on the neural network parameterization of the dual variables. The authors study the convergence and stability of the training procedure, the convergence being proved for the case of discrete distributions. In particular, they give the expression of the gradient of (7) (under a primal form) in the case of discrete regularized OT. This gradient expression is exploited to perform WGAN training by stochastic gradient descent. Notice also that Liu et al. (2019) applied a similar regularized OT framework on empirical distributions (*i.e.* discrete distributions obtained from samples) to learn a generator, with the particularity that the cost function depends on a set of simultaneously-learned features.

Closer to the proposed framework, Chen et al. (2019) consider the semi-dual formulation of OT (6). In their work, the dual variable  $\psi$  is optimized with the stochastic algorithm for semi-discrete OT proposed in (Genevay et al., 2016). Contrary to Seguy et al. (2018), they do not parameterize the dual variable  $\psi$  with a neural network, which hinders the applicability of their method to a very large dataset. Indeed, as shown in (Leclaire & Rabin, 2021), the convergence speed of the ASGD algorithm used for optimizing  $\psi$  decreases when either the dimension  $d$  of samples or the dimension of vector  $\psi$  (equal to the number of training points) increases. In (Chen et al., 2019), the corresponding primal solution  $\pi$  (which, in that case, is supported on a graph) is then used to perform a gradient step on the generative model using the estimated transportation cost. As a result, the whole algorithm of Chen et al. (2019) is not expressed as a direct minimization of (1). Our work, by contrast, shows that the proposed gradient calculation for (7) makes it possible to train the generator  $g_{\theta}$  and the dual variable  $\psi$  with an alternate min-max procedure on the dual cost (3).

---

Closely related to (Chen et al., 2019), Mallasto et al. (2019) propose to parameterize the dual variable  $\psi$  by a neural network and to obtain the second one  $\varphi$  with an approximated  $c$ -transform computed on mini-batches. With such an approximated  $c$ -transform, the pair of dual variables may not satisfy the constraint (4), which led the authors to integrate in the loss a penalty on  $c(x, y) - \varphi(x) - \psi(y)$ . A comparative study on the different ground costs is also realized. One benefit of this approach is that it scales up to a very large database, while keeping a relatively precise way to estimate the Wasserstein cost. On batch strategy, let us also mention the work by Fatras et al. (2020) who consider an alternative Wasserstein cost that is inherently defined as an expectation over mini-batches; this stochastic approximation introduces an estimation bias which is shown in practice to regularize the transportation problem.

### 1.3 Contributions and outline

The main contribution of this paper is to propose a complete set of hypotheses that ensure the validity of the gradient formula (7) for WGAN learning. Our approach is not restricted to the cost  $c(x, y) = \|x - y\|$  (inducing the 1-Wasserstein distance) and involves weak regularity hypotheses on the cost and the generator. Based on this gradient formula, we consider a stochastic algorithm for learning a generative model that can be understood as an alternate optimization algorithm on the semi-dual cost. On the practical side, we provide experiments on generative model learning for MNIST digits generation, which illustrate the impact of the entropic regularization both on numerical results and visual results.

In Section 2 we recall the complete framework for WGAN learning, and in particular we recall well-known results on the dual formulations of OT costs. We also give in Section 2.5 a counter-example based on discrete distributions where the formula (7) does not stand. In Section 3, we show the desired gradient formula in the semi-discrete setting, that is, when  $\nu$  is a distribution supported on a finite set. In Section 4, we adapt the proofs to the case of entropy-regularized OT (with no assumption on  $\nu$ ) and we also adapt the result to the Sinkhorn divergence, introduced in (Genevay et al., 2018) to remove the bias of the regularized OT cost. Section 5 draws a relation of the obtained formula with derivatives in the sense of distributions. Finally, Section 6 contains numerical experiments obtained with an alternate optimization algorithm for WGAN learning.

## 2 The Wasserstein GAN Problem

In this section, we first introduce the primal, dual and semi-dual formulations of the OT cost and we recall useful lemmas about  $c$ -transforms. OT is presented in the general case of entropic regularization with parameter  $\lambda \geq 0$ , thus encompassing the unregularized case  $\lambda = 0$ . Next, we introduce the generative learning problem as well as the regularity hypothesis on the generator, that will be used in the next sections to show the existence of gradients of the OT cost with respect to  $\theta$ . We close this section with a counter-example of measures  $\mu_\theta$  and  $\nu$  for which the desired gradient formula (7) does not hold.

### Notations

Let us here recall the notations used throughout the paper. Let  $\mathcal{X}, \mathcal{Y}$  be compact subsets of  $\mathbb{R}^d$ , and let  $c : \mathcal{X} \times \mathcal{Y} \rightarrow \mathbb{R}$  be a continuous function.

**Definition 1** ( $x$ -regularity). *We say that  $c$  is  **$x$ -regular** if there is  $L > 0$  and an open set  $U$  with  $\mathcal{X} \subset U \subset \mathbb{R}^d$  such that for any  $y \in \mathcal{Y}$ ,  $c(\cdot, y)$  can be extended to a  $L$ -Lipschitz  $\mathcal{C}^1$  function on  $U$ .*

We say that  $\varphi : \mathcal{X} \rightarrow \mathbb{R}$  is  $\mathcal{C}^1$  on  $\mathcal{X}$  if it is  $\mathcal{C}^1$  (i.e. differentiable with continuous derivatives) on an open neighborhood of  $\mathcal{X}$  (the neighborhood will often be the same  $U$  used for the assumption on the cost). Finally,  $\Theta$  is an open subset of  $\mathbb{R}^q$  used to parameterize the generator  $g_\theta$ .

## 2.1 Optimal Transport, Primal and Dual Problems

**Definition 2** (Primal formulation). *For  $\lambda \geq 0$ , the regularized optimal transport cost is defined by*

$$W_\lambda(\mu, \nu) = \inf_{\pi \in \Pi(\mu, \nu)} \int_{\mathcal{X} \times \mathcal{Y}} c \, d\pi + \lambda \text{KL}(\pi | \mu \otimes \nu) \quad (8)$$

where  $\Pi(\mu, \nu)$  is the set of probability distributions on  $\mathcal{X} \times \mathcal{Y}$  with marginals  $\mu$  and  $\nu$ , and where  $\text{KL}(\pi | \mu \otimes \nu) = \int \log(\frac{d\pi}{d\mu \otimes \nu}) d\pi$  if  $\pi$  admits a density  $\frac{d\pi}{d\mu \otimes \nu}$  w.r.t.  $\mu \otimes \nu$  and  $+\infty$  otherwise.

**Theorem 1** (Dual formulation (Santambrogio, 2015; Genevay, 2019; Feydy et al., 2019)). *Strong duality holds in the sense that*

$$W_\lambda(\mu, \nu) = \max_{\varphi, \psi} \int \varphi(x) d\mu(x) + \int \psi(y) d\nu(y) - \int m_\lambda(\varphi(x) + \psi(y) - c(x, y)) d\mu(x) d\nu(y) \quad (9)$$

where, for  $\lambda = 0$ ,  $m_0(t) = 0$  if  $t \geq 0$ , and  $+\infty$  otherwise, and for  $\lambda > 0$ ,  $m_\lambda(t) = \lambda(e^{\frac{t}{\lambda}} - 1)$ . A solution  $(\varphi, \psi)$  of this dual problem is called a pair of Kantorovich potentials. When  $\lambda > 0$ , the solutions of the dual problem are uniquely defined almost everywhere up to an additive constant (i.e. if  $(\varphi, \psi)$  is a solution, then any solution can be written  $(\varphi - k, \psi + k)$  with  $k \in \mathbb{R}$ ). Also, when  $\lambda > 0$ , the primal problem admits a unique solution

$$d\pi(x, y) = \exp\left(\frac{\varphi(x) + \psi(y) - c(x, y)}{\lambda}\right) d\mu(x) d\nu(y). \quad (10)$$

If  $\lambda = 0$  the primal solution  $\pi$  of (8) may not be absolutely continuous w.r.t.  $\mu \otimes \nu$  anymore. In that case, under weak assumptions (for example, when  $\mu, \nu$  admits second-order moments and  $\mu$  is absolutely continuous), one can show (Santambrogio, 2015) that the support of  $\pi$  is actually supported on the graph of an optimal transport map  $T^*$  i.e., the optimal  $\pi$  is the probability distribution of  $(X, T^*(X))$  where  $X$  has distribution  $\mu$  (and in particular  $T^*(X)$  has distribution  $\nu$ ).

For  $\psi \in \mathcal{C}(\mathcal{Y})$ , let us define the regularized  $c$ -transform as in (Feydy et al., 2019)

$$\psi^{c, \lambda}(x) = \text{softmin}_{y \in \mathcal{Y}} c(x, y) - \psi(y) \quad (11)$$

where the softmin operation is defined as

$$\text{softmin}_{y \in \mathcal{Y}} u(y) = \begin{cases} \min_{y \in \mathcal{Y}} u(y) & \text{if } \lambda = 0, \\ -\lambda \log \int e^{-\frac{u(y)}{\lambda}} d\nu(y) & \text{if } \lambda > 0. \end{cases} \quad (12)$$

We also define the analogous operators for the  $x$ -variable (and for simplicity, we use the same notation for  $c$ -transforms of  $x$ -functions or  $y$ -functions). It must be noted that the regularized  $c$ -transform  $\psi^{c, \lambda}$  also depends on  $\nu$  even if  $\nu$  is omitted in the notation.

Given a pair of dual variables  $(\varphi, \psi) \in \mathcal{C}(\mathcal{X}) \times \mathcal{C}(\mathcal{Y})$ , one can see that taking  $c$ -transforms

$$\begin{cases} \tilde{\psi} = \varphi^{c, \lambda} \\ \tilde{\varphi} = \psi^{c, \lambda} \end{cases}, \quad (13)$$

leads to a new pair  $(\tilde{\varphi}, \tilde{\psi})$  of dual variables that have a better dual cost (9) than  $(\varphi, \psi)$ . Therefore, the dual problem can always be restricted to  $c$ -concave functions, that is, functions that can be written as  $c$ -transforms. This is an important point since  $c$ -concave functions inherits some regularity from the cost function (see Lemma 1 below) and can be naturally extrapolated to any  $x \in \mathcal{X}$ .

**Theorem 2** (Semi-dual formulation (Genevay, 2019)). *The dual problem (9) is equivalent to the following semi-dual problem*

$$W_\lambda(\mu, \nu) = \sup_{\psi \in \mathcal{C}(\mathcal{Y})} \int \psi^{c, \lambda}(x) d\mu(x) + \int \psi(y) d\nu(y). \quad (14)$$

A solution  $\psi$  of the semi-dual problem is called a Kantorovich potential. In other words,  $\psi$  is a Kantorovich potential if and only if  $(\psi^{c,\lambda}, \psi)$  is a pair of Kantorovich potentials. By symmetry, we can also formulate a semi-dual problem on the dual variable  $\varphi$ .

Let us consider the exponential scalings of the dual variables  $a = e^{\frac{c}{\lambda}}, b = e^{\frac{\psi}{\lambda}}$ . With this notation, the coupled fixed point equations (13) can be reformulated as a single fixed point equation on  $a$  (or  $b$ ). The corresponding operator for  $a$  (or  $b$ ) can be shown to be a contractive operator on the unit sphere of  $L_+^\infty$  equipped with the Hilbert metric. In the discrete case where  $\mathcal{X}, \mathcal{Y}$  are both finite, iterating this contractive operator corresponds exactly to the Sinkhorn algorithm (Cuturi, 2013).

## 2.2 Continuity of $c$ -transforms

Let us now recall some well-known facts about regularized  $c$ -transforms. For that, we need a modulus of continuity of the cost function, that is, the smallest function  $\omega$  such that

$$\forall x, x' \in \mathcal{X}, \forall y, y' \in \mathcal{Y}, \quad |c(x, y) - c(x', y')| \leq \omega(\|x - x'\| + \|y - y'\|). \quad (15)$$

Since  $c$  is continuous on the compact  $\mathcal{X} \times \mathcal{Y}$ , it is uniformly continuous, thus  $\lim_{\delta \rightarrow 0} \omega(\delta) = 0$ .

**Lemma 1** ((Santambrogio, 2015; Feydy et al., 2019)). *For  $\lambda \geq 0$ , any  $c$ -transform  $\psi^{c,\lambda}$  has a modulus of continuity that is bounded by the modulus of continuity of the cost function.*

*Proof.* If  $u \leq v$  holds pointwise, then  $\text{softmin } u \leq \text{softmin } v$  pointwise. Also, for a constant  $k \in \mathbb{R}$ ,  $\text{softmin}(k + u) = k + \text{softmin}(u)$ . But, from the definition of  $\omega$ , we have

$$c(x, y) - \psi(y) \leq \omega(\|x - x'\|) + c(x', y) - \psi(y). \quad (16)$$

By taking the soft-min, we thus obtain

$$\psi^{c,\lambda}(x) \leq \omega(\|x - x'\|) + \psi^{c,\lambda}(x'), \quad (17)$$

and by symmetry, this leads to

$$|\psi^{c,\lambda}(x) - \psi^{c,\lambda}(x')| \leq \omega(\|x - x'\|). \quad (18)$$

□

**Lemma 2.** *For  $\lambda \geq 0$ , and any  $\psi, \chi \in \mathcal{C}(\mathcal{Y})$ ,  $\|\psi^{c,\lambda} - \chi^{c,\lambda}\|_\infty \leq \|\psi - \chi\|_\infty$ .*

*In other words, the map  $\psi \mapsto \psi^{c,\lambda}$  is 1-Lipschitz for the uniform norm.*

*Proof.* Applying again the monotonicity of the softmin operation to the inequality

$$c(x, y) - \psi(y) \leq \|\psi - \chi\|_\infty + c(x, y) - \chi(y) \quad (19)$$

we get  $\psi^{c,\lambda}(x) \leq \|\psi - \chi\|_\infty + \chi^{c,\lambda}(x)$ , which gives the desired result by a symmetry argument. □

The following lemma states that the optimal dual potentials vary continuously with respect to the input measure as soon as they are unique up to additive constants.

**Lemma 3** ((Feydy et al., 2019)). *Assume that  $\mathcal{X}, \mathcal{Y}$  are compact, and that  $c$  is continuous. Let us fix  $x_0 \in \mathcal{X}$ . Assume that  $\nu$  is fixed, and that  $\mu_n$  converges weak- $\star$  in  $\mathcal{M}_+^1(\mathcal{X})$  to a measure  $\mu$ . Assume furthermore that the Kantorovich potentials associated with  $W_\lambda(\mu, \nu)$  are unique up to an additive constant (which is always the case if  $\lambda > 0$ ). For each  $n$ , let  $\varphi_n$  be a  $c$ -concave Kantorovich potential for  $W_\lambda(\mu_n, \nu)$  such that  $\varphi_n(x_0) = 0$  and let  $\varphi$  be the (necessarily  $c$ -concave) Kantorovich potential for  $W_\lambda(\mu, \nu)$  such that  $\varphi(x_0) = 0$ .*

*Then  $\varphi_n \rightarrow \varphi$  uniformly on  $\mathcal{X}$ .*

*Proof.* By compactness of  $\mathcal{X} \times \mathcal{Y}$ ,  $c$  has a bounded modulus of continuity  $\omega(\delta)$  which tends to zero when  $\delta \rightarrow 0$ . By (18), we obtain that the functions  $|\varphi_n|$  are bounded by  $\sup_{x \in \mathcal{X}} \omega(\|x - x_0\|) < \infty$ . Besides, Lemma 1 also shows that the functions  $\varphi_n$  are uniformly equicontinuous on  $\mathcal{X}$ . Therefore, Arzela-Ascoli theorem ensures that the family  $\{\varphi_n, n \in \mathbb{N}\}$  is relatively compact in  $\mathcal{C}(\mathcal{X})$ .

Now, assume, by contradiction, that  $(\varphi_n)$  does not tend to  $\varphi$  in  $\mathcal{C}(\mathcal{X})$ . Then there would exist  $\varepsilon > 0$  and a subsequence  $(\varphi_{r(n)})$  such that

$$\|\varphi_{r(n)} - \varphi\|_\infty > \varepsilon \quad \forall n. \quad (20)$$

By relative compactness, one can then extract a subsequence  $(\varphi_{r(s(n))})$  which converges in  $\mathcal{C}(\mathcal{X})$  to a function  $\tilde{\varphi}$ . Using the monotonicity of soft-min, this implies that  $(\varphi_{r(s(n))}^{c,\lambda})$  also converges in  $\mathcal{C}(\mathcal{X})$  to  $\tilde{\varphi}^{c,\lambda}$ . Thus

$$W(\mu_n, \nu) = \int \varphi_n d\mu + \int \varphi_n^{c,\lambda} d\nu \xrightarrow{n \rightarrow \infty} \int \tilde{\varphi} d\mu + \int \tilde{\varphi}^{c,\lambda} d\nu. \quad (21)$$

Finally, we also have  $W(\mu_n, \nu) \rightarrow W(\mu, \nu)$  since  $\mu_n \xrightarrow{*} \mu$ . Thus

$$W(\mu, \nu) = \int \tilde{\varphi} d\mu + \int \tilde{\varphi}^{c,\lambda} d\nu \quad (22)$$

and  $\tilde{\varphi}$  is a Kantorovich potential for  $W(\mu, \nu)$  and  $\tilde{\varphi}(x_0) = \lim \varphi_n(x_0) = 0$ . Using the uniqueness assumption, we get  $\tilde{\varphi} = \varphi$  which contradicts (20).  $\square$

### 2.3 Learning a Generative Network

With the main OT concepts now defined, we can turn to the problem of learning a Wasserstein generative adversarial network. Estimating a WGAN from an empirical data distribution  $\nu$  consists in minimizing

$$h_\lambda(\theta) = W_\lambda(\mu_\theta, \nu), \quad (23)$$

where the generated distribution  $\mu_\theta$  is assumed to be the distribution of  $g_\theta(Z)$ , with  $Z$  a random variable. Denoting  $\zeta$  the probability distribution of  $Z$  on the measurable space  $\mathcal{Z}$ , we therefore have that  $\mu_\theta$  is the image measure of  $\zeta$  by the generator  $g_\theta$ , also known as the pushforward  $\mu_\theta = g_\theta \# \zeta$ . More precisely, the notation  $g_\theta$  refers to  $g(\theta, \cdot)$  where  $g : \Theta \times \mathcal{Z} \rightarrow \mathbb{R}^d$  is a function defined on the product of the open set  $\Theta \subset \mathbb{R}^q$  with  $\mathcal{Z}$ . In the following, we give different sets of conditions on  $g$  that allow to compute the derivatives of  $h_\lambda$ .

All the results of this paper are related to the behavior of the function

$$I(\varphi, \theta) = \int_{\mathcal{X}} \varphi d\mu_\theta, \quad (\varphi \in \mathcal{C}(\mathcal{X}), \theta \in \Theta). \quad (24)$$

Indeed, the semi-dual expression of optimal transport gives

$$h_\lambda(\theta) = \max_{\psi \in \mathcal{C}(\mathcal{Y})} I(\psi^{c,\lambda}, \theta) + \int_{\mathcal{Y}} \psi d\nu. \quad (25)$$

In order to study the gradient of  $h_\lambda$ , we thus define  $F_\lambda : \mathcal{C}(\mathcal{Y}) \times \Theta \rightarrow \mathbb{R}$  with

$$\forall \psi \in \mathcal{C}(\mathcal{Y}), \forall \theta \in \Theta, \quad F_\lambda(\psi, \theta) = I(\psi^{c,\lambda}, \theta) = \int_{\mathcal{X}} \psi^{c,\lambda} d\mu_\theta = \mathbb{E}[\psi^{c,\lambda}(g_\theta(Z))], \quad (26)$$

where the expectation is taken with respect to the probability distribution  $\zeta$  of  $Z$ .

Combining all previous definitions, the problem we tackle writes

$$W_\lambda(\mu_\theta, \nu) = h_\lambda(\theta) = \max_{\psi \in \mathcal{C}(\mathcal{Y})} H_\lambda(\psi, \theta) \quad (27)$$

with

$$H_\lambda(\psi, \theta) = F_\lambda(\psi, \theta) + \int_{\mathcal{Y}} \psi d\nu, \quad (28)$$



and our objective is to study the relationship between  $\nabla h_\lambda(\theta)$  and  $\nabla_\theta F_\lambda(\psi, \theta)$ . In the unregularized case, we use simpler notations,  $h$ ,  $W$ ,  $H$ , and  $F$  for  $h_0$ ,  $W_0$ ,  $H_0$  and  $F_0$  respectively, dropping the index  $\lambda$ .

As we will see later, computing the derivatives of  $h_\lambda$  boils down to differentiating under the max, which is allowed by the so-called envelope theorem. This result appears under different forms in the literature (Oyama & Takenawa, 2018). In Appendix A, we recall the version of the envelope theorem that is used in the proofs of the following sections.

However, if we temporarily admit the differentiability of all terms, the computation of the gradient is straightforward:

**Proposition 1.** *Let  $\theta_0$  and  $\psi_0$  satisfying  $h_\lambda(\theta_0) = H_\lambda(\psi_0, \theta_0)$ . If  $h_\lambda$  and  $\theta \mapsto F_\lambda(\psi_0, \theta)$  are both differentiable at  $\theta_0$ , then*

$$\nabla h_\lambda(\theta_0) = \nabla_\theta F_\lambda(\psi_0, \theta_0). \quad (\text{Grad-OT})$$

*Proof.* First, notice that  $H_\lambda(\psi_0, \cdot)$  and  $F_\lambda(\psi_0, \cdot)$  differ by a constant, and thus have same gradients. Using equation 25, for any  $\theta$ ,  $h_\lambda(\theta) \geq H_\lambda(\psi_0, \theta)$  with equality if  $\theta = \theta_0$ . Therefore, the function  $\theta \mapsto h_\lambda(\theta) - H_\lambda(\psi_0, \theta)$  has a minimum at  $\theta = \theta_0$ , and its gradient at  $\theta_0$  vanishes. This gives  $\nabla h_\lambda(\theta_0) = \nabla_\theta H_\lambda(\psi_0, \theta_0)$  and thus the desired result.  $\square$

Now, showing the existence of  $\nabla_\theta F_\lambda(\psi_0, \theta_0)$  consists in differentiating under the expectation in (26). For that, we need the following technical hypothesis.

**Definition 3** (Hypothesis  $(G_\Theta)$ ). *For  $\theta_0 \in \Theta$ , we say that  $g : \Theta \times \mathcal{Z} \rightarrow \mathcal{X}$  satisfies Hypothesis  $(G_{\theta_0})$  if there exists a neighborhood  $V$  of  $\theta_0$  and  $K \in L^1(\mathcal{Z})$  such that almost surely  $\theta \mapsto g(\theta, Z)$  is  $\mathcal{C}^1$  on  $V$  with differential  $\theta \mapsto D_\theta g(\theta, Z)$  and*

$$\forall \theta \in V, \quad \zeta\text{-a.s.} \quad \|g(\theta, Z) - g(\theta_0, Z)\| \leq K(Z)\|\theta - \theta_0\|. \quad (29)$$

*We say that  $g$  satisfies Hypothesis  $(G_\Theta)$  if  $g$  satisfies Hypothesis  $(G_{\theta_0})$  for any  $\theta_0 \in \Theta$ .*

**Remark 1.** *A sufficient condition for Hypothesis  $(G_\Theta)$  is that almost surely,  $\theta \mapsto g(\theta, Z)$  is  $\mathcal{C}^1$  on  $\Theta$  and that there exists  $K \in L^1(\mathcal{Z})$  such that*

$$\forall \theta, \theta' \in \Theta, \quad \zeta\text{-a.s.} \quad \|g(\theta', Z) - g(\theta, Z)\| \leq K(Z)\|\theta' - \theta\|. \quad (30)$$

*Notice that  $\Theta$  is an arbitrary open set, and can thus be reduced to localize the problem in a neighborhood of a fixed point  $\theta \in \Theta$ . The interest of Definition 3 is that it does not impose uniformity in  $\theta_0$  (for both the neighborhood  $V$  and the upper bound  $K(Z)$ ).*

## 2.4 Previous results on the differentiability of OT costs

Outside the recent context of WGAN learning, the regularity and gradient of the entropy-regularized OT cost  $p \mapsto W_\lambda(p, q)$  (with  $\lambda > 0$ ) were expressed in (Cuturi & Peyré, 2016, Prop. 2.3) in the case of discrete distributions  $(p, q)$ . In the same context, the Gâteaux-differentiability of  $(p, q) \mapsto W_\lambda(p, q)$  was proved in (Feydy et al., 2019), which was later extended to continuous differentiability in (Bigot et al., 2019, Prop. 2.3). If  $(p_\theta)$  is a parametric family of distribution supported on the *same* finite set, applying the chain rule gives the differentiability of  $\theta \mapsto W_\lambda(p_\theta, q)$ .

The gradient formula (Grad-OT) is given in the seminal paper on WGAN (Arjovsky et al., 2017) with the assumption that both sides of the equality exist. This formula (extended to more general costs) has been exploited in several papers related to WGAN learning, for example (Liu et al., 2019). In the context of discrete regularized OT, one can find in (Sanjabi et al., 2018, Appendix C) a gradient formula expressed through the primal formulation, with a short proof limited to discrete regularized OT.

## 2.5 A telling counter-example

We now show that the differentiation with respect to  $\theta$  under the max in (25) can fail, even for regular generators  $g$ , thus discarding the formula (Grad-OT). To illustrate this point, let us consider a simple

unregularized OT problem (i.e.  $\lambda = 0$ ) between a Dirac  $\delta_\theta$  located at  $\theta \in \mathbb{R}^d$  and a sum of two Diracs at positions  $y_1 \neq y_2 \in \mathbb{R}^d$ . This setting can be obtained by setting  $Z = 0$  almost surely and defining  $g_\theta(z) = \theta - z$ .

**Proposition 2.** *Let  $\mu_\theta = \delta_\theta$  and  $\nu = \frac{1}{2}\delta_{y_1} + \frac{1}{2}\delta_{y_2}$ .*

*For  $p \geq 1$ , consider the cost  $c(x, y) = \|x - y\|^p$  where  $\|\cdot\|$  is the Euclidean norm on  $\mathbb{R}^d$ .*

*Then*

- $h(\theta) = W(\mu_\theta, \nu)$  is differentiable at any  $\theta \notin \{y_1, y_2\}$  for  $p = 1$ , and at any  $\theta$  for  $p > 1$ ,
- for any  $\psi_{*0} \in \arg \max_\psi H(\psi, \theta_0)$ , the function  $\theta \mapsto F(\psi_{*0}, \theta)$  is **not** differentiable at  $\theta_0$ .

*Hence relation (Grad-OT) never stands.*

*Proof.* The only distribution on  $\mathcal{X} \times \mathcal{Y}$  having marginals  $\mu_\theta, \nu$  is  $\pi = \frac{1}{2}\delta_{(\theta, y_1)} + \frac{1}{2}\delta_{(\theta, y_2)}$ . Therefore, recalling that  $h(\theta) = W(\mu_\theta, \nu)$  and from the definition of the primal problem (8), we have

$$h(\theta) = \frac{1}{2}c(\theta, y_1) + \frac{1}{2}c(\theta, y_2) \quad (31)$$

which gives the first point.

Next, one can explicitly solve the dual problem, which, from the equivalent semi-dual formulation (14) and the definition (26) of  $F$ , reduces to the following optimization problem with respect to  $(\psi(y_1), \psi(y_2)) \in \mathbb{R}^2$

$$\max_{\psi \in \mathbb{R}^2} H(\psi, \theta) \quad \text{with} \quad H(\psi, \theta) = \psi^c(\theta) + \frac{\psi(y_1) + \psi(y_2)}{2}, \quad (32)$$

where, for any  $\psi$ , the  $c$ -transform (12) writes

$$\psi^c(\theta) = \begin{cases} c(\theta, y_1) - \psi(y_1) & \text{if } c(\theta, y_1) - \psi(y_1) \leq c(\theta, y_2) - \psi(y_2), \\ c(\theta, y_2) - \psi(y_2) & \text{otherwise.} \end{cases} \quad (33)$$

Therefore,

$$H(\psi, \theta) = \begin{cases} c(\theta, y_1) + \frac{1}{2}(\psi(y_2) - \psi(y_1)) & \text{if } \psi(y_2) - \psi(y_1) \leq c(\theta, y_2) - c(\theta, y_1), \\ c(\theta, y_2) - \frac{1}{2}(\psi(y_2) - \psi(y_1)) & \text{otherwise.} \end{cases} \quad (34)$$

For a fixed  $\theta_0$ ,  $H(\cdot, \theta_0)$  is maximal at any  $\psi_{*0}$  such that

$$\psi_{*0}(y_2) - \psi_{*0}(y_1) = c(\theta_0, y_2) - c(\theta_0, y_1). \quad (35)$$

Besides, from (34), one sees that  $H(\psi, \cdot)$  is made of two pieces whose gradients can be computed explicitly for the cost  $c(x, y) = \|x - y\|^p$ :

$$\forall \theta \neq y_j, \quad \nabla_\theta c(\theta, y_j) = p\|\theta - y_j\|^{p-2}(\theta - y_j), \quad (36)$$

and, when  $p > 1$  we also have  $\nabla_\theta c(\theta, y_j) = 0$  for  $\theta = y_j$ . Therefore, the gradients of the two pieces agree at no point  $\theta$ : for any  $\theta$ ,  $\nabla_\theta c(\theta, y_1) \neq \nabla_\theta c(\theta, y_2)$ . In particular, the function  $H(\psi, \cdot)$  is not differentiable at the interface  $\mathcal{I} = \{ \theta \in \mathbb{R}^d \mid c(\theta, y_2) - c(\theta, y_1) = \psi(y_2) - \psi(y_1) \}$ . As  $F(\psi, \cdot)$  only differs from  $H(\psi, \cdot)$  by the constant  $\frac{\psi(y_1) + \psi(y_2)}{2}$ , it is also not differentiable on  $\mathcal{I}$ . But then, for  $\psi_{*0}$  satisfying (35),  $\theta_0$  lies on the interface  $\mathcal{I}$ , which gives the second point.  $\square$

We highlight the main pitfall in the previous proof:  $\delta_\theta$  puts some mass on a thin subset where  $\theta \mapsto F(\psi_{*0}, \theta)$  is not smooth. In the following section, we adopt an hypothesis on the generator that avoids this objection.

### 3 Gradient formula in the unregularized semi-discrete setting

In this section, we prove the formula (Grad-OT) in the semi-discrete case, *i.e.* when the target measure  $\nu$  is supported on a finite set of points. Therefore, in this section,  $\mathcal{X} \subset \mathbb{R}^d$  is compact and  $\mathcal{Y}$  is finite with  $J$  points, so that  $\mathcal{C}(\mathcal{Y})$  identifies to  $\mathbb{R}^J$ . We only consider the case  $\lambda = 0$  since more general results are given for regularized OT in the next section. We recall the notations  $W = W_0$ ,  $h = h_0$ , and  $F = F_0$  obtained from the definitions (8), (23) and (26) with  $\lambda = 0$ . Also, for  $\lambda = 0$ , we simply write  $\psi^c = \psi^{c,0}$  the  $c$ -transform of a function  $\psi$ .

Since  $\mathcal{C}(\mathcal{Y})$  identifies to  $\mathbb{R}^J$ , we study the function  $F : \mathbb{R}^J \times \Theta \rightarrow \mathbb{R}$  defined by

$$F(\psi, \theta) = \int_{\mathcal{X}} \psi^c d\mu_{\theta} = \mathbb{E}[\psi^c(g_{\theta}(Z))]. \quad (37)$$

In the following results, we provide an hypothesis on the generator  $g_{\theta}$  that allows to compute the gradient of  $h$  by differentiating under the max in (27). This hypothesis requires the definition of the open Laguerre cells associated with  $\psi$ :

$$\mathbf{L}_{\psi}(y) = \{ x \in \mathcal{X} \mid \forall y' \neq y, c(x, y) - \psi(y) < c(x, y') - \psi(y') \}. \quad (38)$$

It also requires the Laguerre cells boundaries, that is, the set of points for which there is a “tie” in the minimum:

$$A_{\psi} = \mathcal{X} \setminus \bigcup_{y \in \mathcal{Y}} \mathbf{L}_{\psi}(y). \quad (39)$$

For example, if  $c(x, y) = \|x - y\|_2^2$  in  $\mathbb{R}^d$ ,  $A_{\psi}$  is contained in a union of hyperplanes whose directions are orthogonal to the segments  $[y_1, y_2]$ ,  $y_1, y_2 \in \mathcal{Y}$ . In particular  $A_{\psi}$  has zero Lebesgue measure.

By construction, for  $x \in \mathcal{X} \setminus A_{\psi}$ , we can uniquely define the Monge map

$$T_{\psi}(x) = \arg \min_{y \in \mathcal{Y}} c(x, y) - \psi(y). \quad (40)$$

**Lemma 4.** *The map  $(\psi, x) \mapsto T_{\psi}(x)$  is locally constant on  $\mathbb{R}^J \times A_{\psi}$ .*

*Proof.* Let  $(\psi, x) \in \mathbb{R}^J \times A_{\psi}$  and let  $y = T_{\psi}(x)$ . By definition of Laguerre cells,

$$c(x, y) - \psi(y) - \min_{z \in \mathcal{Y} \setminus \{y\}} (c(x, z) - \psi(z)) < 0. \quad (41)$$

The left-hand side is a function that is jointly continuous in  $(\psi, x)$  and is  $< 0$  at  $(\psi, x)$ . Thus, there is a neighborhood  $W$  of  $(\psi, x)$  where it stays negative, that is,

$$\forall (\psi', x') \in W, \quad c(x', y) - \psi'(y) - \min_{z \in \mathcal{Y} \setminus \{y\}} (c(x', z) - \psi'(z)) < 0. \quad (42)$$

Therefore  $T_{\psi'}(x') = y$  on  $W$ . □

**Lemma 5.** *Assume that  $\mathcal{Y}$  is finite with  $J$  points and that  $c$  is  $x$ -regular (see Definition 1). Let  $\psi \in \mathbb{R}^J$ . Then  $\psi^c$  is  $\mathcal{C}^1$  on  $\mathcal{X} \setminus A_{\psi}$  and*

$$\forall x \in \mathcal{X} \setminus A_{\psi}, \quad \nabla \psi^c(x) = \nabla_x c(x, T_{\psi}(x)). \quad (43)$$

*Proof.* First, one can notice that  $\mathcal{X} \setminus A_{\psi} = \bigcup_{y \in \mathcal{Y}} \mathbf{L}_{\psi}(y)$  is an open subset of  $\mathcal{X}$  because the Laguerre cells  $\mathbf{L}_{\psi}(y)$  are open (thanks to the continuity of  $c$ ). Besides, if  $x \in \mathcal{X} \setminus A_{\psi}$ ,  $y = T_{\psi}(x)$  is well-defined and  $x \in \mathbf{L}_{\psi}(y)$ . Thus  $\mathbf{L}_{\psi}(y)$  is an open neighborhood of  $x$  on which we have

$$\forall u \in \mathbf{L}_{\psi}(y), \quad \psi^c(u) = c(u, y) - \psi(y).$$

Therefore, on  $\mathbf{L}_{\psi}(y)$ ,  $\psi^c$  is as regular as  $c(\cdot, y) = c(\cdot, T_{\psi}(x))$  and has the same gradient. □

**Lemma 6.** Assume that  $\mathcal{Y}$  is finite with  $J$  points and that  $c$  is  $x$ -regular. Assume that  $g$  satisfies Hypothesis  $(\mathbf{G}_{\theta_0})$  (Definition 3) at some  $\theta_0 \in \Theta$  and  $V$  the associated neighborhood of  $\theta_0$ . Assume also that for any  $\psi \in \mathbb{R}^J$ ,  $\mu_\theta(A_\psi) = 0$ , i.e. we have  $g(\theta, Z) \in \mathcal{X} \setminus A_\psi$  almost surely.

Then, for any  $\psi \in \mathbb{R}^J$ ,  $\theta \mapsto F(\psi, \theta)$  is differentiable at  $\theta_0$  and

$$\nabla_\theta F(\psi, \theta_0) = \mathbb{E}[D_\theta g(\theta_0, Z)^T \nabla \psi^c(g(\theta_0, Z))] \quad (44)$$

where  $D_\theta g$  is the partial differential of  $g$  with respect to  $\theta$ .

*Proof.* For  $\theta \in \Theta$ , let us denote

$$f(\psi, \theta, Z) = \psi^c(g(\theta, Z)) \quad (45)$$

so that we get  $F(\psi, \theta) = \mathbb{E}[f(\psi, \theta, Z)]$  from (37). The hypotheses on  $g$  and Lemma 5 ensure that  $f(\psi, \cdot, Z)$  is almost surely differentiable at  $\theta_0$  and thanks to the chain-rule, we have

$$\nabla_\theta f(\psi, \theta_0, Z) := D_\theta g(\theta_0, Z)^T \nabla \psi^c(g(\theta_0, Z)). \quad (46)$$

Besides, since  $c$  is  $x$ -regular, the  $c$ -transforms are  $L$ -Lipschitz thanks to Definition 1 and Lemma 1. Therefore, for any  $\theta \in V$ ,

$$\|f(\psi, \theta, Z) - f(\psi, \theta_0, Z)\| \leq L\|g(\theta, Z) - g(\theta_0, Z)\| \leq LK(Z)\|\theta - \theta_0\| \quad \text{a.s.} \quad (47)$$

Besides, replacing  $\theta$  by  $\theta_0 + t(\theta - \theta_0)$  for  $t \in \mathbb{R}$  and letting  $t \rightarrow 0$ , we also get

$$\|\nabla_\theta f(\psi, \theta_0, Z) \cdot (\theta - \theta_0)\| \leq LK(Z)\|\theta - \theta_0\| \quad \text{a.s.} \quad (48)$$

In particular,  $\mathbb{E}[\nabla_\theta f(\psi, \theta_0, Z)]$  exists. Therefore, we have for any  $\theta \in V$ ,

$$\frac{\|f(\psi, \theta, Z) - f(\psi, \theta_0, Z) - \nabla_\theta f(\psi, \theta_0, Z) \cdot (\theta - \theta_0)\|}{\|\theta - \theta_0\|} \leq 2LK(Z) \quad \text{a.s.} \quad (49)$$

When  $\theta \rightarrow \theta_0$ , the left-hand-side tends almost surely to zero, and thus, the dominated convergence theorem ensures that

$$\mathbb{E} \left[ \frac{\|f(\psi, \theta, Z) - f(\psi, \theta_0, Z) - \nabla_\theta f(\psi, \theta_0, Z) \cdot (\theta - \theta_0)\|}{\|\theta - \theta_0\|} \right] \rightarrow 0 \quad (50)$$

and in particular,

$$F(\psi, \theta) - F(\psi, \theta_0) - \mathbb{E}[\nabla_\theta f(\psi, \theta_0, Z)] \cdot (\theta - \theta_0) = o(\|\theta - \theta_0\|). \quad (51)$$

This proves that  $F(\psi, \cdot)$  is differentiable at  $\theta_0$  with

$$\nabla_\theta F(\psi, \theta_0) = \mathbb{E}[\nabla_\theta f(\psi, \theta_0, Z)] = \mathbb{E}[D_\theta g(\theta_0, Z)^T \nabla \psi^c(g_{\theta_0}(Z))]. \quad (52)$$

□

**Theorem 3.** Assume that  $\mathcal{Y}$  is finite with  $J$  points and that  $c$  is  $x$ -regular. Assume that

1. for any  $\theta \in \Theta$ , the Kantorovich potential  $\psi_*$  for  $W(\mu_\theta, \nu)$  (defined in Theorem 2) is unique up to additive constants.
2. for any  $\theta \in \Theta$  and any  $\psi \in \mathbb{R}^J$ ,  $\mu_\theta(A_\psi) = 0$ , i.e. we have  $g(\theta, Z) \in \mathcal{X} \setminus A_\psi$  a.s.
3.  $g$  satisfies Hypothesis  $(\mathbf{G}_\Theta)$  in Definition 3.

Then  $h(\theta) = W(\mu_\theta, \nu)$  is differentiable at any  $\theta \in \Theta$  and

$$\nabla h(\theta) = \nabla_\theta F(\psi_*, \theta) = \mathbb{E} [D_\theta g(\theta, Z)^T \nabla \psi_*^c(g_\theta(Z))] . \quad (53)$$

where all terms are well-defined.

*Proof.* The proof consists in applying the envelope Theorem 7 (recalled in Appendix A) in order to show that  $h$  is differentiable at a fixed  $\theta_0 \in \Theta$ . Let us denote by  $\psi_{*\theta_0}$  the corresponding Kantorovich potential for  $W(\mu_{\theta_0}, \nu)$ .

First, we need to build a selection of Kantorovich potentials that is continuous at  $\theta_0$ . Since we assumed that for any  $\theta \in V$ , the Kantorovich potential for  $W(\mu_\theta, \nu)$  is unique up to additive constants, there is a unique  $\psi_{*\theta} \in \arg \max F(\cdot, \theta)$  such that  $\psi_{*\theta}(y_1) = 0$  (where  $y_1$  is an arbitrary point in  $\mathcal{Y}$ ). The continuity of  $\theta \mapsto \psi_{*\theta}$  then follows from Lemma 3 and the fact that  $\theta \mapsto \mu_\theta$  is weak- $\star$  continuous at  $\theta_0$ . Indeed, Hypothesis 3 implies that when  $\theta \rightarrow \theta_0$ ,  $\mathbb{E}[\|g_\theta(Z) - g_{\theta_0}(Z)\|] \rightarrow 0$ , which means that  $g_\theta(Z) \rightarrow g_{\theta_0}(Z)$  in  $L^1(\mathcal{Z}, \mathbb{R}^d)$  and thus in distribution.

The next step is to show that  $\nabla_\theta F(\psi, \theta)$  exists in the neighborhood of  $(\psi_{*\theta_0}, \theta_0)$ , which is guaranteed by Lemma 6

Finally, we have to demonstrate that  $(\psi, \theta) \mapsto \nabla_\theta F(\psi, \theta)$  is continuous at  $(\psi_{*\theta_0}, \theta_0)$ . For that, we show that the function  $f(\psi, \theta, Z)$  introduced in Equation (45) has a simple expression in the neighborhood of  $(\psi_{*\theta_0}, \theta_0)$ . Indeed, for  $\zeta$ -almost all  $z$ , we can define  $y = T_{\psi_{*\theta_0}}(g_{\theta_0}(z))$ . Lemma 4 gives a neighborhood  $W_1 \times W_2$  of  $(\psi_{*\theta_0}, g_{\theta_0}(z))$  where  $T_\psi(x) = y$ . By continuity of  $g(\cdot, z)$ ,  $U = W_1 \times (g(\cdot, z)^{-1}(W_2))$  is a neighborhood of  $(\psi_{*\theta_0}, \theta_0)$  such that  $\forall (\psi, \theta) \in U$ ,  $T_\psi(g_\theta(z)) = y$ . Thus,

$$\forall (\psi, \theta) \in U, \quad \nabla_\theta f(\psi, \theta, z) = D_\theta g(\theta, z)^T \nabla \psi^c(g_\theta(z)) = D_\theta g(\theta, z)^T \nabla_x c(g_\theta(z), y). \quad (54)$$

In the neighborhood  $U$ , the right-hand side does not depend on  $\psi$  anymore. It is also continuous in  $\theta$  thanks to Hypothesis (G $_\Theta$ ) and the fact that  $c$  is  $x$ -regular. This proves that almost surely  $\nabla_\theta f(\cdot, \cdot, Z)$  is continuous at  $(\psi_{*\theta_0}, \theta_0)$ . Finally, Equation (48) proves that all components of  $\nabla_\theta f(\cdot, \cdot, Z)$  are almost surely bounded by  $LK(Z) \in L^1(\mathcal{Z})$ . Therefore, the dominated convergence theorem ensures that  $\nabla F(\psi, \theta) = \mathbb{E}[\nabla_\theta f(\psi, \theta, Z)]$  is continuous at  $(\psi_{*\theta_0}, \theta_0)$ .

We can thus apply the envelope Theorem 7 which gives the desired result.  $\square$

The hypothesis that  $c$  is  $x$ -regular (see Definition 1) is a quite strong constraint which may not be satisfied in practice even for very usual costs, like the cost  $c(x, y) = \|x - y\|$  in  $\mathbb{R}^d$  for the 1-Wasserstein distance. We now give a similar result with a relaxed hypothesis that encompasses such non-smooth cost functions.

**Theorem 4.** *Assume that  $\mathcal{Y}$  is finite with  $J$  points. Assume that  $c : \mathcal{X} \times \mathcal{Y} \rightarrow \mathbb{R}$  is continuous, and that there is a constant  $L > 0$ , an open set  $U$  with  $\mathcal{X} \subset U \subset \mathbb{R}^d$  and a set  $B \subset \mathcal{X}$  closed in  $\mathbb{R}^d$  such that for any  $y \in \mathcal{Y}$ ,  $c(\cdot, y)$  can be extended to a  $L$ -Lipschitz  $\mathcal{C}^1$  function on  $U \setminus B$ . Assume that  $g$  satisfies the three assumptions of Theorem 3, and assume also that  $\mu_\theta(B) = 0$  for any  $\theta \in \Theta$ .*

*Then  $h(\theta) = W(\mu_\theta, \nu)$  is differentiable at any  $\theta \in \Theta$  with the gradient expression (53).*

Notice that in the case of the Euclidean distance  $c(x, y) = \|x - y\|$  in  $\mathbb{R}^d$ , for any  $y \in \mathcal{Y}$ ,  $c(\cdot, y)$  is  $\mathcal{C}^1$  only on  $\mathbb{R}^d \setminus \{y\}$ . Thus this cost satisfies the condition of the last theorem with  $B = \mathcal{Y}$ , and the resulting hypothesis on the generator reads as  $\mu_\theta(\mathcal{Y}) = 0$  (in addition to the fact that all boundaries  $A_\psi$  of Laguerre cells are also  $\mu_\theta$ -negligible).

*Proof.* One may check that the main parts of the proof given for Theorem 3 are still true with the relaxed hypothesis on the cost. We thus only highlight the minor modifications. First, in Lemma 5, one obtains that the  $c$ -transforms  $\psi^c$  are only  $\mathcal{C}^1$  on  $\mathcal{X} \setminus (A_\psi \cup B)$ . The proof of Lemma 6 is unchanged because  $g_\theta(Z)$  almost surely belongs to  $\mathcal{X} \setminus (A_\psi \cup B)$  where the  $c$ -transforms are differentiable. Finally, in the last step of the proof of Theorem 3, one should appropriately adapt the neighborhood  $W_1 \times W_2$ , which is simply done by replacing  $W_2$  by  $W_2 \cap B^c$ .  $\square$

**Remark 2.** *There exist sufficient conditions that ensure the uniqueness of Kantorovich potentials up to additive constants. Indeed, according to (Santambrogio, 2015, Prop. 7.18), the Kantorovich potential for  $W(\mu_\theta, \nu)$  is unique up to additive constants as soon as the support of  $\mu_\theta$  is the closure of a bounded connected open set. Assuming that  $\mu_\theta = g_\theta \# \zeta$  with  $\zeta$  being the uniform distribution on the hypercube  $Q = [-1, 1]^s$  and  $g_\theta : Q \rightarrow \mathbb{R}^d$  any continuous map, the support of  $\mu_\theta$  is exactly  $g_\theta(Q)$  (see Proposition 3). If moreover the*

image  $g_\theta(\mathring{Q})$  of the interior of  $Q$  is assumed to be open, then  $g_\theta(Q)$  is the closure of  $g_\theta(\mathring{Q})$  which is connected, and thus the uniqueness of Kantorovich potentials follows. Let us mention however that in the case where  $g_\theta$  is given by a neural network,  $g_\theta(\mathring{Q})$  is likely not to be open (for example in the expected case where the image of  $g_\theta$  is included in an hyperplane or a manifold of dimension  $< d$ ). It may also be that the uniqueness of Kantorovich potentials is not ensured for every  $\theta \in \Theta$ . Then, if one is only interested in the local behavior of  $\theta \mapsto W(\mu_\theta, \nu)$  around a given  $\theta_0 \in \Theta$ , one may restrict  $\Theta$  to be an open neighborhood of  $\theta_0$ . Another way to ensure uniqueness of Kantorovich potentials is to work with entropic optimal transport, as we do in the next section.

## 4 Gradient formula for regularized optimal transport

In this section, we provide differentiability results for  $h_\lambda$  in the case of regularized optimal transport. In this setting, as soon as the cost is regular, the regularized  $c$ -transforms are also regular everywhere, while requiring no assumption of the target measure  $\nu$ , thus not restricted to the semi-discrete case. This makes the situation simpler than the unregularized case. Again, we recall the notations  $W_\lambda$ ,  $h_\lambda$ , and  $F_\lambda$  from (8), (23) and (26). In order to obtain the gradient of  $h_\lambda$ , we follow a similar strategy than in Section 3 and study the relation between  $\nabla h_\lambda(\theta)$  and  $\nabla_\theta F_\lambda(\psi, \theta)$ .

### 4.1 Gradient of the regularized $c$ -transforms

**Lemma 7.** *Let  $\lambda > 0$ . Assume that  $c$  is  $x$ -regular (see Definition 1). Let  $\psi \in \mathcal{C}(\mathcal{Y})$ .*

*Then  $\psi^{c,\lambda}$  is  $\mathcal{C}^1$  on  $\mathcal{X}$  and*

$$\forall x \in \mathcal{X}, \quad \nabla \psi^{c,\lambda}(x) = \int_{\mathcal{Y}} \exp\left(\frac{\psi^{c,\lambda}(x) + \psi(y) - c(x, y)}{\lambda}\right) \nabla_x c(x, y) d\nu(y) \quad (55)$$

$$= \frac{\int_{\mathcal{Y}} \exp\left(\frac{\psi(y) - c(x, y)}{\lambda}\right) \nabla_x c(x, y) d\nu(y)}{\int_{\mathcal{Y}} \exp\left(\frac{\psi(y) - c(x, y)}{\lambda}\right) d\nu(y)}. \quad (56)$$

*Proof.* By definition of the regularized  $c$ -transform (see relations (11) and (12)), we have for any  $x \in \mathcal{X}$ ,

$$\forall x \in \mathcal{X}, \quad e^{-\frac{\psi^{c,\lambda}(x)}{\lambda}} = \int_{\mathcal{Y}} e^{\frac{\psi(y) - c(x, y)}{\lambda}} d\nu(y). \quad (57)$$

Since  $c$  is  $x$ -regular, then for any  $y \in \mathcal{Y}$ ,  $x \mapsto e^{\frac{\psi(y) - c(x, y)}{\lambda}}$  is  $\mathcal{C}^1$  on a neighborhood of  $\mathcal{X}$ , and it is Lipschitz with Lipschitz constant  $\frac{L}{\lambda} \exp(\frac{\|\psi\|_\infty + \|c\|_\infty}{\lambda})$ . Therefore, we can differentiate under the integral to get

$$\forall x \in \mathcal{X}, \quad \nabla_x \left( e^{-\frac{\psi^{c,\lambda}(x)}{\lambda}} \right) = -\frac{1}{\lambda} \int_{\mathcal{Y}} \nabla_x c(x, y) e^{\frac{\psi(y) - c(x, y)}{\lambda}} d\nu(y). \quad (58)$$

Expanding the left-hand side, this is equivalent to

$$\forall x \in \mathcal{X}, \quad \nabla \psi^{c,\lambda}(x) e^{-\frac{\psi^{c,\lambda}(x)}{\lambda}} = \int_{\mathcal{Y}} \nabla_x c(x, y) e^{\frac{\psi(y) - c(x, y)}{\lambda}} d\nu(y), \quad (59)$$

which gives the desired formula (55) for  $\nabla \psi^{c,\lambda}$ . The second expression (56) follows from using the definition of  $\psi^{c,\lambda}$  again. Finally, using (56), one can see that all integrated functions are continuous with respect to  $x$ , and bounded. The dominated convergence theorem thus ensures that  $\nabla \psi^{c,\lambda}$  is continuous.  $\square$

### 4.2 Regularity of $F_\lambda$

We recall that Hypothesis  $(G_\Theta)$  is given in Definition 3.

**Lemma 8.** *Let  $\lambda > 0$ . Assume that  $c$  is  $x$ -regular and that  $g$  satisfies Hypothesis  $(G_\Theta)$ .*

*Let  $\psi \in \mathcal{C}(\mathcal{Y})$ . Then the function  $\theta \mapsto F_\lambda(\psi, \theta)$  is differentiable on  $\Theta$  and*

$$\forall \theta_0 \in \Theta, \quad \nabla_\theta F_\lambda(\psi, \theta_0) = \mathbb{E} [D_\theta g(\theta_0, Z)^T \nabla \psi^{c, \lambda}(g(\theta_0, Z))] . \quad (60)$$

*Proof.* As in the proof of Lemma 6, let us introduce

$$f_\lambda(\psi, \theta, Z) = \psi^{c, \lambda}(g_\theta(Z)) \quad (61)$$

so that  $F_\lambda(\psi, \theta) = \mathbb{E}[f_\lambda(\psi, \theta, Z)]$ . Using Hypothesis  $(G_\Theta)$  and Lemma 7,  $f_\lambda(\psi, \cdot, Z)$  is differentiable on  $\Theta$  almost surely, and, thanks to the chain-rule,

$$\nabla_\theta f_\lambda(\psi, \theta, Z) = D_\theta g(\theta, Z)^T \nabla \psi^{c, \lambda}(g_\theta(Z)). \quad (62)$$

Besides, the regularized  $c$ -transforms of a  $x$ -regular cost being still  $L$ -Lipschitz, we have an integrable bound for the finite differences of  $f_\lambda$ . Indeed, for a fixed  $\theta_0 \in \Theta$ , and for any  $\theta$  in the neighborhood  $V$  of  $\theta_0$  given by Hypothesis  $(G_{\theta_0})$ ,

$$\|f_\lambda(\psi, \theta, Z) - f_\lambda(\psi, \theta_0, Z)\| \leq L \|g(\theta, Z) - g(\theta_0, Z)\| \leq LK(Z) \|\theta - \theta_0\| \quad \text{a.s.} \quad (63)$$

The proof can then be ended exactly as the one of Lemma 6. For further use, notice that the previous bound implies

$$\|\nabla_\theta f_\lambda(\psi, \theta, Z)\| = \|D_\theta g(\theta_0, Z)^T \nabla \psi^{c, \lambda}(g(\theta_0, Z))\| \leq LK(Z) \quad \text{a.s.} \quad (64)$$

where  $\|\cdot\|$  is the dual norm associated with  $\|\cdot\|$ .  $\square$

### 4.3 Gradient of the regularized loss

**Theorem 5.** *Assume that  $c$  is  $x$ -regular and that  $g$  satisfies Hypothesis  $(G_\Theta)$ .*

*Then  $h_\lambda$  is  $\mathcal{C}^1$  on  $\Theta$  and*

$$\forall \theta \in \Theta, \quad \nabla_\theta h_\lambda(\theta) = \nabla_\theta F_\lambda(\psi_*, \theta) = \mathbb{E} [D_\theta g(\theta, Z)^T \nabla \psi_*^{c, \lambda}(g(\theta, Z))] , \quad (65)$$

*where  $\psi_* \in \arg \max_\psi H_\lambda(\psi, \theta)$ , and where  $\nabla \psi_*^{c, \lambda}$  is given by (55).*

*Proof.* As in Theorem 3, the proof is based on the envelope Theorem 7 applied on  $h_\lambda(\theta) = W_\lambda(\mu_\theta, \nu) = \max_\psi H_\lambda(\psi, \theta)$  (see (25) and (27)). With the entropic regularization, the uniqueness (up to additive constants) of the solutions of the dual problem directly provides a continuous selection of Kantorovich potentials (thanks to Lemma 3). Besides, Lemma 8 ensures that  $\nabla_\theta F_\lambda(\psi, \theta)$  exists for any  $\psi \in \mathcal{C}(\mathcal{Y})$  and any  $\theta \in \Theta$ . It remains to show that  $\nabla_\theta F_\lambda$  is continuous in  $(\psi, \theta)$ . For that, recall that  $\nabla_\theta F_\lambda(\psi, \theta) = \mathbb{E}[\nabla_\theta f_\lambda(\psi, \theta, Z)]$  where  $f_\lambda$  is defined in (61), and let us fix  $\theta_0 \in \Theta$  and  $\psi_{*0}$  an optimal Kantorovich potential for  $W_\lambda(\mu_{\theta_0}, \nu)$ , and let also  $\psi \in \mathcal{C}(\mathcal{Y})$  be arbitrary. Then, for any  $\theta$  in the neighborhood  $V$  of  $\theta_0$  given by Hypothesis  $(G_{\theta_0})$ ,

$$\|f_\lambda(\psi, \theta, Z) - f_\lambda(\psi, \theta_0, Z)\| \leq L \|g(\theta, Z) - g(\theta_0, Z)\| \leq LK(Z) \|\theta - \theta_0\| \quad \text{a.s.} \quad (66)$$

which ensures that all components of  $\nabla_\theta f_\lambda(\psi, \theta, Z)$  are almost surely bounded by  $LK(Z)$ , an integrable bound which does not depend on  $(\psi, \theta)$ . Since  $\theta \mapsto g(\theta, Z)$  is almost surely  $\mathcal{C}^1$  on  $V$ , using (62), the last thing to show is that  $(\psi, \theta) \mapsto \nabla \psi^{c, \lambda}(g_\theta(Z))$  is almost surely continuous.

Let us fix  $z \in \mathcal{Z}$  for which  $g(\cdot, z)$  is  $\mathcal{C}^1$  on  $V$  and for which (29) holds. Thanks to (59), we can write

$$\nabla \psi^{c, \lambda}(g_\theta(z)) = e^{\frac{\psi^{c, \lambda}(g_\theta(z))}{\lambda}} \int_{\mathcal{Y}} \nabla_x c(g_\theta(z), y) e^{\frac{\psi(y) - c(g_\theta(z), y)}{\lambda}} d\nu(y). \quad (67)$$

For the first term, if  $\psi, \chi \in \mathcal{C}(\mathcal{Y})$  and  $\theta, \tau \in \Theta$ ,

$$|\psi^{c, \lambda}(g_\theta(z)) - \chi^{c, \lambda}(g_\tau(z))| \leq |\psi^{c, \lambda}(g_\theta(z)) - \psi^{c, \lambda}(g_\tau(z))| + |\psi^{c, \lambda}(g_\tau(z)) - \chi^{c, \lambda}(g_\tau(z))| \quad (68)$$

$$\leq L |g_\theta(z) - g_\tau(z)| + \|\psi^{c, \lambda} - \chi^{c, \lambda}\|_\infty \quad (69)$$

$$\leq LK(z) \|\theta - \tau\| + \|\psi - \chi\|_\infty , \quad (70)$$

by virtue of Lemma 2 and thus  $(\psi, \theta) \mapsto e^{\frac{\psi^{c, \lambda}(g_\theta(z))}{\lambda}}$  is continuous. For the integral term,  $\theta \mapsto \nabla_x c(g_\theta(z), y)$  is continuous on  $V$  because  $c$  is  $x$ -regular and because  $\theta \mapsto g(\theta, z)$  is continuous on  $V$ . Additionally,  $(\psi, \theta) \mapsto e^{\frac{\psi(y) - c(g_\theta(z), y)}{\lambda}}$  is a separable product of two terms which are continuous in  $\psi$  and  $\theta$  respectively. Finally, if  $\psi$  is restricted in a neighborhood  $P_0$  of  $\psi_{*0}$  bounded by  $R > 0$ , then all terms under the integral are bounded by  $Le^{\frac{1}{\lambda}(R + \|c\|_\infty)}$ . Again, the dominated convergence theorem implies that  $(\psi, \theta) \mapsto \nabla \psi^{c, \lambda}(g_\theta(z))$  is continuous on  $P_0 \times V$ . Finally, using the bound (64) and the dominated convergence theorem give that  $(\psi, \theta) \mapsto \nabla_\theta F_\lambda(\psi, \theta)$  is continuous on  $P_0 \times V$ . Therefore, all required conditions are satisfied to apply the envelope theorem on  $F_\lambda$ , which gives the desired formula.  $\square$

#### 4.4 Gradient of the Sinkhorn divergence

In this whole paragraph, we assume that  $\mathcal{X} = \mathcal{Y}$ , that  $(x, y) \mapsto c(x, y)$  is  $\mathcal{C}^1$  on  $\mathcal{X} \times \mathcal{X}$  and symmetric (i.e.  $c(x, y) = c(y, x)$  for all  $x, y \in \mathcal{X}$ ) and that it is  $L$ -Lipschitz with respect to  $(x, y)$ :

$$\forall (x, y), (x', y') \in \mathcal{X} \times \mathcal{X}, \quad |c(x, y) - c(x', y')| \leq L(\|x - x'\| + \|y - y'\|). \quad (71)$$

The authors of (Genevay et al., 2018) have shown that learning a generative model based on the Wasserstein cost  $W_\lambda$  induces a bias. For this reason, they propose to use instead the so-called Sinkhorn divergence defined as

$$S_\lambda(\mu, \nu) = W_\lambda(\mu, \nu) - \frac{1}{2} \left( W_\lambda(\mu, \mu) + W_\lambda(\nu, \nu) \right). \quad (72)$$

Since we focus here on the regularized WGAN learning problem, we study the function

$$s_\lambda(\theta) = S_\lambda(\mu_\theta, \nu) = W_\lambda(\mu_\theta, \nu) - \frac{1}{2} \left( W_\lambda(\mu_\theta, \mu_\theta) + W_\lambda(\nu, \nu) \right) \quad (73)$$

and we extend the previous regularity results to this new criterion. The first term of (73) has already been studied, and the last term does not depend on  $\theta$ . It thus remains to study the middle term, and for that we rely on the dual formulation of  $W_\lambda(\mu_\theta, \mu_\theta)$  given by

$$W_\lambda(\mu_\theta, \mu_\theta) = \max_{\chi, \eta \in \mathcal{C}(\mathcal{X})} \mathcal{F}_\lambda(\chi, \eta, \theta) \quad (74)$$

$$\text{where } \mathcal{F}_\lambda(\chi, \eta, \theta) = \int_{\mathcal{X}} \chi d\mu_\theta + \int_{\mathcal{X}} \eta d\mu_\theta - \lambda \int_{\mathcal{X} \times \mathcal{X}} e^{\frac{\chi(x) + \eta(y) - c(x, y)}{\lambda}} d\mu_\theta(x) d\mu_\theta(y) + \lambda. \quad (75)$$

Again, the problem (74) can be restricted to functions which are regularized  $c$ -transforms with respect to  $\mu_\theta$ . Similar to Lemma 8, one can show that such regularized  $c$ -transforms are still  $\mathcal{C}^1$  on  $\mathcal{X}$  and that the gradient can be computed by differentiating under the integral. Besides, because of the symmetry of  $W_\lambda(\mu_\theta, \mu_\theta)$ , one can show that the couple of optimal potentials  $(\chi_*, \eta_*)$  that solve (74) are such that  $\chi_*$  and  $\eta_*$  are equal up to an additive constant

Based on these observations, we can now state the regularity result for the Sinkhorn divergence. Recall that notations  $I$ ,  $F_\lambda$  are defined in (24) and (26).

**Theorem 6.** *Assume that  $c$  is  $x$ -regular and that  $g$  satisfies Hypothesis  $(G_\Theta)$ .*

*Then  $s_\lambda$  is  $\mathcal{C}^1$  on  $\Theta$  and*

$$\forall \theta \in \Theta, \quad \nabla_\theta s_\lambda(\theta) = \nabla_\theta F_\lambda(\psi_*, \theta) - \nabla_\theta I(\chi_*, \theta) \quad (76)$$

$$= \mathbb{E} \left[ D_\theta g(\theta, Z)^T \left( \nabla \psi_*^{c, \lambda}(g(\theta, Z)) - \nabla \chi_*(g(\theta, Z)) \right) \right], \quad (77)$$

where  $\psi_* \in \arg \max_\psi F_\lambda(\psi, \theta)$  and  $(\chi_*, \eta_*) \in \arg \max_{(\chi, \eta)} \mathcal{F}_\lambda(\chi, \eta, \theta)$ .

*Proof.* Again, using the result obtained for  $\theta \mapsto W_\lambda(\mu_\theta, \nu)$  in the last paragraphs, we only have to study the regularity of  $W_\lambda(\mu_\theta, \mu_\theta)$ . Once again, this follows from the envelope theorem recalled in Appendix A. For



that we let  $C$  the set of  $\mathcal{C}^1$  and  $L$ -Lipschitz functions on  $\mathcal{X}$  equipped with the norm  $\|\chi\|_\infty + \|\nabla\chi\|_\infty$  and we use the dual expression (74) that can be written

$$\mathcal{F}_\lambda(\chi, \eta, \theta) = I(\chi, \theta) + I(\eta, \theta) + \lambda - \lambda E_\lambda(\chi, \eta, \theta) \quad (78)$$

$$\text{where } E_\lambda(\chi, \eta, \theta) = \mathbb{E} \left[ \exp \left( \frac{\Gamma(\chi, \eta, \theta, Z, W)}{\lambda} \right) \right] \quad (79)$$

$$\text{and } \Gamma(\chi, \eta, \theta, z, w) = \chi(g_\theta(z)) + \eta(g_\theta(w)) - c(g_\theta(z), g_\theta(w)), \quad (80)$$

and where  $W, Z$  are two independent random variables of distribution  $\zeta$ . For any  $\chi \in C$ , differentiating under the integral as in Lemma 8 gives again that  $I(\chi, \cdot)$  is differentiable on  $\Theta$  with gradient

$$\nabla_\theta I(\chi, \theta) = \mathbb{E} [D_\theta g(\theta, Z)^T \nabla_\chi(g(\theta, Z))]. \quad (81)$$

By the same reasoning, one obtains that  $E_\lambda(\chi, \eta, \cdot)$  is differentiable on  $\Theta$  with gradient

$$\nabla_\theta E_\lambda(\chi, \eta, \theta) = \mathbb{E} \left[ \frac{\nabla_\theta \Gamma(\chi, \eta, \theta, Z, W)}{\lambda} \exp \left( \frac{\Gamma(\chi, \eta, \theta, Z, W)}{\lambda} \right) \right] \quad (82)$$

with

$$\nabla_\theta \Gamma(\chi, \eta, \theta, z, w) = D_\theta g(\theta, z)^T (\nabla_\chi(g(\theta, z)) + \nabla_x c(g(\theta, z), g(\theta, w))) \quad (83)$$

$$+ D_\theta g(\theta, w)^T (\nabla_\eta(g(\theta, w)) + \nabla_y c(g(\theta, z), g(\theta, w))). \quad (84)$$

If  $P_0$  is any bounded set of  $C$ , using a bound similar to (64) and the dominated convergence theorem shows that  $(\chi, \theta) \mapsto I(\chi, \theta)$  is continuous on  $P_0 \times V$ . Similarly, one can see that  $\Gamma$  is Lipschitz in  $\theta$  with a  $(Z, W)$ -integrable bound

$$|\Gamma(\chi, \eta, \theta, z, w) - \Gamma(\chi, \eta, \tau, z, w)| \leq (\|\nabla\chi\|_\infty K(z) + \|\nabla\eta\|_\infty K(w) + L) \|\theta - \tau\|. \quad (85)$$

Since  $\Gamma$  is also bounded by  $\|\chi\|_\infty + \|\psi\|_\infty + \|c\|_\infty$  and since the continuity of  $\Gamma$  with respect to  $(\chi, \eta, \theta)$  can be deduced from (80), this is enough to show that  $(\chi, \eta, \theta) \mapsto \nabla E_\lambda(\chi, \eta, \theta)$  is continuous on  $P_0 \times V$ .

Therefore, we have again all conditions required to apply the envelope theorem: we have a continuous selection of optimal dual variables  $(\chi, \eta)$  and the gradient  $\nabla_\theta \mathcal{F}_\lambda(\chi, \eta, \theta)$  exists and is continuous in  $(\chi, \eta, \theta)$ , and we can differentiate under the max in (74)

$$\nabla_\theta W_\lambda(\mu_\theta, \mu_\theta) = \nabla_\theta \mathcal{F}_\lambda(\chi_*, \eta_*, \theta), \quad (86)$$

where  $(\chi_*, \eta_*)$  is a pair of Kantorovich potentials for  $W_\lambda(\mu_\theta, \mu_\theta)$ . Finally, notice that by symmetry, there exists  $k \in \mathbb{R}$  such that  $\chi_* = \eta_* + k$  and by definition  $E_\lambda(\chi_*, \eta_*, \theta) = 1$  so that

$$\mathcal{F}_\lambda(\chi_*, \eta_*, \theta) = 2I(\chi_*, \theta) - k \quad (87)$$

and therefore

$$\nabla_\theta \mathcal{F}_\lambda(\chi_*, \eta_*, \theta) = 2\nabla I(\chi, \theta) = 2\mathbb{E} [D_\theta g(\theta, Z)^T \nabla_{\chi_*}(g(\theta, Z))]. \quad (88)$$

Putting all terms together, we get the desired formula for  $\nabla_\theta s_\lambda(\theta)$ .  $\square$

## 5 Interpretation with Derivatives in the Sense of Distributions

This section contains a brief discussion on the way to interpret the previously obtained gradient formulae, by taking derivatives in the sense of distributions of  $\theta \mapsto \mu_\theta$ .

In the proofs above, one can notice that a crucial argument is to examine the regularity of the function

$$I(\varphi, \theta) = \int_{\mathcal{X}} \varphi d\mu_\theta. \quad (89)$$

In order to better understand the behavior of this function, it is useful to consider the map  $\theta \mapsto \mu_\theta$  from  $\Theta$  to the space  $\mathcal{D}'(U)$  of distributions on  $U$ , which is the dual of the space  $\mathcal{D}(U)$  of compactly-supported  $\mathcal{C}^\infty$  functions on  $U$  (detailed definitions of these functional spaces can be found in (Hörmander, 2015)). Then, using the duality product on  $\mathcal{D}'(U)$ , one can rewrite

$$I(\varphi, \theta) = \langle \mu_\theta, \varphi \rangle. \quad (90)$$

Therefore, the regularity of (89) can be understood as the regularity of  $\theta \mapsto \mu_\theta$ , after evaluation against  $\varphi$ . In order to benefit from the framework of distribution derivatives, it is useful to work on the product  $\Theta \times U$ .

It is possible to define  $T \in \mathcal{D}'(\Theta \times U)$  by setting

$$\forall \Phi \in \mathcal{D}(\Theta \times U), \quad \langle T, \Phi \rangle = \int_{\Theta} \int_U \Phi(\theta, x) \mu_\theta(dx) d\theta. \quad (91)$$

For any  $\Phi \in \mathcal{D}(\Theta \times U)$  whose support is included in a compact  $K \times L$ ,

$$|\langle T, \Phi \rangle| \leq \|\Phi\|_\infty \int_{\Theta} \int_U d\mu_\theta(x) d\theta \leq \|\Phi\|_\infty \int_K d\theta \quad (92)$$

since  $\mu_\theta$  is a probability distribution on  $U$ . This proves that  $T$  defines a 0-th order distribution on  $\Theta \times U$ . In this context, it is possible to give a meaning to the pointwise evaluation of  $T$  at  $\theta$ , which corresponds to the probability distribution  $T(\theta, \cdot) = \mu_\theta$ .

Now, the distributional derivatives of  $T$  w.r.t. the variable  $\theta_i$  can be written with the limit in the  $\mathcal{D}'$  sense:

$$\partial_{\theta_i} T = \lim_{h \rightarrow 0} \frac{\tau_{-ht_i} T - T}{h} \quad (93)$$

where  $(t_1, \dots, t_p)$  is the canonical basis of  $\mathbb{R}^p$ , and where  $\tau_v T$  is the translation of  $T$  with vector  $v$ , defined by  $\langle \tau_v T, \Phi \rangle = \langle T, \Phi(\cdot + v) \rangle$ . In other words,

$$\forall \Phi \in \mathcal{D}(\Theta \times U), \quad \langle \partial_{\theta_i} T, \Phi \rangle = \lim_{h \rightarrow 0} \int_{\Theta} \int_U \Phi(\theta, x) d\left(\frac{\mu_{\theta+ht_i} - \mu_\theta}{h}\right)(x) d\theta. \quad (94)$$

Since  $T$  is a distribution of order 0, the partial derivative  $\partial_{\theta_i} T$  is a distribution of order  $\leq 1$ . This explains why the expression of  $\partial_{\theta_i} T$  may involve  $\nabla \Phi$ .

In the case where  $\mu_\theta$  is the output distribution of a generative network  $g$  that satisfies Hypothesis (G $_\Theta$ ), then for a separable test function  $\Phi(\theta, x) = \alpha(\theta)\varphi(x)$  with  $\alpha \in \mathcal{D}(\Theta)$  and  $\varphi \in \mathcal{D}(U)$ , we have

$$\int_{\Theta} \int_U \alpha(\theta) \varphi(x) d\left(\frac{\mu_{\theta+ht_i} - \mu_\theta}{h}\right)(x) d\theta = \int_{\Theta} \alpha(\theta) \mathbb{E} \left[ \frac{\varphi(g(\theta + ht_i, Z)) - \varphi(g(\theta, Z))}{h} \right] d\theta. \quad (95)$$

With arguments similar to the last sections, one can show that this limit is well-defined, which gives

$$\langle \partial_{\theta_i} T(\theta, x), \alpha(\theta) \varphi(x) \rangle = \int_{\Theta} \alpha(\theta) \mathbb{E} [D_{\theta_i} g(\theta, Z)^T \nabla \varphi(g(\theta, Z))] d\theta. \quad (96)$$

In other words, the pointwise evaluation of  $\partial_{\theta_i} T$  at  $\theta$  is identified to the distribution of order  $\leq 1$  given by

$$\forall \varphi \in \mathcal{D}(U), \quad \langle \partial_{\theta_i} T(\theta, \cdot), \varphi \rangle = \mathbb{E} [D_{\theta_i} g(\theta, Z)^T \nabla \varphi(g(\theta, Z))]. \quad (97)$$

In conclusion, it is possible to interpret the results of the last sections in terms of distributional derivatives of  $\theta \mapsto \mu_\theta$  and this explains the apparition of  $\nabla \varphi$  in the gradient formulae found in the last sections.

## 6 Experiments

In this section, we provide an algorithm that tackles the practical minimization of (23) in the case of generator learning for handwritten digits, using the MNIST database (LeCun et al., 1998). We do not seek to reach

state-of-the-art results for generator learning with complex databases, but we focus instead on a practical analysis of the convergence of the proposed algorithm (in terms of loss values) and examine the impact of the regularization parameters.

In this setting,  $\mathcal{Y}$  is a finite set of  $J$  data points, so that elements  $\psi \in \mathcal{C}(\mathcal{Y})$  identifies to vectors in  $\mathbb{R}^J$  as in Section 3. On the other hand,  $\mu_\theta$  is the output distribution of the generative network, which can be sampled on demand. In this section,  $\lambda \geq 0$ .

## 6.1 Alternate Algorithm

We aim at solving the optimization problem

$$\min_{\theta \in \Theta} h_\lambda(\theta) = \min_{\theta \in \Theta} \max_{\psi \in \mathcal{C}(\mathcal{Y})} H_\lambda(\psi, \theta) \quad (98)$$

$$\text{where } H_\lambda(\psi, \theta) = \int_{\mathcal{X}} \psi^{c,\lambda} d\mu_\theta + \int_{\mathcal{Y}} \psi d\nu = F_\lambda(\psi, \theta) + \sum_{y \in \mathcal{Y}} \psi(y) \nu(\{y\}) . \quad (99)$$

We adopt an algorithm that alternates between updating  $\theta$  with one gradient step, and updating  $\psi$  with several iterations of a dedicated algorithm.

As we have seen above in (53) and (65), computing the gradient of  $h_\lambda$  requires to compute an optimal dual potential  $\varphi_* = \psi_*^{c,\lambda}$ , which, in general, cannot be done exactly. Instead, we rely on the stochastic algorithm for semi-discrete OT proposed in (Genevay et al., 2016) to approximate the optimal potential. In the proofs of the previous sections, we have written

$$F_\lambda(\psi, \theta) = \mathbb{E}[f_\lambda(\psi, \theta, Z)] \quad (100)$$

where  $f_\lambda(\psi, \theta, z) = \psi^{c,\lambda}(g_\theta(z))$ . It has been shown in (Genevay et al., 2016; Houdard et al., 2022) that  $H_\lambda(\cdot, \theta)$  defined in (99) is a concave function whose supergradient  $\mathcal{D}(\psi, \theta) = \partial_\psi H_\lambda(\psi, \theta)$  can be written as

$$\mathcal{D}(\psi, \theta) = \mathbb{E}[\mathcal{D}(\psi, \theta, Z)] \quad \text{where} \quad \mathcal{D}(\psi, \theta, z) = \partial_\psi \left( f_\lambda(\psi, \theta, z) - \int \psi d\nu \right). \quad (101)$$

This element  $\mathcal{D}(\psi, \theta, z) \in \mathbb{R}^J$  can be computed with an explicit formula given in (Genevay et al., 2016; Houdard et al., 2022) (and in practice is implemented by automatic differentiation).

Therefore, for a current  $\theta$ , we can optimize  $\psi$  with a stochastic supergradient ascent

$$\begin{cases} \tilde{\psi}_k &= \tilde{\psi}_{k-1} + \frac{\gamma}{k^\alpha} \left( \frac{1}{|B_k|} \sum_{z \in B_k} \mathcal{D}(\tilde{\psi}_{k-1}, \theta, z) \right) \\ \psi_k &= \frac{1}{k} (\tilde{\psi}_1 + \dots + \tilde{\psi}_k), \end{cases} \quad (102)$$

where  $\gamma > 0$  is the learning rate,  $\alpha \in (0, 1)$  a parameter,  $B_k$  is a batch containing  $b$  independent samples of the distribution of  $Z$ , and the different batches  $B_k$ 's are also independent. As recalled in (Genevay et al., 2016) and (Galerie et al., 2018), for  $\alpha = 0.5$ , this algorithm has a convergence guarantee in  $\mathcal{O}(\frac{\log k}{\sqrt{k}})$  on the function values. After  $K$  iterations of the inner loop, we obtain an approximation  $\underline{\psi}$  of the optimal dual potential  $\psi_*$  for  $W_\lambda(\mu_\theta, \nu)$ , and we use it to perform the gradient descent step on  $\theta$

$$\nabla_\theta h_\lambda(\theta) \approx \nabla_\theta H_\lambda(\underline{\psi}, \theta) = \nabla_\theta F_\lambda(\underline{\psi}, \theta) = \mathbb{E} \left[ D_\theta g(\theta_0, Z)^T \nabla_{\underline{\psi}} \psi^{c,\lambda}(g(\theta_0, Z)) \right] \quad (103)$$

Actually, using  $\nabla_\theta H_\lambda(\underline{\psi}, \theta)$  as a proxy for  $\nabla_\theta h_\lambda(\theta)$  in the gradient descent step can be simply reinterpreted as saying that we perform an alternate optimization on  $H_\lambda(\psi, \theta)$ , with an inner loop on  $\psi$  at each iteration. Again, the expectation in (103) cannot be computed in closed form so that we realize an approximation by taking another batch  $B'$  on  $z$ :

$$\nabla_\theta H_\lambda(\underline{\psi}, \theta) \approx \frac{1}{|B'|} \sum_{z \in B'} D_\theta g(\theta_0, z)^T \nabla_{\underline{\psi}} \psi^{c,\lambda}(g(\theta_0, z)). \quad (104)$$

The overall algorithm is summarized in Algorithm 1. Notice that in the case of unregularized OT  $\lambda = 0$ , the algorithm is close to the one proposed by Chen et al. (Chen et al., 2019). But all the computations made in the present paper allow to interpret it as a stochastic alternate optimization algorithm on a fixed cost  $H_\lambda(\psi, \theta)$ , thus including naturally the regularized case  $\lambda > 0$ . One benefit of this approach is that the stochastic gradient steps taken on  $\psi$  and  $\theta$  can be implemented by automatic differentiation on a fixed cost  $H_\lambda$ , and the corresponding updates can be implemented with predefined optimizers. In particular, this opens the possibility to adopt other parameterizations of the variable  $\psi$ . Indeed, while the usual stochastic algorithm for semi-discrete OT (Genevay et al., 2016) works on the vector  $(\psi(y_j)) \in \mathbb{R}^J$ , it is also possible to adopt a neural network parameterization  $\psi_t$  of  $\psi$  as in (Seguy et al., 2018). The update of  $\psi$  then translates on an update of the neural network parameters  $t$ , which can be done by backpropagating the gradient of

$$t \mapsto f_\lambda(\psi_t, \theta, z) - \int \psi_t d\nu = \psi_t^{c, \lambda}(g_\theta(z)) - \sum_{y \in \mathcal{Y}} \psi_t(y) \nu(\{y\}) . \quad (105)$$

Finally, let us mention that a limitation of this approach is that computing the gradients of  $H_\lambda$  requires the differentiation of  $f_\lambda(\psi, \theta, z) = \psi^{c, \lambda}(g_\theta(z))$ , for a batch of  $z$  values. The exact computation of  $\psi^{c, \lambda}$  requires to visit all the dataset  $\mathcal{Y}$ , which is prohibitive for a very large database.

---

#### Algorithm 1

---

**Initialization:**  $\psi_0 = 0$ , random initialization of  $\theta$

$n = 1$  **to**  $N$

- Approximate  $\psi_n \approx \arg \max H_\lambda(\cdot, \theta)$ : inner loop with  $K$  iterations of ASGD (102) starting from  $\psi_{n-1}$ , using batches  $B_{n,1}, \dots, B_{n,K}$  of size  $b$  on  $z$
- Update  $\theta$  with one step of ADAM algorithm on  $H_\lambda(\psi_n, \cdot)$  using gradient (104) computed on a batch  $B'_n$  of size  $b$  on  $z$

**end for**

**Output:** estimated generative model parameter  $\theta$

---

**Detailed setting** The following paragraph gathers the parameters and network architectures used in the experiments shown in this section.

- $N = 3000$  iterations on  $\theta$  with ADAM algorithm (with learning rate 0.001)
- $K = 10$  iterations of the inner loop with ASGD algorithm (see learning rates below)
- The cost  $c(x, y)$  is the quadratic cost  $\alpha^{-1} \|x - y\|^2$  normalized by  $\alpha = \frac{1}{J} \sum_{y \in \mathcal{Y}} \|y\|^2$
- For the generator  $g_\theta$  we consider two different architectures:
  - a multilayer perceptron (MLP) with four fully-connected layers; the number of channels for the successive hidden layers is 256, 512, 1024.
  - a Deep Convolutional Adversarial Network (DCGAN) (Radford et al., 2015) adapted for the dimension  $28 \times 28$  of MNIST images, with four deconvolution layers; the number of channels for the successive hidden layers is 256, 128, 64.
- The input of these generators is a random variable  $Z$  following the uniform distribution on  $[-1, 1]^{100}$  (the choice of dimensionality 100 is commonly encountered when fitting a generative network to the MNIST database). All batches of  $Z$  are made of 200 samples.
- The dual variable  $\psi$  is either directly modeled by a vector  $\psi \in \mathbb{R}^J$ , or parameterized by a multilayer perceptron with four fully-connected layers; the number of channels for the successive hidden layers is 512, 256, 128. These two different settings are respectively referred to as SDOT (for semi-dual OT) and SDOTNN (for semi-dual OT with neural network).

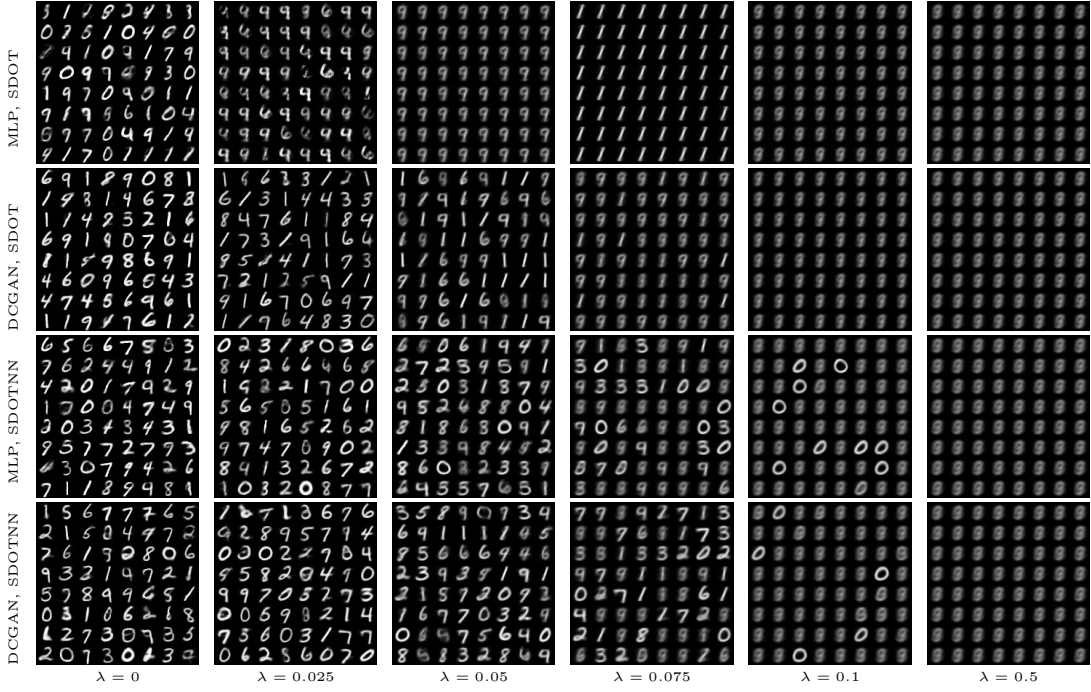


Figure 1: Generation of MNIST digits. We compare here several generative networks trained with different architectures for the generator  $g_\theta$  (MLP or DCGAN) and the dual variable  $\psi$  with no parameterization (SDOT) in the two first rows, or MLP parameterization (SDOTNN) in the two last rows, and varying the parameter  $\lambda$  of entropic regularization.

- The step size strategy for the ASGD inner loop has been chosen in the following manner: for the parameterization SDOT, we use  $\gamma = 5, \alpha = 0.5$  i.e. a step size  $\frac{5}{k^{0.5}}$  while for the parameterization SDOTNN, we use  $\gamma = 0.1, \alpha = 0.8$  i.e. a step size  $\frac{0.1}{k^{0.8}}$ .

## 6.2 Impact of the regularization parameter

In this paragraph, we discuss the behavior of the alternate Algorithm 1 and examine the impact of the regularization parameter  $\lambda$  both on the visual results and the convergence of the loss function. We also discuss the effect of parameterizing the dual variable  $\psi$  by a neural network.

On Fig. 1, we display sampled digits obtained with the generative networks learned with Algorithm 1 run with different settings. One can see that the generators learned with unregularized OT ( $\lambda = 0$ ) produce mostly convincing samples which are slightly more blurry than the images of the database. Some of samples do not exactly resemble a digit but some kind of mixing between different digits, which reflects the fact that the generative network naturally interpolates between the images of the database. The two tested architectures for the generator produce comparable results, with a slight advantage for DCGAN. DCGAN indeed provides cleaner samples that better cover the whole database, thanks to its more complex architecture, well adapted to two-dimensional data.

The visual results deteriorate when the regularization parameter  $\lambda$  grows. For very small  $\lambda$ , the results are still comparable to the unregularized case. For larger  $\lambda$ , the outputs of the generative network seem to concentrate on a blurry average of the database. This can be understood by looking at the gradient formula (65) which involves the gradient of the regularized  $c$ -transform given by (55). When  $\lambda \rightarrow +\infty$ ,  $\nabla \psi^{c,\lambda}(x)$  degenerates to a simple average  $\int_Y \nabla_x c(x, y) d\nu(y)$ . In other words, with the blur created by entropic regularization on the transport plan, the sampled points are pushed towards all the target points in a mixed manner. In contrast, with unregularized transport, each sampled point  $x$  is pushed towards the data point  $T_\psi(x)$  assigned by the current OT map.

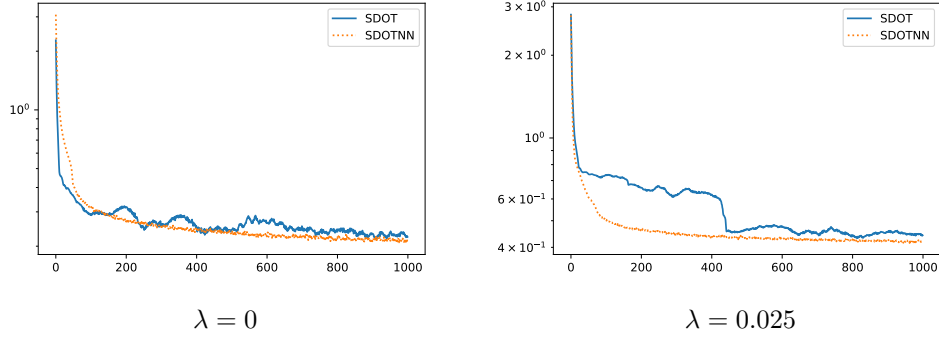


Figure 2: Evolution of the loss  $H_\lambda(\underline{\psi}_\theta, \theta)$  along the iterates of Algorithm 1 to learn a DCGAN for generation of MNIST digits. For each iterate  $\theta$ , the loss is computed by using the current estimate  $\underline{\psi}_\theta$  of the dual variable. For two values of the regularization parameter ( $\lambda = 0$  on the left and  $0.025$  on the right), we compare the OT loss values obtained by parameterizing  $\psi$  directly as a vector in  $\mathbb{R}^J$  (SDOT) or as a neural network (SDOTNN). See the text for comments.

To better understand these visual results, we now examine the behavior of the loss function depending on the adopted setting. Fig. 2 shows the evolution of the loss  $H_\lambda(\underline{\psi}_\theta, \theta)$  (for the current dual variable  $\underline{\psi}_\theta$ , with  $H_\lambda$  defined in (99)) along the iterates of Algorithm 1, when the dual variable  $\psi$  is either parameterized as a vector in  $\mathbb{R}^J$  (SDOT) or a neural network (SDOTNN). One can see that the loss stabilizes in  $\approx 500$  iterations, and that the limit values obtained with both parameterizations (SDOT and SDOTNN) are very similar. It is interesting to notice that the limit value is even lower with the SDOTNN parameterization: since the adopted multilayer perceptron has here  $> 5 \cdot 10^5$  parameters (and is thus much larger than  $J = 6 \cdot 10^4$ ), it is likely that any value  $(\psi(y_j))_{1 \leq j \leq J} \in \mathbb{R}^J$  can be attained with such a parameterization for  $\psi$ . Notice also that the loss decreases in a more stable way with the SDOTNN parameterization: this parameterization is indeed likely to be more robust to the individual changes on  $\psi(y_j)$  when updating the parameters  $\theta$  of the generator.

One can notice that, quite surprisingly, the convergence speed does not improve drastically when using a larger regularization parameter  $\lambda$ . This is confirmed in Fig. 3 where we display results obtained with various regularization parameters  $\lambda$  and the four tested combinations of architectures for the generative network and the dual variable. As expected, increasing the regularization parameter leads to a smoother optimized functional, which reflects in a more stable evolution of the loss. For very small  $\lambda$  ( $\leq 0.025$ ), we observe that the regularization does not improve the convergence speed with respect to  $\theta$ . In this slightly regularized regime, we suggest that the behavior of the loss evolution mostly depends on the chosen architecture.

To further analyze the algorithm, for a fixed generator parameter  $\theta$ , we display in Fig. 4 the evolution of the loss  $\psi \mapsto H_\lambda(\psi, \theta)$  (defined in (99)) during the inner ASGD loop used for optimizing the dual variable  $\psi$ . In order to complete the comparison, we also include the convergence plot obtained with the ADAM algorithm applied on the same problem. These convergence curves reflect the slow convergence rate (in  $\mathcal{O}(\frac{\log k}{\sqrt{k}})$ ) of the ASGD algorithm. In our experiments, we observe that a careful tuning of the learning rate  $\gamma$  in the ASGD algorithm (102) is necessary to obtain a sufficient decrease of the loss. Next, the convergence plots obtained with ASGD are similar with both parameterizations SDOT and SDOTNN. One can observe that for very small regularization, turning to the ADAM algorithm does not improve the convergence speed for the SDOT parameterization. However, we remark that using the ADAM algorithm with the SDOTNN parameterization seems beneficial for all tested regularization parameters: the loss value obtained after 100 iterations is lower than with the SDOT parameterization, and the convergence is much faster.

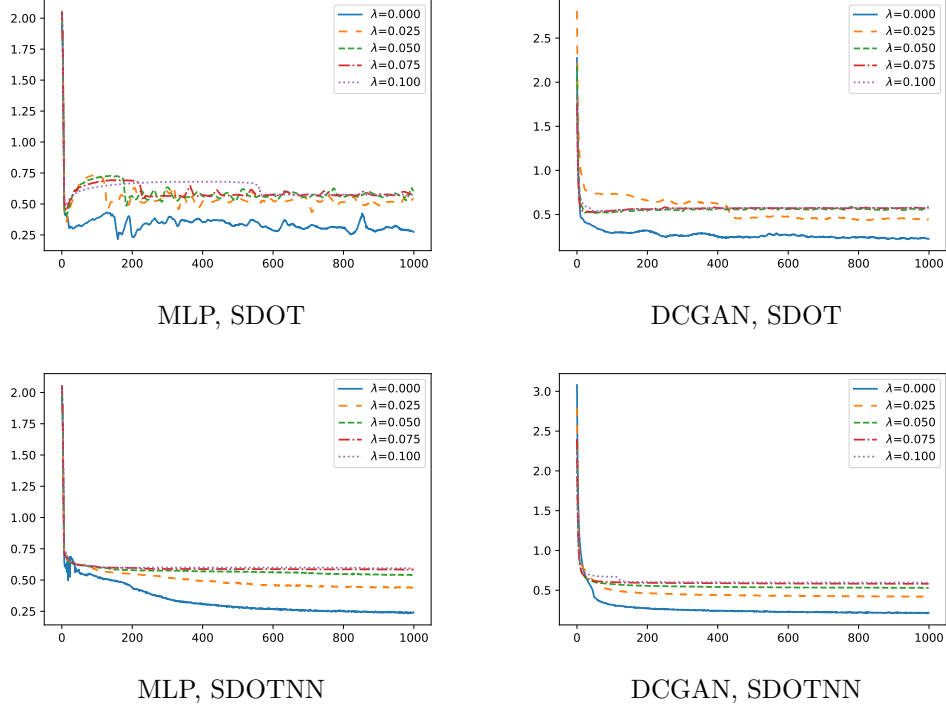


Figure 3: Evolution of the loss  $H_\lambda(\psi_\theta, \theta)$  along the iterates of Algorithm 1, for the four tested combinations of parameterizations of the generator (MLP or DCGAN) and the dual variable (SDOT or SDOTNN). For each iterate  $\theta$ , the loss is computed by using the current estimate  $\psi_\theta$  of the dual variable. Let us recall that the loss function  $H_\lambda$  depends on the regularization parameter  $\lambda$ , which explains why the limit value attained by the algorithm actually increases when  $\lambda \rightarrow 0$ . See the text for additional comments.

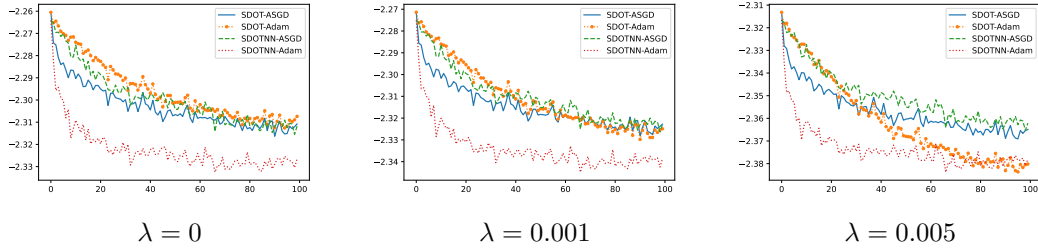


Figure 4: Evolution of the loss  $\psi \mapsto H_\lambda(\psi, \theta)$  along the iterates of the inner loop of Algorithm 1. Here, the parameter  $\theta$  of the DCGAN generator is fixed, i.e. we consider a semi-discrete OT problem between a fixed  $\mu_\theta$  and  $\nu$ . For several values of the regularization parameter, we compare the evolution of the loss when parameterizing  $\psi$  directly as a vector in  $\mathbb{R}^J$  (SDOT) or as a neural network (SDOTNN). For both parameterizations, the optimization is done using either ASGD with decreasing step size ( $\frac{5}{\sqrt{k}}$  for SDOT and  $\frac{0.1}{k^{0.8}}$  for SDOTNN) or ADAM (with learning rate 0.001). See the text for comments.

---

## 7 Conclusion

In this paper we gave new insights on the theory and practice for learning generative networks with regularized Wasserstein distances. On the theoretical side, we proved a gradient formula for the minimized loss in two different frameworks: in the semi-discrete (i.e. when the target distribution  $\nu$  has finite support) without regularization, and in a more general case (with a general  $\nu$ ) with entropic regularization. These results are based on a regularity hypothesis on the generator, and also, in the semi-discrete case, an assumption that the generator does not charge the boundary of Laguerre cells. These hypotheses are helpful to better understand the possible degenerate cases that can be encountered, and we provided such a counterexample.

On the practical side, we showed that an alternate algorithm can approximate the solution of this optimization problem. The inner loop of this algorithm consists in approximating an optimal dual potential for regularized OT with a stochastic optimization algorithm. The results of the previous sections justify the existence of the gradients used in this alternate procedure. Experiments on MNIST digits demonstrate that this algorithm is able to learn a neural network generating relevant images. When approximating the dual variable, convincing visual results are indeed obtained with zero or small regularization parameter  $\lambda$ . For such a small regularization, the smoothing of the targeted loss function is not sufficient to drastically improve the convergence speed of the optimization algorithm for the generator parameters. However, for which concerns the stochastic optimization used to solve the semi-dual OT problem, we observed that it may be beneficial in terms of convergence speed to parameterize the dual variable with a neural network, provided that one uses a well-chosen and carefully tuned algorithm to optimize it. The improvement observed with such a parameterization remains to be explained with a thorough analysis of the ADAM algorithm applied on this semi-discrete OT problem.

The main drawback of the considered algorithm is that the inner loop is based on the computation of a regularized  $c$ -transform and thus requires, at each iteration, to visit all data points (in order to find a kind of biased nearest neighbor). In order to scale up to larger database, it has already been proposed (Mallasto et al., 2019) to approximate the regularized  $c$ -transform with a batch strategy. As a perspective, it would be interesting to see if the errors made at each iteration by this batch strategy can be controlled in order to get a globally stable optimization process.

## Acknowledgments

This study has been carried out with financial support from the French Research Agency through the GOTMI, Mystic and PostProdLEAP projects (ANR-16-CE33-0010-01, ANR-19-CE40-005 and ANR-19-CE23-0027-01) and from the GdR ISIS through the Remoga project.

## References

- Martin Arjovsky, Soumith Chintala, and Léon Bottou. Wasserstein generative adversarial networks. In *International Conference on Machine Learning*, pp. 214–223, 2017.
- Jérémie Bigot, Elsa Cazelles, and Nicolas Papadakis. Central limit theorems for entropy-regularized optimal transport on finite spaces and statistical applications. *Electronic Journal of Statistics*, 13(2):5120–5150, 2019.
- Ashish Bora, Ajil Jalal, Eric Price, and Alexandros G Dimakis. Compressed sensing using generative models. In *International Conference on Machine Learning*, pp. 537–546. PMLR, 2017.
- Yucheng Chen, Matus Telgarsky, Chao Zhang, Bolton Bailey, Daniel Hsu, and Jian Peng. A gradual, semi-discrete approach to generative network training via explicit wasserstein minimization. In *International Conference on Machine Learning*, pp. 1071–1080. PMLR, 2019.
- Lénaïc Chizat. *Unbalanced Optimal Transport: Models, Numerical Methods, Applications*. PhD thesis, Université Paris Dauphine, PSL, 2017.



- 
- Marco Cuturi. Sinkhorn distances: Lightspeed computation of optimal transport. In *Advances in neural information processing systems*, pp. 2292–2300, 2013.
- Marco Cuturi and Gabriel Peyré. A smoothed dual approach for variational wasserstein problems. *SIAM Journal on Imaging Sciences*, 9(1):320–343, 2016.
- Muhammad Fadli Damara, Gregor Kornhardt, and Peter Jung. Solving inverse problems with conditional-gan prior via fast network-projected gradient descent. *arXiv preprint arXiv:2109.01105*, 2021.
- Kilian Fatras, Younes Zine, Rémi Flamary, Remi Gribonval, and Nicolas Courty. Learning with minibatch wasserstein : asymptotic and gradient properties. In Silvia Chiappa and Roberto Calandra (eds.), *Proceedings of the Twenty Third International Conference on Artificial Intelligence and Statistics*, volume 108 of *Proceedings of Machine Learning Research*, pp. 2131–2141. PMLR, 26–28 Aug 2020.
- Jean Feydy, Thibault Séjourné, François-Xavier Vialard, Shun-ichi Amari, Alain Trounev, and Gabriel Peyré. Interpolating between optimal transport and mmd using sinkhorn divergences. In *The 22nd International Conference on Artificial Intelligence and Statistics*, pp. 2681–2690, 2019.
- Bruno Galerne, Arthur Leclaire, and Julien Rabin. A texture synthesis model based on semi-discrete optimal transport in patch space. *SIAM Journal on Imaging Sciences*, 11(4):2456–2493, 2018.
- Aude Genevay. *Entropy-regularized optimal transport for machine learning*. PhD thesis, Université Paris Dauphine, 2019.
- Aude Genevay, Marco Cuturi, Gabriel Peyré, and Francis Bach. Stochastic optimization for large-scale optimal transport. In *Advances in neural information processing systems*, pp. 3440–3448, 2016.
- Aude Genevay, Gabriel Peyré, and Marco Cuturi. Learning generative models with sinkhorn divergences. In Amos Storkey and Fernando Perez-Cruz (eds.), *Proceedings of the Twenty-First International Conference on Artificial Intelligence and Statistics*, volume 84 of *Proceedings of Machine Learning Research*, pp. 1608–1617. PMLR, 2018.
- Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. Generative adversarial nets. In *Advances in neural information processing systems*, pp. 2672–2680, 2014.
- Ishaan Gulrajani, Faruk Ahmed, Martin Arjovsky, Vincent Dumoulin, and Aaron C Courville. Improved training of wasserstein gans. In *Advances in neural information processing systems*, pp. 5767–5777, 2017.
- Paul Hand and Babhru Joshi. Global guarantees for blind demodulation with generative priors. *Advances in Neural Information Processing Systems*, 32, 2019.
- Reinhard Heckel and Mahdi Soltanolkotabi. Denoising and regularization via exploiting the structural bias of convolutional generators. In *International Conference on Learning Representations*, 2020.
- Lars Hörmander. *The analysis of linear partial differential operators I: Distribution theory and Fourier analysis*. Springer, 2015.
- Antoine Houdard, Arthur Leclaire, Nicolas Papadakis, and Julien Rabin. A generative model for texture synthesis based on optimal transport between feature distributions. *Journal of Mathematical Imaging and Vision*, 2022.
- Rakib Hyder and M Salman Asif. Generative models for low-dimensional video representation and reconstruction. *IEEE Transactions on Signal Processing*, 68:1688–1701, 2020.
- Diederik P Kingma and Max Welling. Auto-encoding variational {Bayes}. In *Int. Conf. on Learning Representations*, 2014.
- Arthur Leclaire and Julien Rabin. A stochastic multi-layer algorithm for semi-discrete optimal transport with applications to texture synthesis and style transfer. *Journal of Mathematical Imaging and Vision*, 63(2):282–308, 2021.

- 
- Yann LeCun, Léon Bottou, Yoshua Bengio, and Patrick Haffner. Gradient-based learning applied to document recognition. *Proceedings of the IEEE*, 86(11):2278–2324, 1998.
- Oscar Leong. *Learned Generative Priors for Imaging Inverse Problems*. PhD thesis, Rice University, 2021.
- Dong Liu, Minh Thành Vu, Saikat Chatterjee, and Lars K Rasmussen. Entropy-regularized optimal transport generative models. In *ICASSP 2019-2019 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pp. 3532–3536. IEEE, 2019.
- Anton Mallasto, Jes Frellsen, Wouter Boomsma, and Aasa Feragen. (q, p)-wasserstein gans: Comparing ground metrics for wasserstein gans. *arXiv preprint arXiv:1902.03642*, 2019.
- Sachit Menon, Alexandru Damian, Shijia Hu, Nikhil Ravi, and Cynthia Rudin. Pulse: Self-supervised photo upsampling via latent space exploration of generative models. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 2437–2445, 2020.
- Daisuke Oyama and Tomoyuki Takenawa. On the (non-) differentiability of the optimal value function when the optimal solution is unique. *Journal of Mathematical Economics*, 76:21–32, 2018.
- Gabriel Peyré and Marco Cuturi. Computational optimal transport. *Foundations and Trends® in Machine Learning*, 11(5-6):355–607, 2019.
- Alec Radford, Luke Metz, and Soumith Chintala. Unsupervised representation learning with deep convolutional generative adversarial networks. *arXiv preprint arXiv:1511.06434*, 2015.
- Maziar Sanjabi, Jimmy Ba, Meisam Razaviyayn, and Jason D Lee. On the convergence and robustness of training gans with regularized optimal transport. In *Advances in Neural Information Processing Systems*, pp. 7091–7101, 2018.
- F. Santambrogio. *Optimal transport for applied mathematicians*. Birkhäuser, NY, 2015.
- Vivien Seguy, Bharath Bhushan Damodaran, Remi Flamary, Nicolas Courty, Antoine Rolet, and Mathieu Blondel. Large-scale optimal transport and mapping estimation. In *ICLR 2018-International Conference on Learning Representations*, pp. 1–15, 2018.
- Fahad Shamshad and Ali Ahmed. Compressed sensing-based robust phase retrieval via deep generative priors. *IEEE Sensors Journal*, 21(2):2286–2298, 2020.
- Yujun Shen, Jinjin Gu, Xiaoou Tang, and Bolei Zhou. Interpreting the latent space of gans for semantic face editing. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2020.

## A Technical Results

We first recall the proof of the envelope theorem (Oyama & Takenawa, 2018, Prop. A.1).

**Theorem 7** (Envelope theorem Oyama & Takenawa (2018)). *Let  $X$  be a topological space. Let  $A$  be an open set of a normed vector space  $E$ . Let  $f : X \times A \rightarrow \mathbb{R}$  be a function and let us denote*

$$\forall a \in A, \quad v(a) = \sup_{x \in X} f(x, a). \quad (106)$$

*Let  $s : A \rightarrow X$  be such that for all  $a \in A$ ,  $v(a) = f(s(a), a)$ . Let  $\alpha \in A$  be a point such that*

- *$s$  is continuous at  $\alpha$ ,*
- *the partial differential  $D_a f$  of  $f$  with respect to  $a$  exists in a neighborhood of  $(s(\alpha), \alpha)$ , and is continuous at  $(s(\alpha), \alpha)$ .*

Then  $v$  is differentiable at  $\alpha$  and  $D_av(\alpha) = D_af(s(\alpha), \alpha)$ .

*Proof.* Let  $\xi = s(\alpha)$  and let  $\varepsilon > 0$ . The second hypothesis gives an open neighborhood  $U \times V$  of  $(\xi, \alpha)$  in  $X \times A$  such that for any  $(x, a) \in U \times V$ ,  $f(x, \cdot)$  is differentiable at  $a$  and such that the partial differential  $(x, a) \mapsto D_af(x, a)$  is continuous on  $U \times V$ . By continuity of  $s$ ,  $s^{-1}(U) \cap V$  is an open neighborhood of  $\alpha$  and thus it exists  $\eta > 0$  such that  $\|h\| < \eta$  implies  $\alpha + h \in V$  and  $s(\alpha + h) \in U$ .

By definition of  $v$ , we have for any  $h \in E$  such that  $\|h\| < \eta$ ,

$$f(\xi, \alpha + h) - f(\xi, \alpha) \leq v(\alpha + h) - v(\alpha) \leq f(s(\alpha + h), \alpha + h) - f(s(\alpha + h), \alpha). \quad (107)$$

On the one hand, by definition of  $D_af(\xi, \alpha)$ , there exists  $\eta_1 \in (0, \eta)$  such that  $\|h\| < \eta_1$  implies

$$|f(\xi, \alpha + h) - f(\xi, \alpha) - D_af(\xi, \alpha)h| \leq \varepsilon \|h\|. \quad (108)$$

On the other hand, for  $\|h\| < \eta$ ,  $t \in [0, 1] \mapsto f(s(\alpha + h), \alpha + th)$  is differentiable on  $[0, 1]$  and therefore, there exists  $\theta_h \in (0, 1)$  such that

$$f(s(\alpha + h), \alpha + h) - f(s(\alpha + h), \alpha) = D_af(s(\alpha + h), \alpha + \theta_h h)h. \quad (109)$$

By continuity of  $D_af$ , there is an open neighborhood  $\bar{U} \times \bar{V} \subset U \times V$  such that

$$\forall (x, a) \in U \times V, \quad \|D_af(x, a) - D_af(\xi, \alpha)\| \leq \varepsilon, \quad (110)$$

where  $\|\cdot\|$  denotes the dual norm. Again, by continuity of  $s$ ,  $s^{-1}(\bar{U}) \cap \bar{V}$  is an open neighborhood of  $\alpha$  and thus, there exists  $\eta_2 \in (0, \eta)$  such that  $\|h\| < \eta_2$  implies  $\alpha + h \in s^{-1}(\bar{U}) \cap \bar{V}$ . Therefore, for  $\|h\| < \eta_2$ ,  $(s(\alpha + h), \alpha + \theta_h h) \in \bar{U} \times \bar{V}$  and thus

$$|f(s(\alpha + h), \alpha + h) - f(s(\alpha + h), \alpha) - D_af(\xi, \alpha)h| \quad (111)$$

$$\leq \|D_af(s(\alpha + h), \alpha + \theta_h h) - D_af(\xi, \alpha)\| \|h\| \leq \varepsilon \|h\|. \quad (112)$$

Finally, for  $\|h\| < \min(\eta_1, \eta_2)$  we get

$$|v(\alpha + h) - v(\alpha) - D_af(\xi, \alpha)h| \leq \varepsilon \|h\|, \quad (113)$$

which proves that  $v$  is differentiable at  $\alpha$  and  $D_av(\alpha) = D_af(\xi, \alpha)$ .  $\square$

The next proposition gives the support of a push-forward distribution.

**Proposition 3.** *Let  $Q = [-1, 1]^s$ , and  $g : Q \rightarrow \mathbb{R}^d$  continuous. Let  $Z$  be a random variable with uniform distribution  $\zeta$  on  $Q$  and let  $\mu = g\#\zeta$  be the distribution of  $g(Z)$ . Then the support of  $\mu$  is exactly  $g(Q)$ .*

*Proof.* Since  $g$  is continuous,  $g(Q)$  is compact and in particular closed. Thus  $U = g(Q)^c$  is open, and one has that

$$\mu_\theta(U) = \mathbb{P}(g(Z) \in U) = \zeta(g^{-1}(U)) = 0, \quad (114)$$

because  $g^{-1}(U)$  does not intersect  $Q$ . This proves that  $\text{Supp}(\mu) \subset g_\theta(Q)$ . Now, if  $V$  is an open set such that  $\mu(V) = 0$ , then  $\mathbb{P}(Z \in g_\theta^{-1}(V)) = 0$ , which gives  $g_\theta^{-1}(V) \cap Q = \emptyset$  because  $g_\theta^{-1}(V)$  is open. It follows that  $V \subset g_\theta(Q)^c$ , which proves that  $\text{Supp}(\mu)$  is exactly  $g_\theta(Q)$ .  $\square$