Dávid Terjék¹

Abstract

Variational representations of f-divergences are central to many machine learning algorithms, with Lipschitz constrained variants recently gaining attention. Inspired by this, we define the Moreau-Yosida approximation of *f*-divergences with respect to the Wasserstein-1 metric. The corresponding variational formulas provide a generalization of a number of recent results, novel special cases of interest and a relaxation of the hard Lipschitz constraint. Additionally, we prove that the so-called tight variational representation of fdivergences can be to be taken over the quotient space of Lipschitz functions, and give a characterization of functions achieving the supremum in the variational representation. On the practical side, we propose an algorithm to calculate the tight convex conjugate of *f*-divergences compatible with automatic differentiation frameworks. As an application of our results, we propose the Moreau-Yosida f-GAN, providing an implementation of the variational formulas for the Kullback-Leibler, reverse Kullback-Leibler, χ^2 , reverse χ^2 , squared Hellinger, Jensen-Shannon, Jeffreys, triangular discrimination and total variation divergences as GANs trained on CIFAR-10, leading to competitive results and a simple solution to the problem of uniqueness of the optimal critic.

1. Introduction

Variational representations of divergences between probability measures are central to many machine learning algorithms, such as generative adversarial networks (Nowozin et al., 2016), mutual information estimation (Belghazi et al., 2018) and maximization (Hjelm et al., 2019), and energybased models (Arbel et al., 2021). One class of such measures is the family of f-divergences (Csiszár, 1963; Ali & Silvey, 1966; Csiszár, 1967), generalizing the well-known Kullback-Leibler divergence from information theory. Another is the family of optimal transport distances (Villani, 2008), including the Wasserstein-1 metric. In general, variational representations are supremums of integral formulas taken over sets of functions, such as the Donsker-Varadhan formula (Donsker & Varadhan, 1976) for the Kullback-Leibler divergence or the Kantorovich-Rubinstein formula (Villani, 2008) for the Wasserstein-1 metric. Informally speaking, one can implement (Nowozin et al., 2016; Arjovsky et al., 2017) such a formula by constructing a realvalued neural network called the critic (or discriminator) taking samples from the two probability measures as inputs, which is then trained to maximize the integral formula in order to approximate the supremum, resulting in a learned proxy to the actual divergence of said probability measures. Implementing the Kantorovich-Rubinstein formula in such a way involves restricting the Lipschitz constant of the neural network (Gulrajani et al., 2017; Petzka et al., 2018; Miyato et al., 2018; Adler & Lunz, 2018; Terjék, 2020), which effectively stabilizes the approximation procedure. Recently, Lipschitz regularization has been incorporated (Farnia & Tse, 2018; Zhou et al., 2019; Ozair et al., 2019; Song & Ermon, 2020; Arbel et al., 2021; Birrell et al., 2020) into learning algorithms based on variational formulas of divergences other than the Wasserstein-1 metric, leading to the same empirical effect and a number of theoretical benefits.

Inspired by this, we study Lipschitz-constrained variational representations of f-divergences. We show that existing instances of such variants are special cases of the Moreau-Yosida approximation of f-divergences with respect to the Wasserstein-1 metric. To any divergence and pair of probability measures corresponds a set of optimal critics, which are exactly those functions which achieve the supremum in the variational representation. An optimal critic corresponding to f-divergences is not Lipschitz in general (not even continuous). Since any function represented by a neural network is Lipschitz, when a neural network is trained to approximate such a divergence, its "target", an optimal critic, will never be reached. We show that when the divergence is replaced by its Moreau-Yosida approximation, the corresponding optimal critics are all Lipschitz continuous with uniformly bounded Lipschitz constants, leading to a divergence which is easier to approximate in practice via neural networks. The approximation is parametrized by a

¹Alfréd Rényi Institute of Mathematics, Budapest, Hungary. Correspondence to: Dávid Terjék <dterjek@renyi.hu>.

Proceedings of the 38th International Conference on Machine Learning, PMLR 139, 2021. Copyright 2021 by the author(s).

pair of real numbers, one of which controls the sharpness of the approximation and the Lipschitz constant of optimal critics. The other controls the behavior of the approximation such that a special case induces a hard Lipschitz constraint in the variational representation, and other values induce only a Lipschitz penalty term. While instances of the former already appeared in the literature, the latter is novel to our paper. A special case reduces to a novel, unconstrained variational representation of the Wasserstein-1 metric.

In order to prove these results, we first generalize the socalled tight variational representation of f-divergences to be taken over the space of Lipschitz functions or its quotient space, which is the subspace of functions vanishing at an arbitrary, fixed point. The latter leads to optimal critics being unique, having practical benefits. We additionally characterize the functions achieving the supremum in the variational representation. To apply the results, we propose an algorithm compatible with automatic differentiation frameworks to calculate the tight convex conjugate of f-divergences which in most cases does not admit a closed form, using Newton's method in the forward pass and implicit differentiation in the backward pass.

Finally, to demonstrate the usefulness of our results, we propose the Moreau-Yosida f-GAN, and implement it for the task of generative modeling on CIFAR-10. The experiments show that it is beneficial to use the Moreau-Yosida approximation as a proxy for f-divergences, the novel cases of which often outperform the ones with the hard Lipschitz constraint. On the other hand, the representation over the quotient space leads to a simple solution for the problem of uniqueness of the optimal critic.

To summarize, our contributions are

- a generalization of the tight variational representation of *f*-divergences between probability measures on compact metric spaces along with a characterization of functions achieving the supremum,
- a practical algorithm to calculate the tight convex conjugate of *f*-divergences compatible with automatic differentiation frameworks,
- variational formulas for the Moreau-Yosida approximation of *f*-divergences with respect to the Wasserstein-1 metric, including a relaxation of the hard Lipshcitz constraint and an unconstrained variational representation of the Wasserstein-1 metric, and
- the Moreau-Yosida *f*-GAN implementing the variational formulas for the Kullback-Leibler, reverse Kullback-Leibler, χ^2 , reverse χ^2 , squared Hellinger, Jensen-Shannon, Jeffreys, triangular discrimination and total variation divergences as GANs trained on CIFAR-10, leading to competitive performance.

2. Preliminaries

2.1. Notations

Denote the extended reals $\overline{\mathbb{R}} = \mathbb{R} \cup \{\pm \infty\}$, the nonnegative reals \mathbb{R}_+ , the extended nonnegative reals $\overline{\mathbb{R}}_+ = \mathbb{R}_+ \cup \infty$. The indicator of a set A is denoted i_A with $i_A(x) = 0$ if $x \in A$ and $i_A(x) = \infty$ otherwise. Absolute continuity and singularity of measures is denoted \ll and \bot , the Radon-Nikodym derivative of a measure μ with respect to a nonnegative measure ν such that $\mu \ll \nu$ by $\frac{d\mu}{d\nu}$, the support of a measure μ by $\sup p(\mu)$, a property to hold almost everywhere with respect to a measure μ by μ -a.e. The relative interior of a subset A of a vector space is denoted relint A, which for subsets of \mathbb{R} only differs from the interior for singletons whose relative interior is the singleton itself.

2.2. Convex analysis (Zalinescu, 2002)

Given a topological vector space X, denote its topological dual by X^* , i.e. the set of real-valued continuous linear maps on X, which is a topological vector space itself, and the canonical pairing by $\langle \cdot, \cdot \rangle : X \times X^* \to \mathbb{R}$, which is the continuous bilinear map $((x, x^*) \rightarrow \langle x, x^* \rangle =$ $x^*(x)$). Given a function $f: X \to \overline{\mathbb{R}}$, the set dom f = $\{x \in X : f(x) < \infty\}$ is the effective domain of f. A function f is proper if dom $f \neq \emptyset$ and $f(x) > -\infty$ for all $x \in X$, otherwise it is improper. For a convex function $f : X \to \overline{\mathbb{R}}$, its convex conjugate is $f^* : X^* \to \overline{\mathbb{R}}$ defined by $f^*(x^*) = \sup_{x \in X} \{ \langle x, x^* \rangle - f(x) \}$, and its subdifferential at $x \in X$ is the set $\partial f(x) = \{x^* \in$ $X^* \mid \forall \hat{x} \in X : \langle \hat{x} - x, x^* \rangle \leq f(\hat{x}) - f(x) \}$. The biconjugate f^{**} of f is the conjugate of its conjugate f^* , i.e. $f^{**}(x) = \sup_{x^* \in X^*} \{ \langle x, x^* \rangle - f^*(x^*) \}$, which is equivalent to f if f is proper, convex and lower semicontinuous. In that case, the supremum of the biconjugate representation is achieved precisely at elements of $\partial f(x)$. Conversely, the supremum in the conjugate representation of $f^*(x^*)$ is achieved at elements of $\partial f^*(x^*) = \{x \in X \mid \forall \hat{x}^* \in X^* :$ $\langle x, \hat{x}^* - x^* \rangle \le f^*(\hat{x}^*) - f^*(x^*) \}.$

2.3. f-divergences

Given a proper, convex and lower semicontinuous function¹ $\phi : \mathbb{R} \to \overline{\mathbb{R}}$, a measure μ and a nonnegative measure ν on a measurable space X, the f-divergence $D_{\phi}(\mu \| \nu)$ of μ from ν is defined (Csiszár, 1963; Ali & Silvey, 1966; Csiszár, 1967; Borwein & Lewis, 1993; Csiszár et al., 1999) as

$$\int \phi \circ \frac{d\mu_c}{d\nu} d\nu + \phi'(\infty)\mu_s^+(X) - \phi'(-\infty)\mu_s^-(X).$$
 (1)

Here, $\mu_c \ll \nu, \mu_s \perp \nu$ are the absolutely continuous and singular parts of the Lebesgue decomposition of μ with

¹Originally, f is used in place of ϕ (hence the name), but we reserve the symbol f for other functions.

respect to ν , $\mu_s^+, \mu_s^- \ge 0$ is the Jordan decomposition of the singular part, and $\phi'(\pm \infty) = \lim_{x \to \pm \infty} \frac{\phi(x)}{x} \in \overline{\mathbb{R}}$. The well-known variational representation

$$D_{\phi}(\mu \| \nu) = \sup_{f: X \to \mathbb{R}} \left\{ \int f d\mu - \int \phi^* \circ f d\nu \right\}$$
(2)

can be obtained as the biconjugate of the mapping $(\mu \rightarrow D_{\phi}(\mu \| \nu))$. The so-called tight variational representation

$$D_{\phi}(\mu \| \nu) = \sup_{f: X \to \mathbb{R}} \left\{ \int f d\mu - \inf_{\sup f(X) - \phi'(\infty) \le \gamma} \left\{ \int \phi_{+}^{*} \circ (f - \gamma) d\nu + \gamma \right\} \right\}$$
(3)

with $\phi_+ = \phi + i_{\mathbb{R}_+}$ was obtained in Agrawal & Horel (2020) as the biconjugate of the mapping $(\mu \to D_{\phi}(\mu \| \nu) + i_{P(X)}(\mu))$ (already considered in Ruderman et al. (2012)), and is valid for pairs of probability measures μ, ν .

2.4. Wasserstein-1 distance (Villani, 2008)

Given probability measures μ, ν on a metric space (X, d), the Wasserstein-1 distance of μ and ν is defined as

$$W_1(\mu,\nu) = \inf_{\pi \in \Pi(\mu,\nu)} \int d(x_1,x_2) d\pi(x_1,x_2), \quad (4)$$

where $\Pi(\mu, \nu)$ is the set of probability measures supported on the product space $X \times X$ with marginals μ and ν . It has a well-known variational representation called the Kantorovich-Rubinstein formula

$$W_1(\mu,\nu) = \sup_{\|f\|_L \le 1} \left\{ \int f d\mu - \int f d\nu \right\}, \qquad (5)$$

where

$$||f||_{L} = \sup_{x,y \in X, x \neq y} \left\{ \frac{|f(x) - f(y)|}{d(x,y)} \right\}$$
(6)

is the Lipschitz norm of f. The supremum is achieved by the so-called Kantorovich potentials $f : X \to \mathbb{R}$, unique μ, ν -a.e. up to an additive constant.

2.5. Moreau-Yosida approximation

Let (X, d) be a metric space and $f : X \to \mathbb{R}$ a proper function, and $0 < \lambda, \alpha \in \mathbb{R}$ constants. The Moreau-Yosida approximation of index λ and order α of f is defined (Jost & Li-Jost, 2008; Dal Maso, 1993) as

$$f_{\lambda,\alpha}(x) = \inf_{y \in X} \{ f(y) + \lambda d(x,y)^{\alpha} \}.$$
 (7)

It holds that $\overline{f}(x) = \sup_{\lambda>0} f_{\lambda,\alpha}(x) = \lim_{\lambda\to\infty} f_{\lambda,\alpha}(x)$, where \overline{f} is the greatest lower semicontinuous function with $\overline{f} \leq f$.

3. Lipschitz representation of *f*-divergences

In this work, we consider the set P(X) of probability measures on a compact metric space (X, d), which is itself a compact metric space with the metric W_1 , metrizing the weak convergence of measures. We prove that the tight variational representation of D_{ϕ} from Agrawal & Horel (2020) can be generalized in the sense that the supremum can be taken over the set $Lip(X, x_0)$ of Lipschitz continuous functions on X that vanish at an arbitrary base point $x_0 \in X$. This is a strictly smaller set than the set of bounded and measurable functions over which the supremum was taken originally. To apply convex analytic techniques, we consider the duality between vector spaces of measures and Lipschitz functions. This aspect is detailed in Appendix 8.1. An important property of the choice of vector spaces is that the topology on the space of measures generalizes the usual weak convergence of probability measures (Hanin, 1999). Proofs and more precise statements of our propositions can be found in Appendix 8.2.

Proposition 1. Given probability measures $\mu, \nu \in P(X)$ and a proper, convex and lower semicontinuous function $\phi : \mathbb{R} \to \overline{\mathbb{R}}$ strictly convex at 1 with $\phi(1) = 0$ and $1 \in \text{relint dom } \phi$, the f-divergence D_{ϕ} has the equivalent variational representation

$$D_{\phi}(\mu \| \nu) = \sup_{f \in Lip(X)} \left\{ \int f d\mu - D_{\phi}^{*}(f \| \nu) \right\}$$
$$= \sup_{f \in Lip(X, x_{0})} \left\{ \int f d\mu - D_{\phi}^{*}(f \| \nu) \right\}, \quad (8)$$

with the tight convex conjugate $D^*_{\phi}(\cdot \| \nu) : Lip(X) \to \mathbb{R}$ being

$$D_{\phi}^{*}(f||\nu) = \sup_{\mu \in P(X)} \left\{ \int f d\mu - D_{\phi}(\mu||\nu) \right\}$$
$$= \min_{\sup f(X) - \phi'(\infty) \le \gamma} \left\{ \int \phi_{+}^{*} \circ (f - \gamma) d\nu + \gamma \right\}.$$
(9)

The conjugate $D_{\phi}^{*}(\cdot \|\nu)$ is a topical function (Mohebi, 2005), meaning that $D_{\phi}^{*}(f + C \|\nu) = D_{\phi}^{*}(f \|\nu) + C$ and $D_{\phi}^{*}(f_{1}\|\nu) \geq D_{\phi}^{*}(f_{2}\|\nu)$ both hold for $\forall C \in \mathbb{R}$ and $f_{1} \geq f_{2}$. Based on the constant additivity property, the substitution $D_{\phi}^{*}(f \|\nu) = \int f d\nu + D_{\phi}^{*}(f - \int f d\nu \|\nu)$ leads to

$$\sup_{f\in Lip(X)} \left\{ \int f d\mu - \int f d\nu - D_{\phi}^* \left(f - \int f d\nu \|\nu\right) \right\},\,$$

reinterpreting the variational representation of $D_{\phi}(\mu \| \nu)$ as a penalized variant of maximum mean deviation. A closed form expression for $D_{\phi}^{*}(\cdot \| \nu)$ is available for the Kullback-Leibler divergence with $D_{KL}^{*}(f \| \nu) = \log \int e^{f} d\nu$.

We call functions f_* for which $D_{\phi}(\mu \| \nu) = \int f_* d\mu - D_{\phi}^*(f_* \| \nu)$ holds, i.e. $f_* \in \partial D_{\phi}(\mu \| \nu)$, *Csiszár potentials* of

 μ, ν . This is in analogy with Kantorovich potentials, which are similarly unique μ, ν -a.e. up to an additive constant. In the second variational representation in (8), the additive constant is unique since $f(x_0) = 0$ must hold. The following result is built on Borwein & Lewis (1993, Theorem 2.10).

Proposition 2. Given probability measures $\mu, \nu \in P(X)$, a function $f_* \in Lip(X)$ is a Csiszár potential of μ, ν , i.e. $D_{\phi}(\mu \| \nu) = \int f_* d\mu - D_{\phi}^*(f_* \| \nu)$, if and only if there exists $C \in \mathbb{R}$ such that the conditions

$$\sup f_*(X) + C \le \phi'(\infty), \tag{10}$$

$$\frac{d\mu_c}{d\nu}(x) \in \partial \phi_+^*(f_*(x) + C) \ \nu\text{-a.e.}$$
(11)

and

$$\operatorname{supp}(\mu_s) \subset \{x \in X : f_*(x) + C = \phi'(\infty)\}$$
(12)

hold. Such f_* are unique μ, ν -a.e. up to an additive constant.

If ϕ is of Legendre type (Borwein & Lewis, 1993), then ϕ_+ and ϕ_+^* are both continuously differentiable on int dom ϕ_+ and int dom ϕ_+^* , respectively, while $\phi_+^{*'}$ is increasing, and invertible where its value is positive with its inverse given by the strictly increasing ϕ'_+ . With these, the second condition is equivalent to

$$f_*(x) + C = \phi'_+\left(\frac{d\mu_c}{d\nu}(x)\right) \ \mu_c$$
-a.e. (13)

Informally, this means that f_* is the strictly increasing image of the likelihood ratio. One can then deduce from the Neyman-Pearson lemma (Reid & Williamson, 2011) that for the binary experiment of discriminating samples from μ and ν , the statistical test $(x \to \chi_{[\tau,\infty]}(f_*(x)))$ is a most powerful test for any threshold $\tau \in \mathbb{R}$.

Conversely to the above proposition, given $\nu \in P(X)$ and $f \in Lip(X)$, the same conditions characterize the set of $\mu_* \in P(X)$ for which the supremum is achieved in the conjugate representation of $D_{\phi}^*(\cdot \| \nu)$, i.e. $\mu_* \in$ $\partial D_{\phi}^*(f \| \nu)$. Denoting the optimal γ in (9) by $\gamma_{\phi,\nu}(f)$, for any $\mu_* \in P(X)$ satisfying the conditions in Proposition 2 with $C = -\gamma_{\phi,\nu}(f)$ one has $\mu_* \in \partial D_{\phi}^*(f \| \nu)$. For the Kullback-Leibler divergence, this reduces to the softmax $\mu_* = \frac{1}{\int e^f d\nu} e^f \cdot \nu$. In case X is a finite set, this leads to a family of prediction functions obtained as gradients of $D_{\phi}^*(f \| \nu)$ (Blondel et al., 2020).

We propose an algorithm for the practical evaluation of $D_{\phi}^{*}(\cdot \|\nu)$ when no closed form expression is available in the case when the support of ν is finite² and ϕ is such that ϕ_{+}^{*} is twice differentiable on int dom ϕ_{+}^{*} with non-vanishing

Algorithm 1 Calculate $\gamma_{\phi,\nu}(f)$ and $\nabla_f \gamma_{\phi,\nu}(f)$						
Input: $f, \nu \in \mathbb{R}^n, \phi : \mathbb{R} \to \overline{\mathbb{R}}, 0 < \epsilon, \tau \in \mathbb{R}$						
if $\phi'(\infty) < \infty$ then						
$\gamma = \max(f) - \phi'(\infty) + \epsilon.$						
else						
$\gamma = \langle u, f angle$						
end if						
repeat						
$s = \frac{-\langle \nu, (\phi_+^*)'(f-\gamma) \rangle + 1}{\langle \nu, (\phi_+^*)''(f-\gamma) \rangle}$						
$\gamma = \gamma - s$						
until $ s < \tau$						
$\nabla_f \gamma = \frac{\nu \odot (\phi_+^*)''(f-\gamma)}{\langle \nu, (\phi_+^*)''(f-\gamma) \rangle}$						

second derivative. Assuming that f achieves its maximum on the support of ν and that γ achieving the minimum is unique, finding γ reduces to a finite dimensional problem, i.e. f, ν can be considered as elements of \mathbb{R}^n with n being the number of elements of the support of ν . Based on Newton's method and the implicit function theorem, we propose Algorithm 1 to calculate $\gamma_{\phi,\nu}(f)$ and its gradient³. Then, the conjugate can be calculated as

$$D^*_{\phi}(f \| \nu) = \langle \nu, \phi^*_+(f - \gamma_{\phi,\nu}(f)) \rangle + \gamma_{\phi,\nu}(f).$$
 (14)

The derivation of the algorithm can be found in Appendix 8.3, along with the corresponding functions ϕ_+ , ϕ_+^* and their derivatives for the Kullback-Leibler, reverse Kullback-Leibler, χ^2 , reverse χ^2 , squared Hellinger, Jensen-Shannon, Jeffreys and triangular discrimination divergences. For the Kullback-Leibler divergence, one has the closed form $\gamma_{\phi,\nu}(f) = \log \int e^f d\nu$.

We found that exploiting the constant additivity property by calculating the conjugate as

$$D_{\phi}^{*}(f\|\nu) = D_{\phi}^{*}(f - \max(f)\|\nu) + \max(f)$$
(15)

is beneficial to avoid numerical instabilities. This can be seen as a generalization of the log-sum-exp trick.

4. Moreau-Yosida approximation of *f*-divergences

Since the mapping $(\mu \to D_{\phi}(\mu \| \nu))$ from the metric space $(P(X), W_1)$ to \mathbb{R} is proper and lower semicontinuous, it is an ideal candidate for Moreau-Yosida approximation, for which the infimum is always achieved since $(P(X), W_1)$ is compact if (X, d) is. Given $0 < \lambda, \alpha \in \mathbb{R}$, the Moreau-Yosida approximation of index λ and order α of $D_{\phi}(\cdot \| \nu)$ with respect to W_1 is therefore defined as

$$D_{\phi,\lambda,\alpha}(\mu \| \nu) = \min_{\xi \in P(X)} \{ D_{\phi}(\xi \| \nu) + \lambda W_1(\mu,\xi)^{\alpha} \}.$$
(16)

²Such measures are dense in $(P(X), W_1)$.

 $^{{}^3\}langle\cdot,\cdot\rangle$ and \odot denote the standard dot product and the element-wise product in $\mathbb{R}^n.$

This is still a divergence in the sense that $D_{\phi,\lambda,\alpha}(\mu \| \nu) \ge 0$ with equality if and only if $\mu = \nu$. The original divergence can be recovered as $D_{\phi}(\mu \| \nu) = \sup_{\lambda>0} D_{\phi,\lambda,\alpha}(\mu \| \nu) =$ $\lim_{\lambda\to\infty} D_{\phi,\lambda,\alpha}(\mu \| \nu)$ for any $\alpha > 0$. Moreover, for $\alpha \ge 1$, $D_{\phi,\lambda,\alpha}(\cdot \| \nu)$ is Lipschitz continuous with respect to W_1 . If $\alpha = 1$, the Lipschitz constant is exactly λ . In some cases⁴, variational representations are available.

Proposition 3. Given probability measures $\mu, \nu \in P(X)$, $\lambda > 0, \alpha \ge 1$ and a proper, convex and lower semicontinuous function $\phi : \mathbb{R} \to \overline{\mathbb{R}}$ strictly convex at 1 with $\phi(1) = 0$ and $1 \in \text{relint dom } \phi$, the divergence $D_{\phi,\lambda,\alpha}(\mu \| \nu)$ has the equivalent variational representation

$$\max_{f \in Lip(X,x_0), \|f\|_L \le \lambda} \left\{ \int f d\mu - D^*_{\phi}(f\|\nu) \right\}$$
(17)

if $\alpha = 1$, and

$$\max_{f \in Lip(X,x_0)} \left\{ \int f d\mu - D_{\phi}^*(f \| \nu) - (\alpha - 1)\alpha^{\frac{\alpha}{1-\alpha}} \lambda^{\frac{1}{1-\alpha}} \| f \|_L^{\frac{\alpha}{\alpha-1}} \right\}$$
(18)

if $\alpha > 1$.

In the limit $\alpha \to 1$, (18) converges to (17) in the sense that $\lim_{\alpha\to 1} (\alpha - 1)\alpha^{\frac{\alpha}{1-\alpha}}\lambda^{\frac{1}{1-\alpha}} ||f||_L^{\frac{\alpha}{\alpha-1}} = 0$ if $||f||_L \leq \lambda$ and ∞ otherwise, providing an unconstrained relaxation of the hard constraint $||f||_L \leq \lambda$.

Choosing $\phi = i_{\{1\}}$ (so that $D_{\phi}(\cdot \|\nu) = i_{\{\nu\}}$ and $D_{\phi}^*(f\|\nu) = \int f d\nu$), one has $D_{\phi,\lambda,\alpha}(\mu\|\nu) = \lambda W_1(\mu,\nu)^{\alpha}$, leading to the following unconstrained variational representation of W_1 .

Proposition 4. Given $\mu, \nu \in P(X)$, $\lambda > 0$ and $\alpha > 1$, $W_1(\mu, \nu)$ has the equivalent unconstrained variational representation

$$\left(\frac{1}{\lambda}\max_{f\in Lip(X,x_0)}\left\{\int fd\mu - \int fd\nu - \left(\alpha - 1\right)\alpha^{\frac{\alpha}{1-\alpha}}\lambda^{\frac{1}{1-\alpha}} \|f\|_L^{\frac{\alpha}{\alpha-1}}\right\}\right)^{\frac{1}{\alpha}}.$$
 (19)

The maximum is achieved at $\alpha \lambda W_1(\mu, \nu)^{\alpha-1} f_*$, with f_* being a Kantorovich potential of μ, ν .

As stated, subgradients of the mapping $(\mu \to \lambda W_1(\mu, \nu)^{\alpha})$ are nothing but the Kantorovich potentials f_* achieving the supremum in the Kantorovich-Rubinstein formula, scaled by the coefficient $\alpha \lambda W_1(\mu, \nu)^{\alpha-1}$. This allows the characterization of subgradients of the mapping $(\mu \to D_{\phi,\lambda,\alpha}(\mu \| \nu))$. **Proposition 5.** Given probability measures $\mu, \nu \in P(X)$, $\lambda > 0, \alpha \ge 1$ and a proper, convex and lower semicontinuous function $\phi : \mathbb{R} \to \mathbb{R}$ strictly convex at 1 with $\phi(1) = 0$ and $1 \in$ relint dom ϕ , let $\xi_* \in P(X)$ be a probability measure achieving the minimum in (16), i.e. for which $D_{\phi,\lambda,\alpha}(\mu \| \nu) = D_{\phi}(\xi_* \| \nu) + \lambda W_1(\mu, \xi_*)^{\alpha}$ holds. Then there exists an $f_* \in Lip(X)$ achieving the maximum in (17) if $\alpha = 1$ or (18) if $\alpha > 1$, which is a Csiszár potential of ξ_*, ν and $\alpha \lambda W_1(\mu, \xi_*)^{\alpha-1}$ times a Kantorovich potential of μ, ξ_* at the same time.

These imply that for any $\tau \in \mathbb{R}$, the mapping $(x \to \chi_{[\tau,\infty]}(f_*(x)))$ is a most powerful test for discriminating samples from ξ_* and ν , and that $||f_*||_L = \alpha \lambda W_1(\mu, \xi_*)^{\alpha-1}$. Informally, since ξ_* is close to μ in W_1 , the above mapping can be seen as a Lipschitz regularized version of a most powerful test for discriminating μ and ν .



Figure 1. Multiplier and exponent of $||f||_L$

Consider the reparametrization $\lambda = \frac{1}{\alpha}\beta^{-\alpha}$, so that (18) reduces to

$$D_{\phi,\frac{1}{\alpha}\beta^{-\alpha},\alpha}(\mu\|\nu) = \max_{f\in Lip(X,x_0)} \left\{ \int f d\mu - D_{\phi}^*(f\|\nu) - \frac{\alpha-1}{\alpha} \left(\beta\|f\|_L\right)^{\frac{\alpha}{\alpha-1}} \right\}.$$
 (20)

A plot of the respective values of the multiplier and the exponent for $\beta = 1$ and $\alpha \in [1, 16]$ are visualized in Figure 1. In the limit $\alpha \to \infty$, the multiplier and exponent both converge to 1. On the other hand, one has $\lim_{\alpha \to 1} \frac{\alpha - 1}{\alpha} ||f||_L^{\alpha} = 0$ if $||f||_L \leq 1$ and ∞ otherwise.

An interesting special case is the limit $\alpha \to \infty$, resulting in the minimum of $D_{\phi}(\xi \| \nu)$ with $\xi \in P(X)$ ranging over the

⁴Since the mapping $(\xi \to \lambda W_1(\mu, \xi)^{\alpha})$ is neither convex nor concave if $0 < \alpha < 1$, we could not obtain a variational representation via Fenchel-Rockafellar duality in this case.

Wasserstein-1 ball of radius β centered at μ as

$$\lim_{\alpha \to \infty} D_{\phi, \frac{1}{\alpha} \beta^{-\alpha}, \alpha}(\mu \| \nu) = \min_{\xi \in P(X), W_1(\xi, \mu) \le \beta} \{ D_{\phi}(\xi \| \nu) \}$$

$$= \max_{f \in Lip(X,x_0)} \left\{ \int f d\mu - D_{\phi}^*(f \| \nu) - \beta \| f \|_L \right\}.$$
 (21)

This should be contrasted with (17) corresponding to the $\alpha = 1$ case, which also has a hard constraint, but in the dual formula.

Since the values of the above formulas for a given f are invariant for constant translations f+C, the supremums can equivalently be taken over Lip(X) instead of $Lip(X, x_0)$ in all cases.

5. Moreau-Yosida f-GAN

We propose the Moreau-Yosida *f*-GAN (MY*f*-GAN) as an implementation of the variational formula of the Moreau-Yosida regularization of D_{ϕ} with respect to W_1 . The function *f* in (17) or (18) is parametrized by a neural network called the critic, which is trained to maximize the formula inside the maximum, providing an approximation of the exact value of the divergence. One of the measures μ, ν is represented by the dataset, and the other by a neural network called the generator. The generator transforms samples from a fixed noise distribution into ones resembling the data distribution, and is trained to minimize the divergence approximated by the critic.

Based on the reparametrized formula (20) with the substitution $D_{\phi}^{*}(f||\nu) = \int f d\nu + D_{\phi}^{*} (f - \int f d\nu ||\nu)$, the two minimax games are the following. First let μ be the generated distribution and ν be the data, resulting in the forward (\rightarrow) formulation

$$\min_{\theta_g \in \mathbb{R}^{l}} \max_{\theta_f \in \mathbb{R}^{k}} \mathbb{E}_{(\zeta_n, \nu_n) \sim (P_z, P_d)} \langle g_{\theta_g \#} \zeta_n, f_{\theta_f} \rangle - \langle \nu_n, f_{\theta_f} \rangle
- D_{\phi}^* (f_{\theta_f} - \langle \nu_n, f_{\theta_f} \rangle \| \nu_n)
- \frac{\alpha - 1}{\alpha} (\beta \| f_{\theta_f} \|_{L, g_{\hat{\theta}_g \#} \zeta_n, \nu_n})^{\frac{\alpha}{\alpha - 1}}. \quad (22)$$

Now let μ be the data and ν the generated distribution, leading to the reverse (\leftarrow) formulation

$$\min_{\theta_g \in \mathbb{R}^l} \max_{\theta_f \in \mathbb{R}^k} \mathbb{E}_{(\mu_n, \zeta_n) \sim (P_d, P_z)} \langle \mu_n, f_{\theta_f} \rangle - \langle g_{\theta_g \#} \zeta_n, f_{\theta_f} \rangle
- D_{\phi}^* (f_{\theta_f} - \langle g_{\hat{\theta}_g \#} \zeta_n, f_{\theta_f} \rangle \| g_{\hat{\theta}_g \#} \zeta_n)
- \frac{\alpha - 1}{\alpha} (\beta \| f_{\theta_f} \|_{L, \mu_n, g_{\hat{\theta}_g \#} \zeta_n})^{\frac{\alpha}{\alpha - 1}}. \quad (23)$$

The notation of the minimax games is the following. The functions $f : X \times \mathbb{R}^k \to \mathbb{R}$ and $g : Z \times \mathbb{R}^l \to X$ are the critic and generator neural networks parametrized by weight vectors $\theta_f \in \mathbb{R}^k$ and $\theta_g \in \mathbb{R}^l$, and $f_{\theta_f}, g_{\theta_g}$ are

shorthands for $f(\cdot, \theta_f), g(\cdot, \theta_g)$. The latent space is Z = \mathbb{R}^m . The sample space $X \subset \mathbb{R}^n$ is a compact subset of Euclidean space equipped with the restriction of the metric induced by the Euclidean norm, e.g. $X = [-1, 1]^{3*32*32}$ for CIFAR-10. $P_d \in P(X)$ denotes the data distribution and $P_z \in P(Z)$ the noise distribution, e.g. a standard normal. Empirical measures (corresponding to minibatches) are denoted $\mu_n \sim P$, meaning that $\mu_n = \frac{1}{n} \sum_{i=1}^n \delta_{x_{\mu,i}}$ with $(x_{\mu,i}) \subset X$ being a realization of a sequence of n independent and identical copies of the random variable corresponding to P. The empirical measure corresponding to the generated distribution is obtained as the pushforward $g_{\theta_a \#} \zeta_n$ of the latent empirical measure ζ_n (a minibatch of noise samples) through the generator g_{θ_g} . Empirical means are denoted $\langle \mu_n, f \rangle = \frac{1}{n} \sum_{i=1}^n f(x_{\mu,i})$. The conjugate D_{ϕ}^* is calculated according to (14) using the stabilization trick (15). By $\hat{\theta}_g$ we denote a copy of θ_g , meaning that θ_q is not optimized to minimize terms containing the copy, i.e. the loss function of the generator is $\pm \langle f_{\theta_f}, g_{\theta_a \#} \zeta_n \rangle$. The term $||f_{\theta_f}||_{L,\mu_n,\nu_n}$ denotes a possibly data-dependent estimate of $||f_{\theta_f}||_L$. The minimax games include the case $\lim_{\alpha \to \infty} \frac{\alpha - 1}{\alpha} = \lim_{\alpha \to \infty} \frac{\alpha}{\alpha - 1} = 1.$

Lipschitz norm estimation. Rademacher's theorem (Weaver, 2018) states that if $||f||_L < \infty$ for $f : \mathbb{R}^n \to \mathbb{R}$, then $||(x \to ||\nabla f(x)||_2)||_{\infty} = ||f||_L$ holds, i.e. that the supremum of the function mapping $x \in \mathbb{R}^n$ to the Euclidean norm of the gradient of f at x is equal to the Lipschitz norm of f. Based on this and the gradient penalty of Gulrajani et al. (2017), we propose for $||f_{\theta_f}||_{L,\mu_n,\nu_n}$ the estimator

$$\mathbb{E}_{\upsilon_n \sim \mathcal{U}[0,1)} \max_{x \in \operatorname{supp}(\upsilon_n \mu_n + (1 - \upsilon_n)\nu_n)} \|\nabla f_{\theta_f}(x)\|_2 \quad (24)$$

giving a lower bound to $||f_{\theta_f}||_L$. Here, $\mathcal{U}[0,1)$ is the uniform distribution on [0,1) from which an empirical measure $v_n = \frac{1}{n} \sum_{i=1}^n \delta_{u_i}$ is drawn, and $u_n \mu_n + (1-u)\nu_n = \frac{1}{n} \sum_{i=1}^n \delta_{u_i x \mu_i i} + (1-u_i)x_{\nu,i}$ denotes the corresponding interpolation of μ_n and ν_n . This clearly biased estimator leaves room for improvement. Constructing an unbiased estimator would require assuming a distribution for the random variable representing the value of the gradient norm of the critic, which we leave for future work.

Relaxation of hard Lipschitz constraint. We implement the hard constraint case $\alpha = 1$ by replacing the last term in the minimax games with the one-sided gradient penalty (Gul-rajani et al., 2017; Petzka et al., 2018) $\ell \mathbb{E}_{v_n \sim \mathcal{U}[0,1)} \langle v_n \mu_n + (1 - v_n)\nu_n, (\max\{0, \|\nabla f_{\theta_f}(\cdot)\|_2 - \beta^{-1}\})^2 \rangle$ with the coefficient $0 < \ell \in \mathbb{R}$ controlling the strength of the penalty. This is a widely used method to enforce the hard constraint $\|f_{\theta_f}\|_L \leq \beta^{-1}$. We visualize the maximum, mean and minimum of minibatches of gradient norms of the critic during training in Figure 2 for $\alpha = 1$ with the gradient penalty and $\alpha = 1.05$ with the estimator detailed above. The $\alpha = 1$ case does not enforce the hard constraint, since only the mean

		$\beta = 0$		$\alpha=1.05,\beta=1$		$\alpha=2,\beta=1$		$\alpha=\infty,\beta=0.5\rightarrow0.2$	
D_{ϕ}		IS	FID	IS	FID	IS	FID	IS	FID
KULLBACK-LEIBLER	\rightarrow	7.16	34.12	8.26	13.22	8.33	14.83		
	\leftarrow			8.20	13.85	8.09	13.42	8.20	12.51
REVERSE KULLBACK-LEIBLER	\rightarrow			8.33	12.97	8.30	13.27		
	\leftarrow			8.34	13.24	8.17	13.13	8.09	15.26
χ^2	\rightarrow			8.18	14.17	8.26	13.36		
	\leftarrow			8.37	13.36	8.23	12.95	8.27	13.46
reverse χ^2	\rightarrow			8.47	13.89	8.26	14.59		
	\leftarrow			8.24	14.04	8.45	12.28	8.11	14.17
SQUARED HELLINGER	\rightarrow			8.03	16.41	8.07	16.06		
	\leftarrow			8.25	15.89	8.25	13.93	8.52	12.18
JENSEN-SHANNON	\rightarrow	7.51	30.17	8.30	14.49	8.34	12.71		
	\leftarrow			8.04	16.04	8.37	11.57	8.27	12.58
JEFFREYS	\rightarrow			8.09	13.99	8.21	14.46		
	\leftarrow			8.25	13.32	8.34	13.04		
TRIANGULAR DISCRIMINATION	\rightarrow	6.45	43.14	8.42	13.54	8.08	14.68		
	\leftarrow			8.15	14.28	8.35	12.21	8.09	15.13
TOTAL VARIATION	\rightarrow	7.41	31.09	8.12	15.44	8.28	14.61		
	\leftarrow			8.08	13.53	8.12	13.77	8.05	14.60
TRIVIAL				8.07	15.97	8.04	14.75	6.67	36.48

Table 1. MY f-GAN performance on CIFAR-10



Figure 2. $\|\nabla f(X)\|_2$ with relaxed Lipschitz constraint

of the gradient norms is concentrated around $\beta^{-1} = 1$, and not their maximum. The $\alpha = 1.05$ case, being a relaxation of the hard constraint, empirically behaves very similarly to an ideal hard constraint implementation, in the sense that the maximum of the gradient norms is concentrated around $\beta^{-1} = 1$. This is no surprise in light of Proposition 5, since $||f_*||_L = \alpha \lambda W_1(\mu, \xi_*)^{\alpha-1} = \beta^{-\alpha} W_1(\mu, \xi_*)^{\alpha-1} =$ $\beta^{-(1+\epsilon)} W_1(\mu, \xi_*)^{\epsilon}$ is very close to β^{-1} in practice for small ϵ , such as $\epsilon = 0.05$. We did not observe significant performance differences. This particular experiment used $\ell = 10$ and ϕ corresponding to the Kullback-Leibler divergence, but we observed identical behavior in other hyperparameter settings as well with a range of α close to 1. We argue that using the relaxation with some $\alpha = 1 + \epsilon$ is potentially beneficial for other applications requiring the satisfaction of a hard Lipschitz constraint.

Choice of f-divergence. Quantitative results in terms of Inception Score (IS) and Fréchet Inception Distance (FID) can be seen in Table 1. Missing values in the unregularized case ($\beta = 0$) indicate divergent training, showing that regularization ($\beta > 0$) not only improves performance, but leads to convergent training even in cases when it does not seem possible without regularization. The TRIVIAL case indicates $D_{\phi}(\cdot \| \nu) = i_{\{\nu\}}$, so that the forward and reverse formulations are identical. In this case, $D_{\phi,\frac{1}{\alpha}\beta^{-\alpha},\alpha}(\mu\|\nu)$ reduces to $\frac{1}{\alpha}\beta^{-\alpha}W_1(\mu\|\nu)^{\alpha}$. If $\alpha > 1$, this leads to an unconstrained formulation of the Wasserstein GAN corresponding to Proposition 4. The original, constrained Wasserstein GAN with gradient penalty led to an IS of 8.09 and an FID of 13.40 in our implementation. This is marginally better than the performance of the unconstrained variant as reported in Table 1. As shown in Figure 2, gradient penalty leads to a higher gradient norm than required by the hard constraint, which might lead to the observed marginal performance improvement. Indeed, increasing β leads to better performance for the unconstrained variant, e.g. $\beta = 0.5$ with $\alpha = 2$ led to and IS of 8.14 and an FID of 13.33, which is in turn marginally better than the original, constrained variant. While it is hard to tell from these results which *f*-divergence is the best, it is definitely not the TRIVIAL one.



Figure 3. f(X) for default and quotient critic

Quotient critic. To ensure that $f_{\theta_f} \in Lip(X, x_0)$, we simply modify the forward pass of the critic to return $f_{\theta_f}(x) - f_{\theta_f}(x_0)$ instead of $f_{\theta_f}(x)$. This induces negligible computational overhead since $f_{\theta_f}(x_0)$ can be calculated with a minibatch of size 1, with the choice of x_0 being arbitrary, e.g. the zero vector in our implementation. We call this the quotient critic since $Lip(X, x_0)$ is isomorphic to the quotient space $\frac{Lip(X)}{\mathbb{R}}$. In Figure 3 we visualize the maximum, mean and minimum of the critic output over minibatches of generated samples during training. It is clear that the quotient critic solves the drifting of the output of the critic, which was found to hurt performance in some cases (Karras et al., 2018; Adler & Lunz, 2018). We observed only marginal performance improvement.

Loss function of the generator. The reason for picking the penalized mean deviation form of the variational formulas for this application is that in the reverse case, we found that using $-\langle g_{\theta_g \#} \zeta_n, f_{\theta_f} \rangle$ as the loss function of the generator leads to superior performance than using $-D_{\phi}^*(f_{\theta_f} || g_{\theta_g \#} \zeta_n)$, which cripples performance in most cases. This suggests that gradients of the Csiszár potential f_* might be of greater interest than the gradient of the conjugate $D_{\phi}^*(f_* || \nu)$. The latter is a reweighting of the former, since the gradient of the conjugate is a probability distribution, such as the softmax for the Kullback-Leibler divergence.

Optimal critic has bounded Lipschitz constant. Notice that while the variational formula of f-divergences contains a supremum, the formula of their Moreau-Yosida approximations contains a maximum. This means that in the former case, even if the divergence is finite, the supremum might not be achieved by a Lipschitz function. The variational representation (8) only implies that a sequence of Lipschitz

functions converges to a function achieving the supremum, but the limit is not necessarily Lipschitz continuous, in fact it might not even be continuous. On the other hand, for the Moreau-Yosida approximation, the maximum in (17) or (18) is always achieved by a Lipschitz function. Since any neural network is Lipschitz continuous, we argue that a trained critic can provide a better estimate of the Moreau-Yosida approximation, since its target f_* is not only a Csiszár potential of ξ_*, ν but a scaled Kantorovich potential of μ, ξ_* as well, implying that it has a bounded Lipschitz constant.



Figure 4. $||f_{\theta_f}||_{L,\mu_n,\nu_n}$ during training

The $\alpha = 2$ and $\alpha = \infty$ cases. Since f_* is a Kantorovich potential scaled by the coefficient $\beta^{-\alpha}W_1(\mu,\xi_*)^{\alpha-1}$ and the Lipschitz norm of a Kantorovich potential is 1, the case $\alpha > 1$ can be seen as adaptive Lipschitz regularization, with $||f_*||_L$ decaying during training as μ and ν drift closer and $W_1(\mu,\xi_*)$ becomes smaller. We visualized $\|f_{\theta_f}\|_{L,\mu_n,\nu_n}$ in Figure 4 during training with $\alpha = 1.05, 2, \infty$ and $\beta = 1$. Ideally, the Lipschitz norm of the critic would vanish. This can be observed in the $\alpha = \infty$ case, which leads to finding a generated distribution with Wasserstein-1 distance of $\beta = 1$ from the data distribution, accordingly to (21). The best FID in Table 1 indicates that it can be beneficial to choose $\alpha = 2$ even though the Lipschitz norm does not vanish. While the case $\alpha = \infty$ (where we only consider the case \leftarrow) leads to low performance with high values of β and unstable training with low values of β , we found that decaying β e.g. from 0.5 to 0.2 led to the best IS as can be seen in Table 1⁵. The TRIVIAL case does not perform well in this setting, which is not surprising since the exact value of $D_{\phi,\frac{1}{\alpha}\beta^{-\alpha},\alpha}(\mu\|\nu)$ is ∞ if ν is not contained in the W_1 ball of radius β centered at μ , and 0 otherwise.

⁵Numerical instabilities prevented us from evaluating the Jeffreys divergence in this setting.

Preliminary experiments showed that other values of α behave similarly to the ones we considered, which is why we restricted our attention to the representative values 1.05, 2 and ∞ . The implementation was done in TensorFlow, using the residual critic and generator architectures from Gulrajani et al. (2017). Training was done for 100000 iterations, with 5 gradient descent step per iteration for the critic, and 1 for the generator. Additional results, details of the experimental setup and generated images can be found in Appendix 8.4, along with toy examples validating our approach for approximating f-divergences through the tight variational representations on categorical and Gaussian distributions. The original f-GAN losses (Nowozin et al., 2016) were particularly unstable in our implementation. Training the critic for 1 instead of 5 steps per iteration led to more stability, but even in this case only the χ^2 divergence made it to 100000 iterations without numerical errors, leading to an IS of 6.49 and an FID of 40.64. Source code to reproduce the experiments is available at https://github.com/ renyi-ai/moreau-yosida-f-divergences.

6. Related work

In Farnia & Tse (2018), $D_{\phi,1,1}$ is defined, and a non-tight variational representation is given for symmetric choices of D_{ϕ} . They also prove that $D_{\phi,1,1}$ between the data and generated distributions is a continuous function of the generator parameters, and provide a dual formula for the case $\alpha = 2$ using W_2 instead of W_1 . A future direction is to prove analogous results for general α, λ and W_p . In Birrell et al. (2020), a generalization of $D_{\phi,1,1}$ is defined with arbitrary IPMs instead of W_1 , but their assumptions on ϕ are more restrictive, and they explicitly define $D_{\phi}(\mu \| \nu)$ to be ∞ if $\mu \ll \nu$ does not hold. In Husain et al. (2019), the Lipschitz constrained version of the non-tight variational representation of D_{ϕ} is shown to be a lower bound to the Wasserstein autoencoder objective. In Laschos et al. (2019), it is proved that the supremum in the Donsker-Varadhan formula can equivalently be taken over Lipschitz continuous functions. In Song & Ermon (2020), based on the non-tight representation, another generalization of f-GANs and WGAN is proposed, with the importance weights ranalogous to the gradient of $D_{\phi}^{*}(f \| \nu)$ in our case. Connections to density ratio estimation and sample reweighting are discussed, which apply to our case as well. In Arbel et al. (2021), the Lipschitz constrained version of the Donsker-Varadhan formula is proposed as an objective function for energy-based models. For representation learning by mutual information maximization, Ozair et al. (2019) proposes the Lipschitz constrained version of the Donsker-Varadhan formula as a proxy for mutual information, which is shown to be empirically superior to the unconstrained formulation. In Zhou et al. (2019), it is shown that Lipschitz regularization improves the performance of GANs in general other than

the Wasserstein GAN. The uniqueness of the optimal critic is investigated, and formulas are proposed for which uniqueness holds. We solve the uniqueness problem in another way, by implementing the quotient critic.

To summarize, the recognition of the primal formula being the Moreau-Yosida regularization of D_{ϕ} with respect to W_1 and the case $\alpha \neq 1$ are novel to our paper. This includes the unconstrained variational formula for W_1 . Regarding f-divergences, the tight variational representation over the quotient space $Lip(X, x_0)$ and the characterization of Csiszár potentials are new as well. Additionally, we allow the same generality in terms of the choice of ϕ as Agrawal & Horel (2020). On the practical side, we proposed an algorithm to calculate the tight conjugate $D_{\phi}^{*}(f \| \nu)$ and its gradient. Experimentally, implementations are provided for GANs based on the tight variational representation not only of the Kullback-Leibler divergence, but the reverse Kullback-Leibler, χ^2 , reverse χ^2 , squared Hellinger, Jensen-Shannon, Jeffreys, triangular discrimination and total variation divergences as well.

7. Conclusions

In this paper, we studied the Moreau-Yosida regularization of f-divergences with respect to the Wasserstein-1 metric in a convex duality framework. We presented variational formulas and characterizations of optimal variables, generalizing a number of existing results and leading to novel special cases of interest, and proposed the MY f-GAN as an implementation of the formulas. Future directions include finding the variational formulas for Moreau-Yosida approximation with respect to all Wasserstein-p metrics including the case $0 < \alpha < 1$, improving the estimation of the Lipschitz norm of the critic, making use of the fact that Csiszár-Kantorovich potentials can be seen as Lipschitzregularized statistical tests, e.g. for sample reweighting, and scaling up to higher-dimensional datasets. Additionally, the results can potentially be applied to learning algorithms other than GANs, such as representation learning by mutual information maximization, energy-based models, generalized prediction functions and density ratio estimation.

Acknowledgements

The author was supported by the Hungarian National Excellence Grant 2018-1.2.1-NKP-00008 and by the Hungarian Ministry of Innovation and Technology NRDI Office within the framework of the Artificial Intelligence National Laboratory Program, and would like to thank the Artificial Intelligence Research Group at the Alfréd Rényi Institute of Mathematics, especially Dániel Varga and Diego González-Sánchez, for their generous help.

References

- Adler, J. and Lunz, S. Banach wasserstein GAN. In Bengio, S., Wallach, H. M., Larochelle, H., Grauman, K., Cesa-Bianchi, N., and Garnett, R. (eds.), Advances in Neural Information Processing Systems 31: Annual Conference on Neural Information Processing Systems 2018, NeurIPS 2018, December 3-8, 2018, Montréal, Canada, pp. 6755–6764, 2018.
- Agrawal, R. and Horel, T. Optimal bounds between \$f\$divergences and integral probability metrics. *CoRR*, abs/2006.05973, 2020.
- Ali, S. M. and Silvey, S. D. A general class of coefficients of divergence of one distribution from another. *Journal of the Royal Statistical Society. Series B*, 28(1):131–142, 1966.
- Arbel, M., Zhou, L., and Gretton, A. Generalized energy based models. In *International Conference on Learning Representations*, 2021.
- Arjovsky, M., Chintala, S., and Bottou, L. Wasserstein generative adversarial networks. In Precup, D. and Teh, Y. W. (eds.), Proceedings of the 34th International Conference on Machine Learning, ICML 2017, Sydney, NSW, Australia, 6-11 August 2017, volume 70 of Proceedings of Machine Learning Research, pp. 214–223. PMLR, 2017.
- Belghazi, M. I., Baratin, A., Rajeswar, S., Ozair, S., Bengio, Y., Hjelm, R. D., and Courville, A. C. Mutual information neural estimation. In Dy, J. G. and Krause, A. (eds.), Proceedings of the 35th International Conference on Machine Learning, ICML 2018, Stockholmsmässan, Stockholm, Sweden, July 10-15, 2018, volume 80 of Proceedings of Machine Learning Research, pp. 530–539. PMLR, 2018.
- Birrell, J., Dupuis, P., Katsoulakis, M. A., Pantazis, Y., and Rey-Bellet, L. (f, Γ)-divergences: Interpolating between f-divergences and integral probability metrics. *CoRR*, abs/2011.05953, 2020. URL https://arxiv.org/ abs/2011.05953.
- Blondel, M., Martins, A. F. T., and Niculae, V. Learning with fenchel-young losses. *J. Mach. Learn. Res.*, 21: 35:1–35:69, 2020.
- Borwein, J. M. and Lewis, A. S. Partially-finite programming in l₁ and the existence of maximum entropy estimates. *SIAM J. Optim.*, 3(2):248–267, 1993. doi: 10.1137/0803012.
- Boyd, S. P. and Vandenberghe, L. *Convex Optimization*. Cambridge University Press, 2014. ISBN 978-0-521-83378-3. doi: 10.1017/CBO9780511804441.

- Cobzaş, Ş., Miculescu, R., and Nicolae, A. Lipschitz Functions. Lecture Notes in Mathematics. Springer International Publishing, 2019. ISBN 9783030164881.
- Csiszár, I. Eine informationstheoretische ungleichung und ihre anwendung auf den beweis der ergodizität von markoffschen ketten. A Magyar Tudományos Akadémia Matematikai Kutató Intézetének Közleményei, 8(1–2):85– 108, 1963.
- Csiszár, I. Information-type measures of difference of probability distributions and indirect observations. *Studia Scientiarum Mathematicarum Hungarica*, 2:299–318, 1967.
- Csiszár, I., Gamboa, F., and Gassiat, E. MEM pixel correlated solutions for generalized moment and interpolation problems. *IEEE Trans. Inf. Theory*, 45(7):2253–2270, 1999. doi: 10.1109/18.796367.
- Dal Maso, G. *An Introduction to* Γ*-Convergence*. Birkhäuser Boston, Boston, MA, 1993. ISBN 978-1-4612-0327-8. doi: 10.1007/978-1-4612-0327-8.
- Donsker, M. D. and Varadhan, S. R. S. Asymptotic evaluation of certain markov process expectations for large time—iii. *Communications on Pure and Applied Mathematics*, 29(4):389–461, 1976. doi: 10.1002/cpa. 3160290405.
- Farnia, F. and Tse, D. A convex duality framework for gans. In Bengio, S., Wallach, H. M., Larochelle, H., Grauman, K., Cesa-Bianchi, N., and Garnett, R. (eds.), Advances in Neural Information Processing Systems 31: Annual Conference on Neural Information Processing Systems 2018, NeurIPS 2018, December 3-8, 2018, Montréal, Canada, pp. 5254–5263, 2018.
- Gulrajani, I., Ahmed, F., Arjovsky, M., Dumoulin, V., and Courville, A. C. Improved training of wasserstein gans. In Guyon, I., von Luxburg, U., Bengio, S., Wallach, H. M., Fergus, R., Vishwanathan, S. V. N., and Garnett, R. (eds.), Advances in Neural Information Processing Systems 30: Annual Conference on Neural Information Processing Systems 2017, 4-9 December 2017, Long Beach, CA, USA, pp. 5767–5777, 2017.
- Hanin, L. G. An extension of the kantorovich norm. Contemporary Mathematics, 226:113–130, 1999.
- Hjelm, R. D., Fedorov, A., Lavoie-Marchildon, S., Grewal, K., Bachman, P., Trischler, A., and Bengio, Y. Learning deep representations by mutual information estimation and maximization. In 7th International Conference on Learning Representations, ICLR 2019, New Orleans, LA, USA, May 6-9, 2019. OpenReview.net, 2019.

- Husain, H., Nock, R., and Williamson, R. C. A primal-dual link between gans and autoencoders. In Wallach, H. M., Larochelle, H., Beygelzimer, A., d'Alché-Buc, F., Fox, E. B., and Garnett, R. (eds.), Advances in Neural Information Processing Systems 32: Annual Conference on Neural Information Processing Systems 2019, NeurIPS 2019, December 8-14, 2019, Vancouver, BC, Canada, pp. 413–422, 2019.
- Jost, J. and Li-Jost, X. *Calculus of Variations*. Cambridge Studies in Advanced Mathematics. Cambridge University Press, 2008. ISBN 9780521057127.
- Karras, T., Aila, T., Laine, S., and Lehtinen, J. Progressive growing of gans for improved quality, stability, and variation. In 6th International Conference on Learning Representations, ICLR 2018, Vancouver, BC, Canada, April 30 May 3, 2018, Conference Track Proceedings. OpenReview.net, 2018.
- Laschos, V., Obermayer, K., Shen, Y., and Stannat, W. A fenchel-moreau-rockafellar type theorem on the kantorovich-wasserstein space with applications in partially observable markov decision processes. *Journal of Mathematical Analysis and Applications*, 477(2):1133 – 1156, 2019. ISSN 0022-247X. doi: https://doi.org/10. 1016/j.jmaa.2019.05.004.
- Miyato, T., Kataoka, T., Koyama, M., and Yoshida, Y. Spectral normalization for generative adversarial networks. In 6th International Conference on Learning Representations, ICLR 2018, Vancouver, BC, Canada, April 30 - May 3, 2018, Conference Track Proceedings. OpenReview.net, 2018.
- Mohebi, H. Topical Functions and their Properties in a Class of Ordered Banach Spaces, pp. 343–361. Springer US, Boston, MA, 2005. ISBN 978-0-387-26771-5. doi: 10.1007/0-387-26771-9_12.
- Nowozin, S., Cseke, B., and Tomioka, R. f-gan: Training generative neural samplers using variational divergence minimization. In Lee, D. D., Sugiyama, M., von Luxburg, U., Guyon, I., and Garnett, R. (eds.), Advances in Neural Information Processing Systems 29: Annual Conference on Neural Information Processing Systems 2016, December 5-10, 2016, Barcelona, Spain, pp. 271–279, 2016.
- Ozair, S., Lynch, C., Bengio, Y., van den Oord, A., Levine, S., and Sermanet, P. Wasserstein dependency measure for representation learning. In Wallach, H. M., Larochelle, H., Beygelzimer, A., d'Alché-Buc, F., Fox, E. B., and Garnett, R. (eds.), Advances in Neural Information Processing Systems 32: Annual Conference on Neural Information Processing Systems 2019, NeurIPS 2019, 8-14 December 2019, Vancouver, BC, Canada, pp. 15578–15588, 2019.

- Petzka, H., Fischer, A., and Lukovnikov, D. On the regularization of wasserstein gans. In 6th International Conference on Learning Representations, ICLR 2018, Vancouver, BC, Canada, April 30 - May 3, 2018, Conference Track Proceedings. OpenReview.net, 2018.
- Reid, M. D. and Williamson, R. C. Information, divergence and risk for binary experiments. J. Mach. Learn. Res., 12: 731–817, 2011.
- Ruderman, A., Reid, M. D., García-García, D., and Petterson, J. Tighter variational representations of f-divergences via restriction to probability measures. In *Proceedings of* the 29th International Conference on Machine Learning, ICML 2012, Edinburgh, Scotland, UK, June 26 - July 1, 2012. icml.cc / Omnipress, 2012.
- Song, J. and Ermon, S. Bridging the gap between f-gans and wasserstein gans. In Proceedings of the 37th International Conference on Machine Learning, ICML 2020, 13-18 July 2020, Virtual Event, volume 119 of Proceedings of Machine Learning Research, pp. 9078–9087. PMLR, 2020.
- Tao, T. Several Variable Differential Calculus, pp. 127–161. Springer Singapore, Singapore, 2016. ISBN 978-981-10-1804-6. doi: 10.1007/978-981-10-1804-6_6.
- Terjék, D. Adversarial lipschitz regularization. In 8th International Conference on Learning Representations, ICLR 2020, Addis Ababa, Ethiopia, April 26-30, 2020. Open-Review.net, 2020.
- Villani, C. Optimal transport Old and new, volume 338, pp. xxii+973. 01 2008. doi: 10.1007/978-3-540-71050-9.
- Weaver, N. *Lipschitz Algebras*. WORLD SCIENTIFIC, 2nd edition, 2018. doi: 10.1142/9911.
- Zalinescu, C. *Convex Analysis in General Vector Spaces*. World Scientific, 2002. ISBN 9789812380678.
- Zhou, Z., Liang, J., Song, Y., Yu, L., Wang, H., Zhang, W., Yu, Y., and Zhang, Z. Lipschitz generative adversarial nets. In Chaudhuri, K. and Salakhutdinov, R. (eds.), *Proceedings of the 36th International Conference on Machine Learning, ICML 2019, 9-15 June 2019, Long Beach, California, USA*, volume 97 of *Proceedings of Machine Learning Research*, pp. 7584–7593. PMLR, 2019.

8. Appendix

8.1. Background

In order to establish the dual formulation of the Moreau-Yosida approximation of f-divergences, we will apply techniques from convex analysis, for which we need an appropriate pair of vector spaces that are in duality.

8.1.1. FUNCTIONAL ANALYSIS

We recite a number of definitions and results (without proofs) from functional analysis concerning vector spaces of Lipschitz functions, taken from Cobzaş et al. (2019).

Let (X, d) be a compact metric space, and denote the σ -algebra of its Borel subsets by $\mathcal{B}(X)$. A function $\mu : \mathcal{B}(X) \to \mathbb{R}$ is called a σ -additive measure if $\mu(\bigcup_{i=0}^{\infty} A_i) = \sum_{i=0}^{\infty} \mu(A_i)$ holds for every family $\{A_i : i \in \mathbb{N}\} \subset \mathcal{B}(X)$ of pairwise disjoint elements of $\mathcal{B}(X)$. Any such measure is of bounded variation, i.e. $|\mu|(X) < \infty$ where

$$|\mu|(X) = \sup_{(A_i)_{i \in \{1,\dots,m\}} \text{ is a partition of } X, m \in \mathbb{N}} \left\{ \sum_{i=1}^{m} |\mu(A_i)| \right\}$$
(25)

is the total variation of μ . Denote by $\mathcal{M}(X)$ the set of σ -additive measures on $\mathcal{B}(X)$.

A function $f: X \to \mathbb{R}$ is Lipschitz continuous if there exists a number $M \in [0, \infty)$ such that $|f(x) - f(y)| \le Md(x, y)$ for all $x, y \in X$. The Lipschitz norm of such an f is defined as

$$||f||_{L} = \sup_{x,y \in X, x \neq y} \left\{ \frac{|f(x) - f(y)|}{d(x,y)} \right\}.$$
(26)

Denote by Lip(X) the set of Lipschitz continuous functions $f : X \to \mathbb{R}$. Fixing an arbitrary element $x_0 \in X$, the set $Lip(X, x_0) = \{f \in Lip(X) : f(x_0) = 0\}$ is a Banach space with the norm $\|.\|_L$. For any $\xi \in \mathbb{R}$, $\mathcal{M}(X, \xi) = \{\mu \in \mathcal{M}(X) : \mu(X) = \xi\}$ is a vector subspace of $\mathcal{M}(X)$. With the Kantorovich-Rubinstein norm

$$\|\mu\|_{KR} = \sup_{f \in Lip(X, x_0), \|f\|_L \le 1} \left\{ \int f d\mu \right\},$$
(27)

the pair $(\mathcal{M}(X,0), \|.\|_{KR})$ is a normed vector space.

Theorem. For any $f \in Lip(X, x_0)$ the functional $u_f : \mathcal{M}(X, 0) \to \mathbb{R}$ defined by $u_f(\mu) = \int f d\mu$ is linear and continuous with $||u_f|| = ||f||_L$. Moreover, every continuous linear functional v on $(\mathcal{M}(X, 0), ||.||_{KR})$ is of the form $v(\mu) = u_f(\mu)$ for a uniquely determined function $f \in Lip(X, x_0)$ with $||v|| = ||f||_L$. Consequently, the mapping $f \to u_f$ is an isometric isomorphism of $(Lip(X, x_0), ||.||_L)$ onto the topological dual $(\mathcal{M}(X, 0), ||.||_{KR})^*$, i.e.

$$(Lip(X, x_0), \|.\|_L) \cong (\mathcal{M}(X, 0), \|.\|_{KR})^*.$$
(28)

With the norm

$$||f||_{\max} = \max\{||f||_L, ||f||_\infty\},\tag{29}$$

the pair $(Lip(X), \|.\|_{\max})$ is a Banach space. With the Hanin norm

$$\|\mu\|_{H} = \inf_{\nu \in \mathcal{M}(X,0)} \{ \|\nu\|_{KR} + \|\mu - \nu\|_{TV} \},$$
(30)

the pair $(\mathcal{M}(X), \|.\|_H)$ is a normed vector space. The subspace $\mathcal{M}(X, 0)$ is closed with respect to the topology generated by $\|.\|_H$, and the corresponding subspace topology is equivalent to the topology generated by $\|.\|_{KR}$.

Theorem. For any $f \in Lip(X)$ the functional $u_f : \mathcal{M}(X) \to \mathbb{R}$ defined by $u_f(\mu) = \int f d\mu$ is linear and continuous with $||u_f|| = ||f||_{\max}$. Moreover, every continuous linear functional v on $(\mathcal{M}(X), ||.||_H)$ is of the form $v(\mu) = u_f(\mu)$ for a uniquely determined function $f \in Lip(X)$ with $||v|| = ||f||_{\max}$. Consequently, the mapping $f \to u_f$ is an isometric isomorphism of $(Lip(X), ||.||_{\max})$ onto the topological dual $(\mathcal{M}(X), ||.||_H)^*$, i.e.

$$(Lip(X), \|.\|_{\max}) \cong (\mathcal{M}(X), \|.\|_{H})^{*}.$$
 (31)

Integration is bilinear, i.e. $\int (\alpha f + \beta g) d\mu = \alpha \int f d\mu + \beta \int g d\mu$ and $\int f d(\alpha \mu + \beta \nu) = \alpha \int f d\mu + \beta \int f d\nu$ for any $\alpha, \beta \in \mathbb{R}, f, g \in Lip(X) \text{ and } \mu, \nu \in \mathcal{M}(X).$

The set of nonnegative measures is $\mathcal{M}^+(X) = \{ \mu \in \mathcal{M}(X) \mid \forall A \in \mathcal{B}(X) : \mu(A) \ge 0 \}$. The convex set $P(X) = \{ \mu \in \mathcal{M}(X) \mid \forall A \in \mathcal{B}(X) : \mu(A) \ge 0 \}$. $\mathcal{M}(X,1) \cap \mathcal{M}^+(X)$ is exactly the set of all Borel probability measures on X. It is a compact and complete metric space with respect to the Wasserstein-1 metric

$$W_1(\mu,\nu) = \inf_{\pi \in \Pi(\mu,\nu)} \int d(x_1, x_2) d\pi(x_1, x_2),$$
(32)

where $\Pi(\mu,\nu)$ is the set of probability measures on $X \times X$ with marginals μ and ν , i.e. given any $A \in \mathcal{B}(X)$, the relations $\pi(A \times X) = \mu(A)$ and $\pi(X \times A) = \nu(A)$ hold.

Theorem (Kantorovich-Rubinstein duality). The metric induced by the norm $\|.\|_{KB}$ on P(X) is equivalent to the Wasserstein-1 metric as

$$W_1(\mu,\nu) = \sup_{f \in Lip(X,x_0), \|f\|_L \le 1} \left\{ \int f d\mu - \int f d\nu \right\} = \|\mu - \nu\|_{KR}$$
(33)

for any $\mu, \nu \in P(X)$.

8.1.2. CONVEX ANALYSIS

We recite a number of definitions and results (without proofs) from convex analysis on general vector spaces, taken from Zalinescu (2002).

Let X, Y be separated locally convex topological vector spaces with topological duals X^*, Y^* . For an extended real-valued function $f: X \to \overline{\mathbb{R}}$, the function $f^*: X^* \to \overline{\mathbb{R}}$ defined by

$$f^{*}(x^{*}) = \sup_{x \in X} \{ \langle x, x^{*} \rangle - f(x) \}$$
(34)

is the (convex) conjugate of f, where $\langle \cdot, \cdot \rangle$ is the dual pairing. The conjugate $g^* : X \to \overline{\mathbb{R}}$ of a function $g : X^* \to \overline{\mathbb{R}}$ is defined analogously as

$$g^{*}(x) = \sup_{x^{*} \in X^{*}} \{ \langle x, x^{*} \rangle - g(x^{*}) \},$$
(35)

leading to the notion of the biconjugate $(f^*)^* = f^{**}$ of f, which is the greatest lower semicontinuous convex function with $f^{**} \le f.$

If
$$0 < \alpha \in \mathbb{R}$$
, then
 $(\alpha f(\cdot))^*(x^*) = \alpha f^*(\alpha^{-1}x^*).$

$$(\alpha f(\cdot))^*(x^*) = \alpha f^*(\alpha^{-1}x^*), \tag{36}$$

if $0 \neq \beta \in \mathbb{R}$, then

$$(f(\beta \cdot))^*(x^*) = f^*(\beta^{-1}x^*), \tag{37}$$

if $x_0 \in X$, then

$$(f(x_0 + \cdot))^*(x^*) = f^*(x^*) - \langle x_0, x^* \rangle,$$
(38)

and the Young-Fenchel inequality states that

$$\forall x \in X, \forall x^* \in X^* : f(x) + f^*(x^*) \ge \langle x, x^* \rangle.$$
(39)

If $(X, \|\cdot\|)$ and $(X^*, \|\cdot\|_*)$ are normed spaces and $f: X \to \overline{\mathbb{R}}$ is defined by $f(x) = \|x\|$, then

$$f^{*}(x^{*}) = \begin{cases} 0 \text{ if } ||x^{*}||_{*} \leq 1, \\ \infty \text{ otherwise} \end{cases}$$
(40)

is the indicator function of the unit ball of the dual space, and if $\psi : \mathbb{R}_+ \to \overline{\mathbb{R}}_+$ is such that $\psi(0) = 0$ and $f : X \to \overline{\mathbb{R}}$ is defined by $f(x) = \psi(||x||)$, then

$$f^*(x^*) = \psi^{\#}(\|x^*\|_*), \tag{41}$$

where $\psi^{\#}: \mathbb{R}_+ \to \overline{\mathbb{R}}_+$ is defined by

$$\psi^{\#}(s) = \sup_{0 \le t \in \mathbb{R}} \{ st - \psi(t) \}.$$
(42)

Given a function $f: X \to \overline{\mathbb{R}}$, the set

$$\operatorname{dom} f = \{x \in X : f(x) < \infty\}$$
(43)

is the effective domain of f. A function f is proper if dom $f \neq \emptyset$ and $f(x) > -\infty$ for all $x \in X$, otherwise it is improper. A function f is convex if

$$\forall x, y \in X, \forall \lambda \in [0, 1] : f(\lambda x + (1 - \lambda)y) \le \lambda f(x) + (1 - \lambda)f(y)$$
(44)

holds, and strictly convex if (44) holds with \leq replaced by <. An $f : \mathbb{R} \to \overline{\mathbb{R}}$ is strictly convex at $x_0 \in \text{dom } f$ if $\lambda f(x_1) + (1 - \lambda)f(x_2) > f(x_0)$ holds for $\forall \lambda \in (0, 1)$ and $\forall x_1, x_2 \in \text{dom } f$ such that $\lambda x_1 + (1 - \lambda)x_2 = x_0$, unless $x_1 = x_2 = x_0$.

A function f is lower semicontinuous at $x_0 \in X$ if $f(x_0) \leq \liminf_{x \to x_0} f(x)$, and f is lower semicontinuous if it is lower semicontinuous at $\forall x_0 \in X$.

Theorem (Fenchel-Moreau biconjugation). Let X be a separated locally convex topological vector space with topological dual X^* , and $f: X \to \overline{\mathbb{R}}$ a function. Then $f^{**} \leq f$, and the relation

$$f = f^{**}. (45)$$

i.e. f is equivalent to its biconjugate (the conjugate of its conjugate) holds if and only if either f is proper, lower semicontinuous and convex, or f is constant $\pm \infty$.

Given a function $f: X \to \overline{\mathbb{R}}$ and $\hat{x} \in X$, the subdifferential of f at \hat{x} is defined as the set

$$\partial f(\hat{x}) = \{ x^* \in X^* : \forall x \in X : \langle x - \hat{x}, x^* \rangle \le f(x) - f(\hat{x}) \}.$$

$$\tag{46}$$

Any $x \in \partial f(\hat{x})$ is called a subgradient of f at \hat{x} . It is possible that $\partial f(\hat{x}) = \emptyset$, which is always the case if $f(\hat{x}) = \pm \infty$. It holds that if f is proper and $f(x) \in \mathbb{R}$, (39) becomes an equality if and only if $x^* \in \partial f(x)$, or equivalently $x \in \partial f^*(x^*)$. It follows that if f is proper, convex and lower semicontinuous, then

$$f(x) = f^{**}(x) = \sup_{x^* \in X^*} \left\{ \langle x, x^* \rangle - f^*(x^*) \right\} = \langle x, \hat{x}^* \rangle - f^*(\hat{x}^*) \iff \hat{x}^* \in \partial f(x)$$
(47)

and

$$f^*(x^*) = \sup_{x \in X} \left\{ \langle x, x^* \rangle - f(x) \right\} = \langle \hat{x}, x^* \rangle - f(\hat{x}) \iff \hat{x} \in \partial f^*(x^*)$$
(48)

both hold.

The adjoint of a linear map $A : X \to Y$ is the linear map $A^* : Y^* \to X^*$ such that $\langle Ax, y^* \rangle = \langle x, A^*y^* \rangle$ holds for $\forall x \in X, y^* \in Y^*$.

Theorem. ⁶ Let $A : X \to \mathbb{R}$ be a continuous linear map (so that $A \in X^*$) with adjoint $A^* : \mathbb{R} \to X^*$, and $f : X \to \overline{\mathbb{R}}, g : \mathbb{R} \to \overline{\mathbb{R}}$ be proper convex functions, and consider the proper convex function $h : X \to \overline{\mathbb{R}}$ defined as h(x) = f(x) + g(Ax). If dom $f \cap A^{-1}(\operatorname{dom} g) \neq \emptyset$ and $0 \in \operatorname{relint}(A(\operatorname{dom} f) - \operatorname{dom} g)$, then it holds that

$$h^{*}(x^{*}) = \min_{\gamma \in \mathbb{R}} \left\{ f^{*}(x^{*} - A^{*}\gamma) + g^{*}(\gamma) \right\}$$
(49)

and

$$\partial h(x) = \partial f(x) + A^*(\partial g(Ax)).$$
(50)

Theorem (Fenchel-Rockafellar duality). ⁷ Let $f, g : X \to \overline{\mathbb{R}}$ be proper convex functions for which $\exists x_0 \in \text{dom } f \cap \text{dom } g$ such that g is continuous at x_0 . Then it holds that

$$\inf_{x \in X} \left\{ f(x) + g(x) \right\} = \max_{x^* \in X^*} \left\{ -f^*(x^*) - g^*(-x^*) \right\}$$
(51)

and

$$\exists \hat{x} \in X : \inf_{x \in X} \left\{ f(x) + g(x) \right\} = f(\hat{x}) + g(\hat{x}) \iff \exists \hat{x}^* \in X^* : -\hat{x}^* \in \partial f(\hat{x}) \land \hat{x}^* \in \partial g(\hat{x}).$$
(52)

⁶This theorem is Zalinescu (2002, Theorem 2.8.3(viii)) with $Y = \mathbb{R}$.

⁷This theorem is Zalinescu (2002, Corollary 2.8.5) and condition Zalinescu (2002, Theorem 2.8.3(iii)) with $Y = X, A : X \to X$ the identity and replacing the dual variable x^* with $-x^*$.

8.1.3. MOREAU-YOSIDA APPROXIMATION

Let (X, d) be a metric space and $f : X \to \mathbb{R}$ a proper function. Given $0 < \lambda, \alpha \in \mathbb{R}$, the Moreau-Yosida approximation (Dal Maso, 1993; Jost & Li-Jost, 2008) of index λ and order α of f is defined as

$$f_{\lambda,\alpha}(x) = \inf_{y \in X} \{ f(y) + \lambda d(x, y)^{\alpha} \}.$$
(53)

It holds that

$$\overline{f}(x) = \sup_{\lambda > 0} f_{\lambda,\alpha}(x) = \lim_{\lambda \to \infty} f_{\lambda,\alpha}(x).$$
(54)

where \overline{f} is the greatest lower semicontinuous function with $\overline{f} \leq f$.

Theorem. If $0 < \alpha \le 1$, then $(f_{\lambda_1,\alpha})_{\lambda_2,\alpha} = f_{\min(\lambda_1,\lambda_2),\alpha}$, and $f_{\lambda,\alpha}$ is the greatest function among those $g \le f$ for which

$$\forall x, y \in X : |g(x) - g(y)| \le \lambda d(x, y)^{\alpha}$$
(55)

(i.e. g is Hölder continuous with exponent α and Hölder constant λ) holds.

The functions $f_{\lambda,1}$ *satisfy*

$$\forall x, y \in X : |f_{\lambda,1}(x) - f_{\lambda,1}(y)| \le \lambda d(x, y)$$
(56)

(i.e. they are Lipschitz continuous with Lipschitz constant λ).

If $\alpha \geq 1$, f is non-negative and $0 < r \in \mathbb{R}$, $0 \leq M \in \mathbb{R}$ are constants, then there exists a constant $0 < c(\alpha, \lambda, M, r) \in R$ such that for $\forall z \in X$ it holds that if $f_{\lambda,\alpha}(z) \leq M$, then

$$\forall x, y \in X, d(x, z) \le r, d(y, z) \le r : |f_{\lambda, \alpha}(x) - f_{\lambda, \alpha}(y)| \le cd(x, y)$$
(57)

(*i.e.* $f_{\lambda,\alpha}$ is locally Lipschitz continuous).

8.2. Proofs

Proposition 6. Given $\nu \in \mathcal{M}^+(X)$ and a proper, convex and lower semicontinuous function $\phi : \mathbb{R} \to \overline{\mathbb{R}}$ strictly convex at 1 with $\phi(1) = 0$, the function $I_{\phi,\nu} : (\mathcal{M}(X), \|.\|_H) \to \overline{\mathbb{R}}$ defined by

$$I_{\phi,\nu}(\mu) = D_{\phi}(\mu \| \nu) \tag{58}$$

is proper, convex and lower semicontinuous, and its conjugate $I_{\phi,\nu}^* : (Lip(X), \|.\|_{\max}) \to \overline{\mathbb{R}}$ is

$$I_{\phi,\nu}^*(f) = \begin{cases} \int \phi^* \circ f d\nu \ if \ f(X) \subseteq [\phi'(-\infty), \phi'(\infty)],\\ \infty \ otherwise. \end{cases}$$
(59)

Proof. By Agrawal & Horel (2020, Proposition 4.2.6), one has

$$\sup_{\mu \in \mathcal{M}(X)} \left\{ \int f d\mu - I_{\phi,\nu}(\mu) \right\} = \begin{cases} \int \phi^* \circ f d\nu \text{ if } f(X) \subseteq [\phi'(-\infty), \phi'(\infty)], \\ \infty \text{ otherwise.} \end{cases}$$
(60)

for any bounded and measurable $f : X \to \mathbb{R}$. Any $f \in Lip(X)$ is bounded and measurable, hence the claimed conjugate relation. Clearly $I_{\phi,\nu}$ is convex and proper. For it to be lower semicontinuous, by (45) we only need to show that $I_{\phi,\nu}^{**} \ge I_{\phi,\nu}$, i.e. that there exists a sequence (f_n) in Lip(X) such that $\lim_{n\to\infty} \int f_n d\mu - I_{\phi,\nu}^*(f_n) \ge I_{\phi,\nu}(\mu)$.

By Borwein & Lewis (1993, Theorem 2.7) and (45), if $\operatorname{supp}(\nu) = X$, then

$$I_{\phi,\nu}(\mu) = \sup_{f \in C(X)} \left\{ \int f d\mu - \int \phi^* \circ f d\nu \right\}$$
(61)

holds with C(X) being the space of continuous functions on X. By Agrawal & Horel (2020, Lemma 3.2.11), the closure of dom ϕ^* is $[\phi'(-\infty), \phi'(\infty)]$, so that for $\int \phi^* \circ f d\nu < \infty$ to hold, $f(X) \subseteq [\phi'(-\infty), \phi'(\infty)]$ is necessary, meaning that the above supremum is equivalent to

$$\sup_{f \in C(X), f(X) \subseteq [\phi'(-\infty), \phi'(\infty)]} \left\{ \int f d\mu - \int \phi^* \circ f d\nu \right\} = \sup_{f \in C(X)} \left\{ \int f d\mu - I^*_{\phi, \nu}(f) \right\},\tag{62}$$

i.e. there exists a sequence (f_n) in C(X) such that $\lim_{n\to\infty} \int f_n d\mu - I^*_{\phi,\nu}(f_n) = I_{\phi,\nu}(\mu)$.

In the general case with $\operatorname{supp}(\nu) \subset X$, since the support of ν is closed by definition, being a closed subset of a compact metric space, it is a compact metric space itself with the restriction of the metric d. Consider the decomposition $\mu = \mu_1 + \mu_2$ defined as $\mu_1(A) = \mu(A \cap \operatorname{supp}(\nu))$ and $\mu_2(A) = \mu(A \setminus \operatorname{supp}(\nu))$. Then one has

$$D_{\phi}(\mu \| \nu) = D_{\phi}(\mu_1 \| \nu) + \phi'(\infty) \mu_2^+(X) - \phi'(-\infty) \mu_2^-(X),$$
(63)

and by the above considerations there exists a sequence (f_n) in $C(\operatorname{supp}(\nu))$ such that $\lim_{n\to\infty} \int f_n d\mu_1 - I^*_{\phi,\nu}(f_n) = I_{\phi,\nu}(\mu_1) = D_{\phi}(\mu_1 \| \nu)$ and $f_n(X) \subseteq [\phi'(-\infty), \phi'(\infty)]$ for $\forall n$. By the regularity of the measures μ_2^+, μ_2^- , one has

$$\mu_2^{\pm}(A) = \sup_{B \in \mathcal{B}(X), B \subseteq A, B \text{ compact}} \left\{ \mu_2^{\pm}(B) \right\}$$
(64)

for $\forall A \in \mathcal{B}(X)$, i.e. there exist sequences (B_n^{\pm}) in $\mathcal{B}(X)$ such that $\lim_{n \to \infty} \mu_2^{\pm}(B_n) = \mu_2^{\pm}(X)$ with (B_n^{\pm}) compact, and therefore closed. By the definition of μ_2 , we can assume without loss of generality that $B_n^{\pm} \cap \operatorname{supp}(\nu) = \emptyset$ for $\forall n$, and by the definition of the Jordan decomposition, $B_n^+ \cap B_n^- = \emptyset$ can be assumed as well. Define $\tilde{f}_n : \operatorname{supp}(\nu) \cup B_n^+ \cup B_n^- \to \mathbb{R}$ as

$$\tilde{f}_n = \begin{cases} f_n(x) \text{ if } x \in \operatorname{supp}(\nu), \\ \beta_n^+ \text{ if } x \in B_n^+, \\ \beta_n^- \text{ if } x \in B_n^- \end{cases}$$
(65)

with sequences $(\beta_n^{\pm}) \subset \mathbb{R} \cap [\phi'(-\infty), \phi'(\infty)]$ such that $\lim_{n \to \infty} \beta_n^{\pm} = \phi'(\pm \infty)$ and $\forall n : \phi'(-\infty) < \beta_n^{\pm} < \phi'(\infty)$. Since \tilde{f}_n is clearly continuous, by the Tietze extension theorem, there exists a continuous extension $\hat{f}_n \in C(X)$ agreeing with \tilde{f}_n on $\operatorname{supp}(\nu) \cup B_n^+ \cup B_n^-$ with $\hat{f}_n(X) \subseteq [\phi'(-\infty), \phi'(\infty)]$, for which one has

$$\lim_{n \to \infty} \int \hat{f}_n d\mu - I^*_{\phi,\nu}(\hat{f}_n) = \lim_{n \to \infty} \int \hat{f}_n d\mu - \int \phi^* \circ \hat{f}_n d\nu$$

$$= \lim_{n \to \infty} \int f_n d\mu_1 + \int \hat{f}_n d\mu_2 - \int \phi^* \circ f_n d\nu = D_{\phi}(\mu_1 \| \nu) + \lim_{n \to \infty} \int \hat{f}_n d\mu_2$$

$$= D_{\phi}(\mu_1 \| \nu) + \lim_{n \to \infty} \beta_n^+ \mu_2^+(B_n^+) - \beta_n^- \mu_2^-(B_n^-) + \int_{X \setminus B_n^+} \hat{f}_n d\mu_2^+ - \int_{X \setminus B_n^-} \hat{f}_n d\mu_2^-. \quad (66)$$

If $\phi'(\pm \infty)$ are finite, then \hat{f}_n is bounded uniformly independent of n, implying that

$$\lim_{n \to \infty} \int_{X \setminus B_n^+} \hat{f}_n d\mu_2^+ - \int_{X \setminus B_n^-} \hat{f}_n d\mu_2^- = 0.$$
(67)

If one of $\phi'(\pm \infty)$ is infinite, say $\phi'(\infty)$, then

$$\lim_{n \to \infty} \beta_n^+ \mu_2^+(B_n^+) - \beta_n^- \mu_2^-(B_n^-) + \int_{X \setminus B_n^+} \hat{f}_n d\mu_2^+ - \int_{X \setminus B_n^-} \hat{f}_n d\mu_2^- = \begin{cases} \infty \text{ if } \mu_2^+ \neq 0, \\ 0 \text{ otherwise}, \end{cases}$$
(68)

with the case $\phi'(-\infty) = -\infty$ following similarly. If $\phi'(\pm\infty)$ are both infinite, then the above limit is ∞ if $\mu_2 \neq 0$ and 0 otherwise, meaning that in any case

$$\lim_{n \to \infty} \beta_n^+ \mu_2^+(B_n^+) - \beta_n^- \mu_2^-(B_n^-) + \int_{X \setminus B_n^+} \hat{f}_n d\mu_2^+ - \int_{X \setminus B_n^-} \hat{f}_n d\mu_2^- = \phi'(\infty) \mu_2^+(X) - \phi'(-\infty) \mu_2^-(X), \tag{69}$$

so that one has

$$\lim_{n \to \infty} \int \hat{f}_n d\mu - I^*_{\phi,\nu}(\hat{f}_n) = D_\phi(\mu_1 \| \nu) + \phi'(\infty)\mu_2^+(X) - \phi'(-\infty)\mu_2^-(X) = D_\phi(\mu \| \nu) = I_{\phi,\nu}(\mu).$$
(70)

This proves that

$$I_{\phi,\nu}(\mu) \le \sup_{f \in C(X)} \left\{ \int f d\mu - I_{\phi,\nu}^*(f) \right\}$$
(71)

holds for ν with supp $(\nu) \subset X$. Since X is compact, by the Stone-Weierstrass theorem Lip(X) is dense in C(X), hence

$$\sup_{f \in C(X)} \left\{ \int f d\mu - I_{\phi,\nu}^*(f) \right\} = \sup_{f \in Lip(X)} \left\{ \int f d\mu - I_{\phi,\nu}^*(f) \right\} = I_{\phi,\nu}^{**}(\mu), \tag{72}$$

giving the claim $I_{\phi,\nu}(\mu) \leq I_{\phi,\nu}^{**}(\mu)$.

We get the non-tight variational representation over Lip(X) as a corollary by (45).

Corollary 1. Given $\nu \in \mathcal{M}^+(X)$, $\mu \in \mathcal{M}(X)$ and a proper, convex and lower semicontinuous function $\phi : \mathbb{R} \to \overline{\mathbb{R}}$ strictly convex at 1 with $\phi(1) = 0$, one has

$$D_{\phi}(\mu \| \nu) = \sup_{f \in Lip(X), f(X) \subseteq [\phi'(-\infty), \phi'(\infty)]} \left\{ \int f d\mu - \int \phi^* \circ f d\nu \right\}.$$
(73)

Proposition 7. Given $\nu \in P(X)$ and a proper, convex and lower semicontinuous function $\phi : \mathbb{R} \to \overline{\mathbb{R}}$ strictly convex at 1 with $\phi(1) = 0$ and $1 \in \text{relint dom } \phi$, the function $D_{\phi,\nu} : (\mathcal{M}(X), \|.\|_H) \to \overline{\mathbb{R}}$ defined by

$$D_{\phi,\nu}(\mu) = D_{\phi}(\mu \| \nu) + i_{P(X)}(\mu)$$
(74)

is proper, convex and lower semicontinuous, and its conjugate $D^*_{\phi,\nu}$: $(Lip(X), \|.\|_{\max}) \to \overline{\mathbb{R}}$ is

$$D^*_{\phi,\nu}(f) = \min_{\gamma \in \mathbb{R}, \gamma \ge \sup f(X) - \phi'(\infty)} \left\{ \int \phi^*_+ \circ (f - \gamma) d\nu + \gamma \right\},\tag{75}$$

for which $D^*_{\phi,\nu}(f+C) = D^*_{\phi,\nu}(f) + C$ and $D^*_{\phi,\nu}(f_1) \leq D^*_{\phi,\nu}(f_2)$ holds for $\forall C \in \mathbb{R}$ and $f_1 \leq f_2$, meaning that $D^*_{\phi,\nu}$ is a topical function.

Proof. By Agrawal & Horel (2020, Lemma 4.3.1), one has $D_{\phi,\nu}(\mu) = I_{\phi+\nu}(\mu) + i_{\{1\}}(\mu(X))$ with $\phi_+ = \phi + i_{\mathbb{R}_+}$. For the constant function $1 \in Lip(X)$, the map $(\mu \to \mu(X) = \langle \mu, 1 \rangle = \int_X 1d\mu) : \mathcal{M}(X) \to \mathbb{R}$ is linear and continuous, and its adjoint is clearly $(\gamma \to (x \to \gamma)) : \mathbb{R} \to Lip(X)$, mapping constants in \mathbb{R} to the corresponding constant functions in Lip(X). Since $i_{\{1\}} : \mathbb{R} \to \overline{\mathbb{R}}$ is the indicator function of the set $\{1\}$ with its conjugate $i_{\{1\}}^* = (s \to \sup_{t \in \{1\}} \{st\} = s)$ being the identity function, by (49) one has

$$D^{*}_{\phi,\nu}(f) = \min_{\gamma \in \mathbb{R}} \{ I^{*}_{\phi_{+},\nu}(f-\gamma) + \gamma \},$$
(76)

where the minimum can be equivalently taken over those $\gamma \in \mathbb{R}$ for which $I_{\phi_+,\nu}^*(f-\gamma)$ can be finite, i.e. for which $\forall x \in X : \phi'_+(-\infty) \leq f(x) - \gamma \leq \phi'_+(\infty)$ holds for $\forall x \in X$, with the first half being vacuous since $\phi'_+(-\infty) = -\infty$, leading to the claimed conjugate relation. Since $D_{\phi,\nu}$ is the sum of two lower semicontinuous functions, it is itself lower semicontinuous. It is clearly proper and convex as well.

To see that the conditions of (49) hold, notice that $\nu \in \text{dom } I_{\phi_+,\nu}$ always holds, while $\text{dom } i_{\{1\}} = \{1\}$ so that $P(X) \subset (\mu \to \mu(X))^{-1} \text{dom } i_{\{1\}} = \{\mu \in \mathcal{M}(X) : \mu(X) = 1\}$, meaning that $\nu \in \text{dom } I_{\phi_+,\nu} \cap (\mu \to \mu(X))^{-1} \text{dom } i_{\{1\}} \neq \emptyset$. Since $1 \in \text{relint } \text{dom } \phi$ by assumption, either $\text{dom } \phi = \{1\}$, or $\text{dom } \phi$ contains a neighborhood of 1. In the former case, $\phi = i_{\{1\}}$, and one has $\text{dom } I_{\phi_+,\nu} \cap (\mu \to \mu(X))^{-1} \text{dom } i_{\{1\}} = \{\nu\}$, so that $(\mu \to \mu(X)) \text{dom } I_{\phi_+,\nu} - \text{dom } i_{\{1\}} = \{1\} - \{1\} = \{0\}$, for which relint $\{0\} = \{0\}$, and the condition holds.

For other choices of ϕ , there exists $0 < a \in \mathbb{R}$ such that $(1 - a, 1 + a) \subseteq \operatorname{dom} \phi$, and one has for $\forall b \in (1 - a, 1 + a)$ that $I_{\phi_+,\nu}(b\nu) = \int \phi_+(b)d\nu = \phi_+(b) < \infty$, so that $\{b\nu : b \in (1 - a, 1 + a)\} \subseteq \operatorname{dom} I_{\phi_+,\nu}$. This implies that $(1 - a, 1 + a) \subseteq (\mu \to \mu(X)) \operatorname{dom} I_{\phi_+,\nu}$, further implying that $(-a, a) \subseteq (\mu \to \mu(X)) \operatorname{dom} I_{\phi_+,\nu} - \operatorname{dom} \phi_+$, for which $0 \in \operatorname{relint}(\mu \to \mu(X)) \operatorname{dom} I_{\phi_+,\nu} - \operatorname{dom} \phi_+$ clearly holds, proving that the conditions of (49) hold.

The constant additivity property follows from

$$D_{\phi,\nu}^{*}(f+C) = \sup_{\mu \in \mathcal{M}(X)} \int (f+C)d\mu - D_{\phi,\nu}(\mu) = \sup_{\mu \in P(X)} \int (f+C)d\mu - D_{\phi,\nu}(\mu)$$
$$= \sup_{\mu \in P(X)} \int fd\mu + \int Cd\mu - D_{\phi,\nu}(\mu) = \sup_{\mu \in P(X)} \int fd\mu + C - D_{\phi,\nu}(\mu) = D_{\phi,\nu}^{*}(f) + C, \quad (77)$$

and the other from

$$f_1 \le f_2 \implies \int f_1 d\mu - D_{\phi,\nu}(\mu) \le \int f_2 d\mu - D_{\phi,\nu}(\mu) \tag{78}$$

for $\forall \mu \in P(X)$. These properties define topical functions (Mohebi, 2005).

Proposition 8. Given $\nu \in P(X)$, $\omega \in \mathcal{M}(X)$ with $\omega(X) = 1$ and a proper, convex and lower semicontinuous function $\phi: \mathbb{R} \to \overline{\mathbb{R}}$ strictly convex at 1 with $\phi(1) = 0$ and $1 \in \text{relint dom } \phi$, the function $D_{\phi,\nu,\omega}: (\mathcal{M}(X,0), \|.\|_{KR}) \to \overline{\mathbb{R}}$ defined by

$$D_{\phi,\nu,\omega}(\mu) = D_{\phi}(\mu + \omega \| \nu) + i_{P(X)}(\mu + \omega)$$
(79)

is proper, convex and lower semicontinuous, and its conjugate $D^*_{\phi,\nu,\omega}$: $(Lip(X, x_0), \|.\|_L) \to \overline{\mathbb{R}}$ is

$$D^*_{\phi,\nu,\omega}(f) = \min_{\gamma \in \mathbb{R}, \gamma \ge \sup f(X) - \phi'(\infty)} \left\{ \int \phi^*_+ \circ (f - \gamma) d\nu + \gamma \right\} - \int f d\omega.$$
(80)

Proof. By the previous proposition and (38), the conjugate relation

$$(\mu \to D_{\phi,\nu}(\mu+\omega))^* = \left(f \to D^*_{\phi,\nu}(f) - \int f d\omega\right) \tag{81}$$

holds. Notice that for $D_{\phi,\nu}(\mu+\omega)$ to be finite, $\mu(X) = 0$ must hold, since $(\mu+\omega)(X) = 1$ is needed, and $\omega(X) = 1$ by assumption, hence for $f \in Lip(X)$, one has

$$\sup_{\mu \in \mathcal{M}(X)} \left\{ \int f d\mu - D_{\phi,\nu}(\mu + \omega) \right\} = \sup_{\mu \in \mathcal{M}(X,0)} \left\{ \int f d\mu - D_{\phi,\nu}(\mu + \omega) \right\},\tag{82}$$

which is exactly the definition of the value at f of the conjugate of the restriction of $(\mu \to D_{\phi,\nu}(\mu + \omega))$ to $\mathcal{M}(X,0)$. This restriction is clearly proper and convex, and lower semicontinuous as well by being the restriction of a lower semicontinuous function to a closed subspace.

As a corollary, we get the following dual representation of D_{ϕ} on the space of probability measures, which is the tightest to date, in the sense that the supremum is taken over a set of functions that is a proper subset of those of the previous dual representation (Agrawal & Horel, 2020), which included all bounded and measurable functions. Our representation is over the space of Lipschitz functions vanishing at an arbitrary base point.

Corollary 2. Given $\mu, \nu \in P(X)$ and a proper, convex and lower semicontinuous function $\phi : \mathbb{R} \to \mathbb{R}$ strictly convex at 1 with $\phi(1) = 0$ and $1 \in \text{relint dom } \phi$, $D_{\phi}(\mu \| \nu)$ has the equivalent variational representation

$$\sup_{f \in Lip(X,x_0)} \left\{ \int f d\mu - \min_{\gamma \in \mathbb{R}, \gamma \ge \sup f(X) - \phi'(\infty)} \left\{ \int \phi_+^* \circ (f - \gamma) d\nu + \gamma \right\} \right\}.$$
(83)

Proof. Let $\mu \in P(X)$. By (45) and the previous proposition,

$$D_{\phi,\nu,\omega}(\mu-\omega) = D_{\phi,\nu,\omega}^{**}(\mu-\omega)$$

$$= \sup_{f \in Lip(X,x_0)} \left\{ \int f d(\mu-\omega) - \min_{\sup f(X) - \phi'(\infty) \le \gamma} \left\{ \int \phi_+^* \circ (f-\gamma) d\nu + \gamma \right\} + \int f d\omega \right\}$$
(84)
$$\square$$
Ids. giving the claim.

holds, giving the claim.

Proposition 9. Given $\mu, \nu \in P(X)$ and a proper, convex and lower semicontinuous function $\phi : \mathbb{R} \to \overline{\mathbb{R}}$ strictly convex at 1 with $\phi(1) = 0$ and $1 \in \text{relint dom } \phi$, the relation

$$D_{\phi}(\mu \| \nu) = \int f_* d\mu - \min_{\gamma \in \mathbb{R}, \gamma \ge \sup f_*(X) - \phi'(\infty)} \left\{ \int \phi_+^* \circ (f_* - \gamma) d\nu + \gamma \right\}$$
(85)

holds for $f_* \in Lip(X)$ if and only if there exists $C \in \mathbb{R}$ such that the conditions

$$\sup f_*(X) + C \le \phi'(\infty),\tag{86}$$

$$\frac{d\mu_c}{d\nu}(x) \in \partial\phi_+^*(f_*(x) + C) \text{ almost everywhere with respect to }\nu$$
(87)

and

$$\operatorname{supp}(\mu_s) \subset \{x \in X : f_*(x) + C = \phi'(\infty)\}$$
(88)

hold. If ϕ is of Legendre type, the second condition is equivalent to

$$f_*(x) + C = \phi'_+ \left(\frac{d\mu_c}{d\nu}(x)\right) \text{ almost everywhere with respect to } \mu_c.$$
(89)

Proof. By Borwein & Lewis (1993, Theorem 2.10), given $\mu \in \mathcal{M}(X)$ and $\nu \in \mathcal{M}^+(X)$ with $\operatorname{supp}(\nu) = X$, one has

$$D_{\phi}(\mu \| \nu) + \int \phi^* \circ f_* d\nu = \int f_* d\mu$$
(90)

for $f_*: X \to \mathbb{R}$ continuous if and only if

$$f_*(x) \in [\phi'(-\infty), \phi'(\infty)] \text{ for } \forall x \in X,$$
(91)

$$\frac{d\mu_c}{d\nu}(x) \in \partial\phi^*(f_*(x)) \text{ almost everywhere with respect to }\nu,$$
(92)

$$\operatorname{supp} \mu_s^- \subset \{f_*(x) = \phi'(-\infty)\}$$
(93)

and

$$\operatorname{supp} \mu_s^+ \subset \{f_*(x) = \phi'(\infty)\}$$
(94)

hold.

In the general case with $\operatorname{supp}(\nu) \subset X$, since the support of ν is closed by definition, it is a compact metric space itself with the restriction of the metric d. Consider again the decomposition $\mu = \mu_1 + \mu_2$ defined as $\mu_1(A) = \mu(A \cap \operatorname{supp}(\nu))$ and $\mu_2(A) = \mu(A \setminus \operatorname{supp}(\nu))$. Then one has

$$D_{\phi}(\mu \| \nu) = D_{\phi}(\mu_1 \| \nu) + \phi'(\infty)\mu_2^+(X) - \phi'(-\infty)\mu_2^-(X),$$
(95)

while the optimality conditions above imply that

$$D_{\phi}(\mu_1 \| \nu) + \int \phi^* \circ f_* d\nu = \int f_* d\mu_1$$
(96)

for $f_*: X \to \mathbb{R}$ continuous if and only if

$$f_*(x) \in [\phi'(-\infty), \phi'(\infty)] \text{ for } \forall x \in \operatorname{supp}(\nu),$$
(97)

$$\frac{d\mu_{1c}}{d\nu}(x) \in \partial \phi^*(f_*(x)) \text{ almost everywhere with respect to } \nu,$$
(98)

$$\operatorname{supp} \mu_{1s}^{-} \subset \{f_*(x) = \phi'(-\infty)\}$$
(99)

and

$$supp \,\mu_{1s}^+ \subset \{f_*(x) = \phi'(\infty)\}$$
(100)

hold. For

$$D_{\phi}(\mu \| \nu) + \int \phi^* \circ f_* d\nu = \int f_* d\mu \tag{101}$$

to hold, one needs additionally that

$$\phi'(\infty)\mu_2^+(X) - \phi'(-\infty)\mu_2^-(X) = \int f_* d\mu_2, \tag{102}$$

which holds exactly if

$$\operatorname{supp} \mu_{2s}^- \subset \{f_*(x) = \phi'(-\infty)\}$$
 (103)

and

$$\sup \mu_{2s}^{+} \subset \{f_{*}(x) = \phi'(\infty)\}$$
(104)

hold. To summarize, since $\mu_{1c} = \mu_c$ and $\mu_{1s} + \mu_{2s} = \mu_s$, the optimality conditions for $\mu \in \mathcal{M}(X)$ and $\nu \in \mathcal{M}^+(X)$ with ν not necessarily having full support are the same as cited above from (Borwein & Lewis, 1993).

Now consider the tight representation, which follows by taking the conjugate of $(\mu \to D_{\phi}(\mu \| \nu) + i_{P(X)}(\mu) = D_{\phi_{+}}(\mu \| \nu)) + i_{\{1\}}(\int 1d\mu)$ through (49). Substituting into (50), one has

$$\partial \left(D_{\phi_+}(\cdot \|\nu) + i_{\{1\}} \left(\int 1 d \cdot \right) \right) (\mu) = \partial D_{\phi_+}(\cdot \|\nu)(\mu) + (\gamma \to (x \to \gamma)) \left(\partial i_{\{1\}} \left(\int 1 d \mu \right) \right), \tag{105}$$

which gives the claim since $\partial i_{\{1\}}(1) = \mathbb{R}$, $\phi'_+(-\infty) = -\infty$, $\phi'_+(\infty) = \phi'(\infty)$, $\mu_s^- = 0$ and subdifferentials are exactly those f_* for which the supremum is achieved.

If ϕ is of Legendre type (Borwein & Lewis, 1993), then ϕ_+ and ϕ_+^* are both continuously differentiable on their respective domains, while ϕ_+^*' is increasing, and invertible on the subset of dom ϕ_+^* where its value is positive with its inverse given by the strictly increasing ϕ'_+ by Borwein & Lewis (1993, Lemma 2.6). With these, the second condition is equivalent to

$$f_*(x) + C = \phi'_+ \left(\frac{d\mu_c}{d\nu}(x)\right) \ \mu_c$$
-a.e. (106)

Since we only consider the case of compact sample spaces, the infimum defining the Moreau-Yosida approximation turns into a minimum.

Proposition 10. If the metric space (X, d) is compact, the infimum defining the Moreau-Yosida approximation of any f-divergence with respect to the Wasserstein-1 distance is always achieved as

$$\inf_{\xi \in P(X)} \left\{ D_{\phi}(\xi \| \nu) + \lambda W_1(\mu, \xi)^{\alpha} \right\} = \min_{\xi \in P(X)} \left\{ D_{\phi}(\xi \| \nu) + \lambda W_1(\mu, \xi)^{\alpha} \right\}$$
(107)

for any $0 < \lambda, \alpha \in \mathbb{R}$.

Proof. If (X, d) is compact, then $(P(X), W_1)$ is compact as well. Since $(\xi \to D_{\phi}(\xi \| \nu))$ is lower semicontinuous and $(\xi \to W_1(\mu, \xi))$ is continuous, the sum $(\xi \to D_{\phi}(\xi \| \nu) + \lambda W_1(\mu, \xi)^{\alpha})$ is lower semicontinuous. The proposition follows from the fact that a lower semicontinuous function on a compact metric space always has a minimum.

A number of properties follow from the theory of Moreau-Yosida approximation.

Proposition 11. For any $0 < \alpha \in \mathbb{R}$, one has $D_{\phi}(\mu \| \nu) = \sup_{\lambda > 0} D_{\phi,\lambda,\alpha}(\mu \| \nu) = \lim_{\lambda \to \infty} D_{\phi,\lambda,\alpha}(\mu \| \nu)$. Moreover,

- if $0 < \alpha < 1$, then $D_{\phi,\lambda,\alpha}(\cdot \|\nu)$ is Hölder continuous with respect to W_1 with exponent α and Hölder constant λ ,
- if $\alpha = 1$, $D_{\phi,\lambda,\alpha}(\cdot \| \nu)$ is Lipschitz continuous with respect to W_1 with Lipschitz constant λ , and
- if $\alpha > 1$, then $D_{\phi,\lambda,\alpha}(\cdot \|\nu)$ is locally Lipschitz continuous with respect to W_1 , hence by $(P(X), W_1)$ being compact $D_{\phi,\lambda,\alpha}(\cdot \|\nu)$ is (globally) Lipschitz continuous.

Proof. The proposition follows from Theorem 8.1.3.

To obtain the dual representations of the Moreau-Yosida approximations of $D_{\phi}(\cdot \|\nu)$ with respect to the Wasserstein-1 distance, we need the convex conjugates of the functions mapping probability measures $\xi + \omega$ to λ times the α th power of their Wasserstein-1 distance from a given probability measure μ , which by (33) is equivalent to $\lambda \|\xi + \omega - \mu\|_{KR}^{\alpha}$. We consider the cases $0 < \alpha < 1$, $\alpha = 1$ and $\alpha > 1$ separately.

We will need the following lemma.

Lemma 1. Let $\psi : \mathbb{R}_+ \to \overline{\mathbb{R}}$ be such that

$$\psi(t) = \lambda t^{\alpha}.\tag{108}$$

If $1 < \alpha \in \mathbb{R}$, then

$$\psi^{\#}(s) = (\alpha - 1)\alpha^{\frac{\alpha}{1-\alpha}}\lambda^{\frac{1}{1-\alpha}}s^{\frac{\alpha}{\alpha-1}},$$
(109)

and if $0 < \alpha < 1$, then

$$\psi^{\#}(s) = \begin{cases} 0 \text{ if } s = 0, \\ \infty \text{ otherwise.} \end{cases}$$
(110)

Proof. By (42), $\psi^{\#}(s) = \sup_{0 \le t \in \mathbb{R}} \{st - \lambda t^{\alpha}\}$. Since $\frac{\partial}{\partial t}st - \lambda t^{\alpha} = s - \alpha\lambda t^{\alpha-1}$ and $\frac{\partial}{\partial t}^2 st - \lambda t^{\alpha} = (1 - \alpha)\alpha\lambda t^{\alpha-2}$, one has an extremum at $t = \frac{s}{\alpha\lambda}^{\frac{1}{\alpha-1}}$ by letting $\frac{\partial}{\partial t} = 0$, which is a maximum if $1 < \alpha \in \mathbb{R}$ and a minimum if $0 < \alpha < 1$ by the second derivative test. This implies the proposition.

Remark 1 (The case $0 < \alpha < 1$). By (110) and (41), it holds that

$$(\xi \to \lambda \| \xi + \omega - \mu \|_{KR}^{\alpha})^* = \left(f \to \begin{cases} 0 \text{ if } \| f \|_L = 0, \\ \infty \text{ otherwise} \end{cases} \right), \tag{111}$$

which implies that the mapping $(\xi \to \lambda \| \xi + \omega - \mu \|_{KR}^{\alpha})^*$ is not convex by (45) (as it is clearly continuous and proper). Hence it is not possible to obtain a dual representation of $\inf_{\xi \in P(X)} \{D(\xi \| \nu) + \lambda W(\mu, \xi)^{\alpha}\}$ by Fenchel-Rockafellar duality when $0 < \alpha < 1$. Another approach would be to use Toland-Singer duality, which would require the mapping $(\xi \to \lambda \| \xi + \omega - \mu \|_{KR}^{\alpha})^*$ to be concave, but this is also not the case, since it is the composition of a convex and a concave nondecreasing function (Boyd & Vandenberghe, 2014, Section 3.2.3).

Proposition 12. Given $\mu, \omega \in P(X)$ and $0 < \lambda \in \mathbb{R}$, let the function $W_{\mu,\omega,\lambda,1} : (\mathcal{M}(X,0), \|.\|_{KR}) \to \mathbb{R}$ be defined by

$$W_{\mu,\omega,\lambda,1}(\xi) = \lambda \|\xi + \omega - \mu\|_{KR}.$$
(112)

Then, the function $W_{\mu,\lambda,1}$ is proper, convex and continuous, and its convex conjugate $W^*_{\mu,\omega,\lambda,1} : (Lip(X,x_0), \|.\|_L) \to \mathbb{R}$ is

$$W^*_{\mu,\omega,\lambda,1}(f) = \begin{cases} \int f d(\mu - \omega) \ if \|f\|_L \le \lambda, \\ \infty \ otherwise. \end{cases}$$
(113)

Proof. By (40),

$$(\xi \to \|\xi\|_{KR})^* = \left(f \to \begin{cases} 0 \text{ if } \|f\|_L \le 1, \\ \infty \text{ otherwise} \end{cases}\right).$$
(114)

By (38),

$$(\xi \to \|\xi + \omega - \mu\|_{KR})^* = \left(f \to \begin{cases} -\int f d(\omega - \mu) \text{ if } \|f\|_L \le 1, \\ \infty \text{ otherwise} \end{cases} \right).$$
(115)

By (36),

$$(\xi \to \lambda \|\xi + \omega - \mu\|_{KR})^* = \left(f \to \begin{cases} -\lambda \int \lambda^{-1} f d(\omega - \mu) \text{ if } \|\lambda^{-1} f\|_L \le 1, \\ \infty \text{ otherwise} \end{cases} \right), \tag{116}$$

which is equivalent to the proposed conjugate.

Proposition 13. Given $\mu, \omega \in P(X)$, $1 < \alpha \in \mathbb{R}$ and $0 < \lambda \in \mathbb{R}$, let the function $W_{\mu,\omega,\lambda,\alpha} : (\mathcal{M}(X,0), \|.\|_{KR} \to \mathbb{R}$ be defined by

$$W_{\mu,\omega,\lambda,\alpha}(\xi) = \lambda \|\xi + \omega - \mu\|_{KR}^{\alpha}.$$
(117)

Then, the function $W_{\mu,\omega,\lambda,\alpha}$ is proper, convex and continuous, and its convex conjugate $W^*_{\mu,\omega,\lambda,\alpha} : (Lip(X,x_0), \|.\|_L) \to \mathbb{R}$ is

$$W^*_{\mu,\omega,\lambda,\alpha}(f) = \int f d(\mu - \omega) + \alpha^{\frac{\alpha}{1-\alpha}} \lambda^{\frac{1}{1-\alpha}} (\alpha - 1) \|f\|_L^{\frac{\alpha}{\alpha-1}}.$$
(118)

Proof. By (109) and (41), it holds that

$$(\xi \to \lambda \|\xi\|_{KR}^{\alpha})^* = \left(f \to (\alpha - 1)\alpha^{\frac{\alpha}{1-\alpha}}\lambda^{\frac{1}{1-\alpha}} \|f\|_L^{\frac{\alpha}{\alpha-1}}\right).$$
(119)

By (38),

$$(\xi \to \lambda \|\xi + \omega - \mu\|_{KR}^{\alpha})^* = \left(f \to (\alpha - 1)\alpha^{\frac{\alpha}{1 - \alpha}} \lambda^{\frac{1}{1 - \alpha}} \|f\|_L^{\frac{\alpha}{\alpha - 1}} - \int f d(\omega - \mu) \right), \tag{120}$$

which is equivalent to the proposed conjugate.

We obtain the unconstrained variational representation of W_1 as a corollary.

Corollary 3. Given $\mu, \nu \in P(X)$, $1 < \alpha \in \mathbb{R}$ and $0 < \lambda \in \mathbb{R}$, one has the variational representation

$$\lambda W_1(\mu,\nu)^{\alpha} = \sup_{f \in Lip(X,x_0)} \left\{ \int f d\mu - \int f d\nu - \alpha^{\frac{\alpha}{1-\alpha}} \lambda^{\frac{1}{1-\alpha}} (\alpha-1) \|f\|_L^{\frac{\alpha}{\alpha-1}} \right\},\tag{121}$$

where the supremum is achieved by $\alpha \lambda W_1(\mu, \nu)^{\alpha-1} f_*$ with f_* being a Kantorovich potential of μ, ν .

Proof. The variational representation follows from the previous proposition and (45). For the other claim, one has

$$\sup_{f \in Lip(X,x_0)} \left\{ \int f d\mu - \int f d\nu - \alpha^{\frac{\alpha}{1-\alpha}} \lambda^{\frac{1}{1-\alpha}} (\alpha - 1) \|f\|_L^{\frac{\alpha}{\alpha-1}} \right\}$$

$$= \sup_{\beta \in \mathbb{R}_+, f \in Lip(X,x_0), \|f\|_L = 1} \left\{ \int \beta f d\mu - \int \beta f d\nu - \alpha^{\frac{\alpha}{1-\alpha}} \lambda^{\frac{1}{1-\alpha}} (\alpha - 1) \|\beta f\|_L^{\frac{\alpha}{\alpha-1}} \right\}$$

$$= \sup_{\beta \in \mathbb{R}_+} \left\{ \beta \sup_{f \in Lip(X,x_0), \|f\|_L = 1} \left\{ \int f d\mu - \int f d\nu \right\} - \alpha^{\frac{\alpha}{1-\alpha}} \lambda^{\frac{1}{1-\alpha}} (\alpha - 1) \beta^{\frac{\alpha}{\alpha-1}} \right\}$$

$$= \sup_{\beta \in \mathbb{R}_+} \left\{ \beta W_1(\mu, \nu) - \alpha^{\frac{\alpha}{1-\alpha}} \lambda^{\frac{1}{1-\alpha}} (\alpha - 1) \beta^{\frac{\alpha}{\alpha-1}} \right\}. \quad (122)$$

Equating the derivative with respect to β to 0 and solving the resulting equation gives $\beta = \alpha \lambda W_1(\mu, \nu)^{\alpha-1}$.

Now we have all the information we need in order to invoke Fenchel-Rockafellar duality to obtain the dual representations. **Proposition 14.** Given $\mu, \nu \in P(X)$, $0 < \lambda \in \mathbb{R}$ and a proper, convex and lower semicontinuous function $\phi : \mathbb{R} \to \overline{\mathbb{R}}$ strictly convex at 1 with $\phi(1) = 0$ and $1 \in \text{relint dom } \phi$, one has

$$\min_{\xi \in P(X)} \left\{ D_{\phi}(\xi \| \nu) + \lambda W_1(\mu, \xi) \right\} = \max_{f \in Lip(X, x_0), \|f\|_L \le \lambda} \left\{ \int f d\mu - \min_{\sup f(X) - \phi'(\infty) \le \gamma} \left\{ \int \phi^*_+ \circ (f - \gamma) d\nu + \gamma \right\} \right\}, \quad (123)$$

and for $\xi_* \in P(X)$ such that $\min_{\xi \in P(X)} \{ D_{\phi}(\xi \| \nu) + \lambda W_1(\mu, \xi) \} = D_{\phi}(\xi_* \| \nu) + \lambda W_1(\mu, \xi_*)$, there exists $f_* \in Lip(X, x_0)$ achieving the maximum such that f_* is a Csiszár potential of ξ_*, ν and λ times a Kantorovich potential of μ, ξ_* .

Proof. Substituting $D_{\phi,\nu,\omega}$, $D^*_{\phi,\nu,\omega}$, $W_{\mu,\omega,\lambda,1}$ and $W^*_{\mu,\omega,\lambda,1}$ into (51), for which the condition $\exists \xi_0 \in \text{dom } D_{\phi,\nu,\omega} \cap \text{dom } W_{\mu,\omega,\lambda,1}$ and $W_{\mu,\omega,\lambda,1}$ is continuous at ξ_0 clearly holds with $\xi_0 = \nu - \omega$, gives

$$\inf_{\xi \in \mathcal{M}(X,0)} \{ D_{\phi,\nu,\omega}(\xi) + W_{\mu,\omega,\lambda,1}(\xi) \} = \max_{f \in Lip(X,x_0)} \{ -D^*_{\phi,\nu,\omega}(f) - W^*_{\mu,\omega,\lambda,1}(-f) \}.$$
 (124)

Since $W_{\phi,\nu,\omega}(\xi) = \infty$ unless $\xi + \omega \in P(X)$, the infimum is equivalent to

$$\inf_{\xi \in P(X) - \omega} \left\{ D_{\phi}(\xi + \omega \| \nu) + \lambda \| \xi + \omega - \mu \|_{KR} \right\},\tag{125}$$

which is further equivalent to

$$\inf_{\xi \in P(X)} \{ D_{\phi}(\xi \| \nu) + \lambda \| \xi - \mu \|_{KR} \}.$$
(126)

Since $D^*_{\mu,\omega,\lambda,1}(f) = \infty$ unless $||f||_L \leq \lambda$, the maximum is equivalent to

$$\max_{f \in Lip(X,x_0), \|f\|_L \le \lambda} \left\{ -\min_{\sup f(X) - \phi'(\infty) \le \gamma} \left\{ \int \phi_+^* \circ (f - \gamma) d\nu + \gamma \right\} + \int f d\omega - \int -f d(\mu - \omega) \right\}.$$
 (127)

These, together with (33) and (107) give the claim.

By (52), there exists $f_* \in Lip(X, x_0)$ such that f_* is a Csiszár potential of ξ_*, ν and λ times a Kantorovich potential of μ, ξ_* . It achieves the maximum since $\int f_* d\mu - \min_{\sup f_*(X) - \phi'(\infty) \leq \gamma} \left\{ \int \phi_+^* \circ (f_* - \gamma) d\nu + \gamma \right\} = \int f_* d\mu - \int f_* d\xi_* + \int f_* d\xi_* - \min_{\sup f_*(X) - \phi'(\infty) \leq \gamma} \left\{ \int \phi_+^* \circ (f_* - \gamma) d\nu + \gamma \right\} = \lambda W_1(\mu, \xi_*) + D_{\phi}(\xi_* \| \nu).$

Proposition 15. Given $\mu, \nu \in P(X)$, $1 < \alpha \in \mathbb{R}$, $0 < \lambda \in \mathbb{R}$ and a proper, convex and lower semicontinuous function $\phi : \mathbb{R} \to \overline{\mathbb{R}}$ strictly convex at 1 with $\phi(1) = 0$ and $1 \in \text{relint dom } \phi$, one has

$$\min_{\xi \in P(X)} \left\{ D_{\phi}(\xi \| \nu) + \lambda W(\mu, \xi)^{\alpha} \right\} = \max_{f \in Lip(X, x_0)} \left\{ \int f d\mu - \min_{\sup f(X) - \phi'(\infty) \le \gamma} \left\{ \int \phi_+^* \circ (f - \gamma) d\nu + \gamma \right\} - (\alpha - 1) \alpha^{\frac{\alpha}{1 - \alpha}} \lambda^{\frac{1}{1 - \alpha}} \| f \|_L^{\frac{\alpha}{\alpha - 1}} \right\}, \quad (128)$$

and for $\xi_* \in P(X)$ such that $\min_{\xi \in P(X)} \{ D_{\phi}(\xi \| \nu) + \lambda W_1(\mu, \xi)^{\alpha} \} = D_{\phi}(\xi_* \| \nu) + \lambda W_1(\mu, \xi_*)^{\alpha}$, there exists $f_* \in Lip(X, x_0)$ achieving the maximum such that f_* is a Csiszár potential of ξ_*, ν and $\alpha \lambda W_1(\mu, \xi_*)^{\alpha-1}$ times a Kantorovich potential of μ, ξ_* .

Proof. Substituting $D_{\phi,\nu,\omega}$, $D^*_{\phi,\nu,\omega}$, $W_{\mu,\omega,\lambda,\alpha}$ and $W^*_{\mu,\omega,\lambda,\alpha}$ into (51), for which the condition $\exists \xi_0 \in \text{dom } D_{\phi,\nu,\omega} \cap \text{dom } W_{\mu,\omega,\lambda,\alpha}$ and $W_{\mu,\omega,\lambda,\alpha}$ is continuous at ξ_0 clearly holds with $\xi_0 = \nu - \omega$, gives

$$\inf_{\xi \in \mathcal{M}(X,0)} \left\{ D_{\phi,\nu,\omega}(\xi) + W_{\mu,\omega,\lambda,\alpha}(\xi) \right\} = \max_{f \in Lip(X,x_0)} \left\{ -D^*_{\phi,\nu,\omega}(f) - W^*_{\mu,\omega,\lambda,\alpha}(-f) \right\}.$$
(129)

Since $D_{\phi,\nu,\omega}(\xi) = \infty$ unless $\xi + \omega \in P(X)$, the infimum is equivalent to

$$\inf_{\xi \in P(X) - \omega} \left\{ D_{\phi}(\xi + \omega \| \nu) + \lambda \| \xi + \omega - \mu \|_{KR}^{\alpha} \right\},\tag{130}$$

which is further equivalent to

$$\inf_{\xi \in P(X)} \{ D_{\phi}(\xi \| \nu) + \lambda \| \xi - \mu \|_{KR}^{\alpha} \}.$$
(131)

These, together with (33) and (107) give the claim. The function $f_* \in Lip(X, x_0)$ achieving the maximum follows similarly as in the proof of the previous proposition.

8.3. Practical evaluation and differentiation of $\gamma_{\phi,\nu}(f)$

Postponing the general case as future work, we restrict our attention to evaluating $\gamma_{\phi,\nu}(f)$ when the support of ν is discrete, i.e. there exists $n \in \mathbb{R}$, $\{a_i : 1 \le i \le n, \sum_i a_i = 1\} \subset [0, 1]$ and $\operatorname{supp} \nu = \{x_i : 1 \le i \le n\} \subset X$ such that $\nu = \sum_i a_i \delta_{x_i}$ can be expressed as a convex combination of Dirac measures. Given $f \in Lip(X, x_0)$, we represent ν as a vector in the n-dimensional simplex in \mathbb{R}^n defined by (a_i) , and f as a vector in \mathbb{R}^n defined by $(f(x_i))$. The minimization problem is then reduced to

$$\min_{\max(f)-\phi'(\infty)\leq\gamma}\left\{\langle\nu,\phi_{+}^{*}(f-\gamma)\rangle+\gamma\right\}$$
(132)

with ϕ_+^* being applied element-wise and $\langle \cdot, \cdot \rangle$ being the standard dot product. The first derivative test gives

$$-\langle \nu, (\phi_+^*)'(f-\gamma) \rangle + 1 = 0.$$
(133)

Assuming the solution is unique, we define $\gamma_{\phi,\nu}(f)$ implicitly as

$$\gamma_{\phi,\nu}(f) = \gamma' \iff -\langle \nu, (\phi_+^*)'(f - \gamma') \rangle + 1 = 0.$$
(134)

Two cases offer closed-form solution. One is the Kullback-Leibler divergence $\phi(x) = x \log x$, for which we get $\gamma_{(x \to x \log x),\nu}(f) = \log \langle \nu, e^f \rangle$, i.e. the term from the Donsker-Varadhan formula. The other is the total variation divergence corresponding to $\phi(x) = |x - 1|$, for which the mapping

$$\gamma \to \langle \nu, -\chi_{(-\infty, -1)}(f - \gamma) + (f - \gamma)\chi_{[-1, 1]}(f - \gamma) \rangle + \gamma$$
(135)

is nondecreasing for $\gamma \geq \max(f) - 1$ by its derivative $\langle \nu, -\chi_{[-1,1]}(f-\gamma) \rangle + 1$ being nonnegative, implying that the optimal value is $\gamma_{(x \to |x-1|),\nu}(f) = \max(f) - 1$, and the conjugate is $D_{(x \to |x-1|)}(f \| \nu) = \langle \nu, -\chi_{(-\infty,-1)}(f - \max(f) + 1) + (f - \max(f) + 1)\chi_{[-1,1]}(f - \max(f) + 1) \rangle + \max(f) - 1$.

For other choices of ϕ considered, it seems that no closed-form solution is available. Instead, we calculate $\gamma_{\phi,\nu}(f)$ by Newton's method, for which we need the derivative of the function whose zeroes define the values of $\gamma_{\phi,\nu}(f)$, which is

$$\langle \nu, (\phi_+^*)''(f-\gamma) \rangle. \tag{136}$$

Then, Newton's method suggests that the iteration

$$\gamma_{n+1} = \gamma_n - \frac{-\langle \nu, (\phi_+^*)'(f-\gamma) \rangle + 1}{\langle \nu, (\phi_+^*)''(f-\gamma) \rangle}$$
(137)

converges to $\gamma_{\phi,\nu}(f)$. For the initial value, the choice $\gamma_0 = \max f - \phi'(\infty) + \epsilon$ works for some small $\epsilon > 0$ tuned manually if $\phi'(\infty) < \infty$. The the other cases with $\phi'(\infty) = \infty$, we found the choice $\gamma_0 = \langle \nu, f \rangle$ to work just fine.

To integrate this implicit function into automatic differentiation frameworks, we need to be able to compute the gradients $\nabla_f \gamma_{\phi,\nu}(f)$. Instead of the potentially unstable method of backpropagating through the iterations of Newton's method, we use the implicit function theorem (Tao, 2016), which tells us that

$$\nabla_f \gamma_{\phi,\nu}(f) = \frac{-\nabla_f(-\langle \nu, (\phi_+^*)'(f-\gamma')\rangle + 1)}{\frac{d}{d\gamma}(-\langle \nu, (\phi_+^*)'(f-\gamma')\rangle + 1)} = \frac{\nu \odot (\phi_+^*)''(f-\gamma)}{\langle \nu, (\phi_+^*)''(f-\gamma)\rangle}$$
(138)

with \odot denoting element-wise product in \mathbb{R}^n . Notice that in all cases, $\nabla_f \gamma_{\phi,\nu}(f)$ is in the standard simplex, e.g. for the Kullback-Leibler divergence one has the softmax $\nabla_f \gamma_{\phi,\nu}(f) = \frac{\nu \odot e^f}{\langle \nu, e^f \rangle}$ as the gradient.

We implemented the implicit functions with Newton's method in the forward pass and the backward pass formula given by the implicit function theorem for the Kullback-Leibler, reverse Kullback-Leibler, χ^2 , reverse χ^2 , squared Hellinger, Jensen-Shannon, Jeffreys and triangular discrimination divergences. To test the validity of the approach, we optimized an $f \in \mathbb{R}^n$ with gradient descent to maximize the corresponding variational formulas with random categorical distributions μ, ν over an alphabet of size n, and found that the resulting value for the divergences matched that of the closed-form solution with high accuracy. We found that calculating $\gamma_{\phi,\nu}(f)$ as $\gamma_{\phi,\nu}(f) = \gamma_{\phi,\nu}(f - \max(f)) + \max(f)$ is beneficial to avoid numerical instabilities, which can be seen as a generalization of the log-sum-exp trick. We detail the functions derived from ϕ corresponding to the listed f-divergences defined by functions of Legendre type below, as well as the corresponding Csiszár potentials.

8.3.1. KULLBACK-LEIBLER DIVERGENCE

$$\phi_{+}(x) = \begin{cases} x \log(x) - x + 1 \text{ if } x \ge 0, \\ \infty \text{ otherwise.} \end{cases}$$
(139)

$$\partial \phi_{+}(x) = \begin{cases} \{\log(x)\} \text{ if } x > 0, \\ \emptyset \text{ otherwise.} \end{cases}$$
(140)

$$\phi'(\infty) = \infty \tag{141}$$

$$\phi_+^*(x) = e^x - 1 \tag{142}$$

$$\phi_{+}^{*\,\prime}(x) = e^x \tag{143}$$

$$\phi_{+}^{*\,''}(x) = e^x \tag{144}$$

$$f_*(x) + C = \log\left(\frac{d\mu_c}{d\nu}(x)\right)$$
 almost everywhere with respect to μ_c (145)

8.3.2. Reverse Kullback-Leibler divergence

$$\phi_{+}(x) = \begin{cases} x - 1 - \log(x) \text{ if } x \ge 0, \\ \infty \text{ otherwise.} \end{cases}$$
(146)

$$\partial \phi_{+}(x) = \begin{cases} \left\{ \frac{x-1}{x} \right\} & \text{if } x > 0, \\ \emptyset & \text{otherwise.} \end{cases}$$
(147)

$$\phi'(\infty) = 1 \tag{148}$$

$$\phi_{+}^{*}(x) = \begin{cases} -\log(1-x) \text{ if } x \le 1, \\ \infty \text{ otherwise.} \end{cases}$$
(149)

$$\phi_{+}^{*\,\prime}(x) = \frac{1}{1-x} \tag{150}$$

$$\phi_{+}^{*\,''}(x) = \frac{1}{(1-x)^2} \tag{151}$$

$$f_*(x) + C = \begin{cases} \frac{d\mu_c}{d\nu}(x) - 1\\ \frac{d\mu_c}{d\nu}(x) \end{cases} \text{ almost everywhere with respect to } \mu_c, \\ 1 \text{ if } x \in \operatorname{supp}(\mu_s) \end{cases}$$
(152)

8.3.3. χ^2 divergence

$$\phi_{+}(x) = \begin{cases} (x-1)^{2} \text{ if } x \ge 0, \\ \infty \text{ otherwise.} \end{cases}$$
(153)

$$\partial \phi_{+}(x) = \begin{cases} \{2x - 2\} \text{ if } x \ge 0, \\ \emptyset \text{ otherwise.} \end{cases}$$
(154)

$$\phi'(\infty) = \infty \tag{155}$$

$$\phi_{+}^{*}(x) = \begin{cases} \frac{1}{4}x^{2} + x \text{ if } x \ge -2, \\ -1 \text{ otherwise.} \end{cases}$$
(156)

$$\phi_{+}^{*\,\prime}(x) = \begin{cases} \frac{1}{2}x + 1 \text{ if } x \ge -2, \\ 0 \text{ otherwise.} \end{cases}$$
(157)

$$\phi_{+}^{*\,''}(x) = \begin{cases} \frac{1}{2} \text{ if } x \ge -2, \\ 0 \text{ otherwise.} \end{cases}$$
(158)

$$f_*(x) + C = 2\frac{d\mu_c}{d\nu}(x) - 2$$
 almost everywhere with respect to μ_c (159)

8.3.4. Reverse χ^2 divergence

$$\phi_{+}(x) = \begin{cases} \frac{1}{x} + x - 2 \text{ if } x \ge 0, \\ \infty \text{ otherwise.} \end{cases}$$
(160)

$$\partial \phi_{+}(x) = \begin{cases} \left\{ 1 - \frac{1}{x^{2}} \right\} & \text{if } x > 0, \\ \emptyset & \text{otherwise.} \end{cases}$$
(161)

$$\phi'(\infty) = 1 \tag{162}$$

$$\phi_{+}^{*}(x) = \begin{cases} 2 - 2\sqrt{1-x} \text{ if } x \le 1, \\ \infty \text{ otherwise.} \end{cases}$$
(163)

$$\phi_{+}^{*}'(x) = \frac{1}{\sqrt{1-x}} \tag{164}$$

$$\phi_{+}^{*\,''}(x) = \frac{1}{2\sqrt{1-x^3}} \tag{165}$$

$$f_*(x) + C = \begin{cases} 1 - \frac{1}{\left(\frac{d\mu_c}{d\nu}(x)\right)^2} \text{ almost everywhere with respect to } \mu_c, \\ 1 \text{ if } x \in \text{supp}(\mu_s) \end{cases}$$
(166)

8.3.5. Squared Hellinger divergence

$$\phi_{+}(x) = \begin{cases} (\sqrt{x} - 1)^{2} \text{ if } x \ge 0, \\ \infty \text{ otherwise.} \end{cases}$$
(167)

$$\partial \phi_{+}(x) = \begin{cases} \left\{ 1 - \frac{1}{\sqrt{x}} \right\} & \text{if } x > 0, \\ \emptyset & \text{otherwise.} \end{cases}$$
(168)

$$\phi'(\infty) = 1 \tag{169}$$

$$\phi_{+}^{*}(x) = \begin{cases} \frac{x}{1-x} \text{ if } x \le 1, \\ \infty \text{ otherwise.} \end{cases}$$
(170)

$$\phi_{+}^{*\,\prime}(x) = \frac{1}{(1-x)^2} \tag{171}$$

$$\phi_{+}^{*\,\prime\prime}(x) = \frac{2}{(1-x)^3} \tag{172}$$

$$f_*(x) + C = \begin{cases} 1 - \frac{1}{\sqrt{\frac{d\mu_c}{d\nu}(x)}} \text{ almost everywhere with respect to } \mu_c, \\ 1 \text{ if } x \in \text{supp}(\mu_s) \end{cases}$$
(173)

8.3.6. JENSEN-SHANNON DIVERGENCE

$$\phi_{+}(x) = \begin{cases} x \log(x) - (x+1) \log(\frac{x+1}{2}) \text{ if } x \ge 0, \\ \infty \text{ otherwise.} \end{cases}$$
(174)

$$\partial \phi_{+}(x) = \begin{cases} \{\log(x) - \log(x+1) + \log(2)\} & \text{if } x > 0, \\ \emptyset & \text{otherwise.} \end{cases}$$
(175)

$$\phi'(\infty) = \log(2) \tag{176}$$

$$\phi_{+}^{*}(x) = \begin{cases} -\log(2 - e^{x}) \text{ if } x \le \log(2), \\ \infty \text{ otherwise.} \end{cases}$$
(177)

$$\phi_{+}^{*}'(x) = \frac{1}{2e^{-x} - 1} \tag{178}$$

$$\phi_{+}^{*\,''}(x) = \frac{2e^x}{(e^x - 2)^2} \tag{179}$$

$$f_*(x) + C = \begin{cases} \log\left(\frac{d\mu_c}{d\nu}(x)\right) - \log\left(\frac{d\mu_c}{d\nu}(x) + 1\right) + \log(2) \text{ almost everywhere with respect to } \mu_c, \\ \log(2) \text{ if } x \in \operatorname{supp}(\mu_s) \end{cases}$$
(180)

8.3.7. JEFFREYS DIVERGENCE

$$\phi_{+}(x) = \begin{cases} (x-1)\log(x) \text{ if } x \ge 0, \\ \infty \text{ otherwise.} \end{cases}$$
(181)

$$\partial \phi_+(x) = \begin{cases} \left\{ \log(x) - \frac{1}{x} + 1 \right\} & \text{if } x > 0, \\ \emptyset & \text{otherwise.} \end{cases}$$
(182)

$$\phi'(\infty) = \infty \tag{183}$$

$$\phi_{+}^{*}(x) = x + W(e^{1-x}) + \frac{1}{W(e^{1-x})} - 2$$
(184)

$$\phi_{+}^{*}'(x) = \frac{1}{W(e^{1-x})} \tag{185}$$

$$\phi_{+}^{*\,''}(x) = \frac{1}{W(e^{1-x})} - \frac{1}{W(e^{1-x}) + 1}$$
(186)

$$f_*(x) + C = \log\left(\frac{d\mu_c}{d\nu}(x)\right) - \frac{1}{\frac{d\mu_c}{d\nu}(x)} + 1 \text{ almost everywhere with respect to } \mu_c$$
(187)

W denotes the principal branch of the Lambert W function, also called the product logarithm, defined implicitly by the relation $W(x)e^{W(x)} = x$. Similarly to the proposed conjugates, it is computed by Newton's method and its gradient by the implicit function theorem.

8.3.8. TRIANGULAR DISCRIMINATION DIVERGENCE

$$\phi_{+}(x) = \begin{cases} \frac{(x-1)^2}{x+1} & \text{if } x \ge 0, \\ \infty & \text{otherwise.} \end{cases}$$
(188)

$$\partial \phi_{+}(x) = \begin{cases} \left\{ \frac{(x-1)(x+3)}{(x+1)^{2}} \right\} & \text{if } x \ge 0, \\ \emptyset & \text{otherwise.} \end{cases}$$
(189)

$$\phi'(\infty) = 1 \tag{190}$$

$$\phi_{+}^{*}(x) = \begin{cases} -1 \text{ if } x < -3, \\ (\sqrt{1-x}-1)(\sqrt{1-x}-3) \text{ if } -3 \le x \le 1, \\ \infty \text{ otherwise.} \end{cases}$$
(191)

$$\phi_{+}^{*'}(x) = \begin{cases} 0 \text{ if } x < -3, \\ \frac{2}{\sqrt{1-x}} - 1 \text{ if } -3 \le x \le 1 \end{cases}$$
(192)

$$\phi_{+}^{*\,''}(x) = \begin{cases} 0 \text{ if } x < -3, \\ \frac{1}{(\sqrt{1-x})^3} \text{ if } -3 \le x \le 1 \end{cases}$$
(193)

$$f_*(x) + C = \begin{cases} \frac{\left(\frac{d\mu_c}{d\nu}(x) - 1\right) \left(\frac{d\mu_c}{d\nu}(x) + 3\right)}{\left(\frac{d\mu_c}{d\nu}(x) + 1\right)^2} \text{ almost everywhere with respect to } \mu_c, \\ 1 \text{ if } x \in \operatorname{supp}(\mu_s) \end{cases}$$
(194)

8.4. Experiments

8.4.1. MY f-GAN ON CIFAR-10

The implementation was done in TensorFlow based on the official codebase of Adler & Lunz (2018), with the critic and generator architectures being faithful reimplementations of the residual architecture from Gulrajani et al. (2017). We used both the train and test parts of the CIFAR-10 dataset with randomly flipping images and adding uniform noise as augmentation. Minibatch size was 128, which for the critic included 64 real and 64 generated samples. We used the ADAM optimizer with parameters $\beta_1 = 0$, $\beta_2 = 0.9$ and constant learning rate 2×10^{-4} , and trained the model for 100000 iterations with 5 gradient descent step per iteration for the critic, and 1 for the generator. We monitored performance by evaluating the Inception Score at every 1000 iteration during training on 10000 generated samples, and once at the end of training. We applied exponential moving averaging to the generator weights θ_g with coefficient 0.9999. The Lambert W function implementation was based on https://github.com/jackd/lambertw. Trainings were done on GeForce 2080Ti GPUs running at 200-250W and took around 12 hours to complete, leading to an estimated 2.4-3kWh power consumption. Computing the conjugates via the proposed algorithm introduced a computational overhead that induced 15% longer training time at worst (for the Jeffreys divergence) compared to closed-form conjugates. Due to the large number of hyperparameters, we ran each setting once. Generated images can be seen in Figure 5, Figure 6, Figure 7, Figure 8, Figure 9, Figure 10, Figure 11, Figure 12, Figure 13 and Figure 14, with missing images denoting failed training such as discriminator collapse or numerical instabilities.

8.4.2. 1D GAUSSIAN DISTRIBUTIONS

For a pair $\mathcal{N}(\mu_1, \sigma_1), \mathcal{N}(\mu_2, \sigma_2)$ of 1-dimensional Gaussian distributions, the corresponding probability distribution functions are the Radon-Nikodym derivatives with respect to the Lebesgue measure, hence by the chain rule one has

$$\frac{d\mathcal{N}(\mu_1,\sigma_1)}{d\mathcal{N}(\mu_2,\sigma_2)}(x) = \frac{\sigma_2}{\sigma_1} e^{\frac{1}{2}\left(\left(\frac{x-\mu_2}{\sigma_2}\right)^2 - \left(\frac{x-\mu_1}{\sigma_1}\right)^2\right)},$$
(195)

so that Csiszár potentials can be calculated in closed form if ϕ is of Legendre type. We have implemented a toy example to demonstrate that the proposed algorithm for calculating the conjugates enables training neural networks to approximate



Figure 5. MY f-GAN generated images with D_{ϕ} being the Kullback-Leibler divergence

f-divergences based on the tight variational representations in the sense that the trained neural network closely approximates the corresponding Csiszár potential in the case of 1-dimensional Gaussian distributions. Results are visualized in Figure 15, showing the probability distribution functions of the two Gaussians, the exact Csiszár potential and the output of the trained neural network. For the Jeffreys divergence, numerical instabilities prevented us from obtaining the desired result, so that only the exact Csiszár potential is visualized. Close approximation of the exact Csiszár potential is evident in areas of higher density. It must be emphasized that the neural network is not explicitly trained to approximate the Csiszár potential, only implicitly, by maximizing the tight variational formula, so that this experiment could be seen as a validation of both the algorithm for computing the conjugates and of the characterization of Csiszár potentials.

8.4.3. CATEGORICAL DISTRIBUTIONS

For a pair μ, ν of categorical distributions over an alphabet of size n, the potential f can be considered an element of \mathbb{R}^n . We implemented a toy example to approximate D_{ϕ} in this case by optimizing an $f \in \mathbb{R}^n$ to maximize the formula inside the supremum in the tight variational representation, and found that the approximation accurately recovers the exact value of the divergence to at least 4 decimals, except for the reverse χ^2 divergence for which the approximation procedure is slightly less accurate. For the Kullback-Leibler and total variation divergences even though the conjugates are available in closed form, we implemented the proposed algorithm as well, and found that it works as well as the closed forms in this scenario. The generalized log-sum-exp trick is necessary in some cases to stabilize the algorithm.



Figure 6. MYf-GAN generated images with D_{ϕ} being the reverse Kullback-Leibler divergence



Figure 7. MYf-GAN generated images with D_ϕ being the χ^2 divergence



Figure 8. MYf-GAN generated images with D_{ϕ} being the reverse χ^2 divergence



Figure 9. MYf-GAN generated images with D_ϕ being the squared Hellinger divergence



Figure 10. MY f-GAN generated images with D_{ϕ} being the Jensen-Shannon divergence



Figure 11. MY f-GAN generated images with D_{ϕ} being the Jeffreys divergence



Figure 12. MYf-GAN generated images with D_{ϕ} being the triangular discrimination divergence



Figure 13. MYf-GAN generated images with D_{ϕ} being the total variation divergence



Figure 14. MY f-GAN generated images with D_{ϕ} being the trivial divergence



Figure 15. Csiszár potentials and trained critics of $\mathcal{N}(-1, 0.3), \mathcal{N}(0.5, 0.6)$