
A Convex Duality Framework for GANs

Farzan Farnia*
farnia@stanford.edu

David Tse*
dntse@stanford.edu

Abstract

Generative adversarial network (GAN) is a minimax game between a generator mimicking the true model and a discriminator distinguishing the samples produced by the generator from the real training samples. Given an unconstrained discriminator able to approximate any function, this game reduces to finding the generative model minimizing a divergence measure, e.g. the Jensen-Shannon (JS) divergence, to the data distribution. However, in practice the discriminator is constrained to be in a smaller class \mathcal{F} such as neural nets. Then, a natural question is how the divergence minimization interpretation changes as we constrain \mathcal{F} . In this work, we address this question by developing a convex duality framework for analyzing GANs. For a convex set \mathcal{F} , this duality framework interprets the original GAN formulation as finding the generative model with minimum JS-divergence to the distributions penalized to match the moments of the data distribution, with the moments specified by the discriminators in \mathcal{F} . We show that this interpretation more generally holds for f-GAN and Wasserstein GAN. As a byproduct, we apply the duality framework to a hybrid of f-divergence and Wasserstein distance. Unlike the f-divergence, we prove that the proposed hybrid divergence changes continuously with the generative model, which suggests regularizing the discriminator's Lipschitz constant in f-GAN and vanilla GAN. We numerically evaluate the power of the suggested regularization schemes for improving GAN's training performance.

1 Introduction

Learning a probability model from data samples is a fundamental task in unsupervised learning. The recently developed generative adversarial network (GAN) [1] leverages the power of deep neural networks to successfully address this task across various domains [2]. In contrast to traditional methods of parameter fitting like maximum likelihood estimation, the GAN approach views the problem as a *game* between a *generator* G whose goal is to generate fake samples that are close to the real data training samples and a *discriminator* D whose goal is to distinguish between the real and fake samples. The generator creates the fake samples by mapping from random noise input.

The following minimax problem is the original GAN problem, also called *vanilla GAN*, introduced in [1]

$$\min_{G \in \mathcal{G}} \max_{D \in \mathcal{F}} \mathbb{E}[\log D(\mathbf{X})] + \mathbb{E}[\log(1 - D(G(\mathbf{Z})))] \quad (1)$$

Here \mathbf{Z} denotes the generator's noise input, \mathbf{X} represents the random vector for the real data distributed as $P_{\mathbf{X}}$, and \mathcal{G} and \mathcal{F} respectively represent the generator and discriminator function sets. Implementing this minimax game using deep neural network classes \mathcal{G} and \mathcal{F} has led to the state-of-the-art generative model for many different tasks.

To shed light on the probabilistic meaning of vanilla GAN, [1] shows that given an unconstrained discriminator D , i.e. if \mathcal{F} contains all possible functions, the minimax problem (1) will reduce to

$$\min_{G \in \mathcal{G}} \text{JSD}(P_{\mathbf{X}}, P_G(\mathbf{z})), \quad (2)$$

*Department of Electrical Engineering, Stanford University.

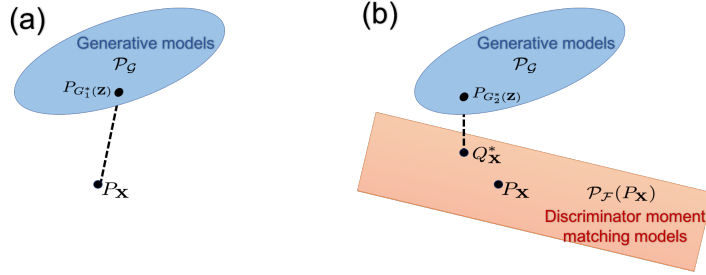


Figure 1: (a) Divergence minimization in (2) between $P_{\mathbf{X}}$ and generative models $\mathcal{P}_{\mathcal{G}}$ for unconstrained \mathcal{F} , (b) Divergence minimization in (3) between generative models $\mathcal{P}_{\mathcal{G}}$ and discriminator moment matching models $\mathcal{P}_{\mathcal{F}}(P_{\mathbf{X}})$.

where JSD denotes the Jensen-Shannon (JS) divergence. The optimization problem (2) can be interpreted as finding the closest generative model to the data distribution $P_{\mathbf{X}}$ (Figure 1a), where distance is measured using the JS-divergence. Various GAN formulations were later proposed by changing the divergence measure in (2): f-GAN [3] generalizes vanilla GAN by minimizing a general f-divergence; Wasserstein GAN (WGAN) [4] considers the first-order Wasserstein (the earth-mover’s) distance; MMD-GAN [5, 6, 7] considers the maximum mean discrepancy; Energy-based GAN [8] minimizes the total variation distance as discussed in [4]; Quadratic GAN [9] finds the distribution minimizing the second-order Wasserstein distance.

However, GANs trained in practice differ from this minimum divergence formulation, since their discriminator is not optimized over an unconstrained set and is constrained to smaller classes such as neural nets. As shown in [10], constraining the discriminator is in fact necessary to guarantee good generalization properties for GAN’s learned model. Then, how does the minimum divergence interpretation (2) change as we constrain \mathcal{F} ? A standard approach used in [10, 11] is to view the maximum discriminator objective as an \mathcal{F} -based distance between distributions. For unconstrained \mathcal{F} , the \mathcal{F} -based distance reduces to the original divergence measure, e.g. the JS-divergence in vanilla GAN.

While \mathcal{F} -based distances have been shown to be useful for analyzing GAN’s generalization properties [10], their connection to the original divergence measure remains unclear for a constrained \mathcal{F} . Then, what is the interpretation of GAN minimax game with a constrained discriminator? In this work, we address this question by interpreting the dual problem to the discriminator optimization. To analyze the dual problem, we develop a convex duality framework for general divergence minimization problems. We apply the duality framework to the f-divergence and optimal transport cost families, providing interpretation for f-GAN, including vanilla GAN minimizing JS-divergence, and Wasserstein GAN.

Specifically, we generalize the interpretation for unconstrained \mathcal{F} in (2) to any linear space discriminator set \mathcal{F} . For this class of discriminator sets, we interpret vanilla GAN as the following JS-divergence minimization between two sets of probability distributions, the set of generative models and the set of discriminator moment-matching distributions (Figure 1b),

$$\min_{G \in \mathcal{G}} \min_{Q \in \mathcal{P}_{\mathcal{F}}(P_{\mathbf{X}})} \text{JSD}(P_G(\mathbf{z}), Q). \quad (3)$$

Here $\mathcal{P}_{\mathcal{F}}(P_{\mathbf{X}})$ contains any distribution Q satisfying the moment matching constraint $\mathbb{E}_Q[D(\mathbf{X})] = \mathbb{E}_P[D(\mathbf{X})]$ for all discriminator D ’s in \mathcal{F} . More generally, we show that a similar interpretation applies to GANs trained over any convex discriminator set \mathcal{F} . We further discuss the application of our duality framework to neural net discriminators with bounded Lipschitz constant. While a set of neural network functions is not necessarily convex, we prove any convex combination of Lipschitz-bounded neural nets can be approximated by uniformly combining boundedly-many neural nets. This result applied to our duality framework suggests considering a uniform mixture of multiple neural nets as the discriminator.

As a byproduct, we apply the duality framework to the minimum sum hybrid of f-divergence and the first-order Wasserstein (W_1) distance, e.g. the following hybrid of JS-divergence and W_1 distance:

$$d_{\text{JSD}, W_1}(P_1, P_2) := \min_Q W_1(P_1, Q) + \text{JSD}(Q, P_2). \quad (4)$$

We prove that this hybrid divergence enjoys a continuous behavior in distribution P_1 . Therefore, the hybrid divergence provides a remedy for the discontinuous behavior of the JS-divergence when optimizing the generator parameters in vanilla GAN. [4] observes this issue with the JS-divergence in vanilla GAN and proposes to instead minimize the continuously-changing W_1 distance in WGAN. However, as empirically demonstrated in [12] vanilla GAN with Lipschitz-bounded discriminator remains the state-of-the-art method for training deep generative models in several benchmark tasks. Here, we leverage our duality framework to prove that the hybrid d_{JSD, W_1} , which possesses the same continuity property as in W_1 distance, is in fact the divergence measure minimized in vanilla GAN with 1-Lipschitz discriminator. Our analysis hence provides an explanation for why regularizing the discriminator’s Lipschitz constant via gradient penalty [13] or spectral normalization [12] improves the training performance in vanilla GAN. We then extend our focus to the hybrid of f-divergence and the second-order Wasserstein (W_2) distance. In this case, we derive the f-GAN (e.g. vanilla GAN) problem with its discriminator being adversarially trained using Wasserstein risk minimization [14]. We numerically evaluate the power of these families of hybrid divergences in training vanilla GAN.

2 Divergence Measures

2.1 Jensen-Shannon divergence

The Jensen-Shannon divergence is defined in terms of the KL-divergence (denoted by KL) as

$$\text{JSD}(P, Q) := \frac{1}{2} \text{KL}(P \| M) + \frac{1}{2} \text{KL}(Q \| M)$$

where $M = \frac{P+Q}{2}$ is the mid-distribution between P and Q . Unlike the KL-divergence, the JS-divergence is symmetric $\text{JSD}(P, Q) = \text{JSD}(Q, P)$ and bounded $0 \leq \text{JSD}(P, Q) \leq 1$.

2.2 f-divergence

The f-divergence family [15] generalizes the KL and JS divergence measures. Given a convex lower semicontinuous function f with $f(1) = 0$, the f-divergence d_f is defined as

$$d_f(P, Q) := \mathbb{E}_P \left[f \left(\frac{q(\mathbf{X})}{p(\mathbf{X})} \right) \right] = \int p(\mathbf{x}) f \left(\frac{q(\mathbf{x})}{p(\mathbf{x})} \right) d\mathbf{x}. \quad (5)$$

Here \mathbb{E}_P denotes expectation over distribution P and p, q denote the density functions for distributions P, Q , respectively. The KL-divergence and the JS-divergence are members of the f-divergence family, corresponding to respectively $f_{\text{KL}}(t) = t \log t$ and $f_{\text{JSD}}(t) = \frac{t}{2} \log t - \frac{t+1}{2} \log \frac{t+1}{2}$.

2.3 Optimal transport cost, Wasserstein distance

The optimal transport cost for cost function $c(\mathbf{x}, \mathbf{x}')$, which we denote by OT_c , is defined as

$$\text{OT}_c(P, Q) := \inf_{M \in \Pi(P, Q)} \mathbb{E}[c(\mathbf{X}, \mathbf{X}')], \quad (6)$$

where $\Pi(P, Q)$ contains all couplings with marginals P, Q . The Kantorovich duality [16] shows that for a non-negative lower semi-continuous cost c ,

$$\text{OT}_c(P, Q) = \max_{D \text{ c-concave}} \mathbb{E}_P[D(\mathbf{X})] - \mathbb{E}_Q[D^c(\mathbf{X})], \quad (7)$$

where we use D^c to denote D ’s c-transform defined as $D^c(\mathbf{x}) := \sup_{\mathbf{x}'} D(\mathbf{x}') - c(\mathbf{x}, \mathbf{x}')$ and call D c-concave if D is the c-transform of a valid function. Considering the norm-based cost $c_q(\mathbf{x}, \mathbf{x}') = \|\mathbf{x} - \mathbf{x}'\|^q$ with $q \geq 1$, the q th order Wasserstein distance W_q is defined based on the c_q optimal transport cost as

$$W_q(P, Q) := \text{OT}_{c_q}(P, Q)^{1/q} = \inf_{M \in \Pi(P, Q)} \mathbb{E}[\|\mathbf{X} - \mathbf{X}'\|^q]^{1/q}. \quad (8)$$

An important special case is the first-order Wasserstein (W_1) distance corresponding to the difference norm cost $c_1(\mathbf{x}, \mathbf{x}') = \|\mathbf{x} - \mathbf{x}'\|$. Given cost function c_1 , a function D is c-concave if and only if D is 1-Lipschitz, and the c-transform $D^c = D$ for any 1-Lipschitz D . Therefore, the Kantorovich duality (7) implies that

$$W_1(P, Q) = \max_{D \text{ 1-Lipschitz}} \mathbb{E}_P[D(\mathbf{X})] - \mathbb{E}_Q[D(\mathbf{X})]. \quad (9)$$

Another notable special case is the second-order Wasserstein (W_2) distance, corresponding to the difference norm-squared cost $c_2(\mathbf{x}, \mathbf{x}') = \|\mathbf{x} - \mathbf{x}'\|^2$.

3 Divergence minimization in GANs: a convex duality framework

In this section, we develop a convex duality framework for analyzing divergence minimization problems conditioned to moment-matching constraints. Our framework generalizes the duality framework developed in [17] for the f-divergence family.

For a general divergence measure $d(P, Q)$, we define d 's conjugate over distribution P , which we denote by d_P^* , as the following mapping from real-valued functions of \mathbf{X} to real numbers

$$d_P^*(D) := \sup_Q \mathbb{E}_Q[D(\mathbf{X})] - d(P, Q). \quad (10)$$

Here the supremum is over all distributions on \mathbf{X} with support set \mathcal{X} . We later show the following theorem, which is based on the above definition, recovers various well-known GAN formulations, when applied to divergence measures discussed in Section 2.

Theorem 1. *Suppose divergence $d(P, Q)$ is non-negative, lower semicontinuous and convex in distribution Q . Consider a convex set of continuous functions \mathcal{F} and assume support set \mathcal{X} is compact. Then,*

$$\begin{aligned} & \min_{G \in \mathcal{G}} \max_{D \in \mathcal{F}} \mathbb{E}_{P_{\mathbf{X}}}[D(\mathbf{X})] - d_{P_{G(\mathbf{Z})}}^*(D) \\ &= \min_{G \in \mathcal{G}} \min_Q \left\{ d(P_{G(\mathbf{Z})}, Q) + \max_{D \in \mathcal{F}} \{ \mathbb{E}_{P_{\mathbf{X}}}[D(\mathbf{X})] - \mathbb{E}_Q[D(\mathbf{X})] \} \right\}. \end{aligned} \quad (11)$$

Proof. We defer the proof to the Appendix. \square

Theorem 1 interprets (11)'s LHS minimax problem as searching for the closest generative model to the distributions penalized to share the same moments specified by \mathcal{F} with $P_{\mathbf{X}}$. The following corollary of Theorem 1 shows if we further assume that \mathcal{F} is a linear space, then the penalty term penalizing moment mismatches can be moved to the constraints. This reduction reveals a divergence minimization problem between generative models and the following set $\mathcal{P}_{\mathcal{F}}(P)$ which we call the set of discriminator moment matching distributions,

$$\mathcal{P}_{\mathcal{F}}(P) := \{ Q : \forall D \in \mathcal{F}, \mathbb{E}_Q[D(\mathbf{X})] = \mathbb{E}_P[D(\mathbf{X})] \}. \quad (12)$$

Corollary 1. *In Theorem 1 suppose \mathcal{F} is further a linear space, i.e. for any $D_1, D_2 \in \mathcal{F}$ and $\lambda \in \mathbb{R}$ we have $D_1 + \lambda D_2 \in \mathcal{F}$. Then,*

$$\min_{G \in \mathcal{G}} \max_{D \in \mathcal{F}} \mathbb{E}_{P_{\mathbf{X}}}[D(\mathbf{X})] - d_{P_{G(\mathbf{Z})}}^*(D) = \min_{G \in \mathcal{G}} \min_{Q \in \mathcal{P}_{\mathcal{F}}(P_{\mathbf{X}})} d(P_{G(\mathbf{Z})}, Q). \quad (13)$$

In next section, we apply this duality framework to divergence measures discussed in Section 2 and show how to derive various GAN problems through the developed framework.

4 Duality framework applied to different divergence measures

4.1 f-divergence: f-GAN and vanilla GAN

Theorem 2 shows the application of Theorem 1 to f-divergences. We use f^* to denote f 's convex-conjugate [18], defined as $f^*(u) := \sup_t ut - f(t)$. Note that Theorem 2 applies to any f-divergence d_f with non-decreasing convex-conjugate f^* , which holds for all f-divergence examples discussed in [3] with the only exception of Pearson χ^2 -divergence.

Theorem 2. *Consider f-divergence d_f where the corresponding f has a non-decreasing convex-conjugate f^* . In addition to Theorem 1's assumptions, suppose \mathcal{F} is closed to adding constant functions, i.e. $D + \lambda \in \mathcal{F}$ if $D \in \mathcal{F}$, $\lambda \in \mathbb{R}$. Then, the minimax problem in the LHS of (11) and (13), will reduce to*

$$\min_{G \in \mathcal{G}} \max_{D \in \mathcal{F}} \mathbb{E}[D(\mathbf{X})] - \mathbb{E}[f^*(D(G(\mathbf{Z})))]. \quad (14)$$

Proof. We defer the proof to the Appendix. \square

The minimax problem (14) is in fact the f-GAN problem [3]. Theorem 2 hence reveals that f-GAN searches for the generative model minimizing f-divergence to the distributions matching moments specified by \mathcal{F} to the moments of true distribution.

Example 1. Consider the JS-divergence, i.e. f -divergence corresponding to $f_{\text{JSD}}(t) = \frac{t}{2} \log t - \frac{t+1}{2} \log \frac{t+1}{2}$. Then, (14) up to additive and multiplicative constants reduces to

$$\min_{G \in \mathcal{G}} \max_{D \in \mathcal{F}} \mathbb{E}[D(\mathbf{X})] + \mathbb{E}[\log(1 - \exp(D(G(\mathbf{Z})))]. \quad (15)$$

Moreover, if for function set $\tilde{\mathcal{F}}$ the corresponding $\mathcal{F} = \{D : D(\mathbf{x}) = -\log(1 + \exp(\tilde{D}(\mathbf{x})))\}$, $\tilde{D} \in \tilde{\mathcal{F}}$ is a convex set, then (15) will reduce to the following minimax game which is the vanilla GAN problem (1) with sigmoid activation applied to the discriminator output,

$$\min_{G \in \mathcal{G}} \max_{\tilde{D} \in \tilde{\mathcal{F}}} \mathbb{E}\left[\log \frac{1}{1 + \exp(\tilde{D}(\mathbf{X}))}\right] + \mathbb{E}\left[\log \frac{\exp(\tilde{D}(\mathbf{X}))}{1 + \exp(\tilde{D}(\mathbf{X}))}\right]. \quad (16)$$

4.2 Optimal Transport Cost: Wasserstein GAN

Theorem 3. Let divergence d be optimal transport cost OT_c where c is a non-negative lower semicontinuous cost function. Then, the minimax problem in the LHS of (11) and (13) reduces to

$$\min_{G \in \mathcal{G}} \max_{D \in \mathcal{F}} \mathbb{E}[D(\mathbf{X})] - \mathbb{E}[D^c(G(\mathbf{Z}))]. \quad (17)$$

Proof. We defer the proof to the Appendix. \square

Therefore the minimax game between G and D in (17) can be viewed as minimizing the optimal transport cost between generative models and the distributions matching moments over \mathcal{F} with $P_{\mathbf{X}}$'s moments. The following example applies this result to the first-order Wasserstein distance and recovers the WGAN problem [4] with a constrained 1-Lipschitz discriminator.

Example 2. Let the optimal transport cost in (17) be the W_1 distance, and suppose \mathcal{F} is a convex subset of 1-Lipschitz functions. Then, the minimax problem (17) will reduce to

$$\min_{G \in \mathcal{G}} \max_{D \in \mathcal{F}} \mathbb{E}[D(\mathbf{X})] - \mathbb{E}[D(G(\mathbf{Z}))]. \quad (18)$$

Therefore, the moment-matching interpretation also holds for WGAN: for a convex set \mathcal{F} of 1-Lipschitz functions WGAN finds the generative model with minimum W_1 distance to the distributions penalized to share the same moments over \mathcal{F} with the data distribution. We discuss two more examples in the Appendix: 1) for the indicator cost $c_I(\mathbf{x}, \mathbf{x}') = \mathbb{I}(\mathbf{x} \neq \mathbf{x}')$ corresponding to the total variation distance we draw the connection to the energy-based GAN [8], 2) for the second-order cost $c_2(\mathbf{x}, \mathbf{x}') = \|\mathbf{x} - \mathbf{x}'\|^2$ we recover [9]'s quadratic GAN formulation under the LQG setting assumptions, i.e. linear generator, quadratic discriminator and Gaussian input data.

5 Duality framework applied to neural net discriminators

We applied the duality framework to analyze GAN problems with convex discriminator sets. However, a neural net set $\mathcal{F}_{nn} = \{f_{\mathbf{w}} : \mathbf{w} \in \mathcal{W}\}$, where $f_{\mathbf{w}}$ denotes a neural net function with fixed architecture and weights \mathbf{w} in feasible set \mathcal{W} , does not generally satisfy this convexity assumption. Note that a linear combination of several neural net functions in \mathcal{F}_{nn} may not remain in \mathcal{F}_{nn} .

Therefore, we apply the duality framework to \mathcal{F}_{nn} 's convex hull, which we denote by $\text{conv}(\mathcal{F}_{nn})$, containing any convex combination of neural net functions in \mathcal{F}_{nn} . However, a convex combination of infinitely-many neural nets from \mathcal{F}_{nn} is characterized by infinitely-many parameters, which makes optimizing the discriminator over $\text{conv}(\mathcal{F}_{nn})$ computationally intractable. In the following theorem, we show that although a function in $\text{conv}(\mathcal{F}_{nn})$ is a combination of infinitely-many neural nets, that function can be approximated by uniformly combining boundedly-many neural nets in \mathcal{F}_{nn} .

Theorem 4. Suppose any function $f_{\mathbf{w}} \in \mathcal{F}_{nn}$ is L -Lipschitz and bounded as $|f_{\mathbf{w}}(\mathbf{x})| \leq M$. Also, assume that the k -dimensional random input \mathbf{X} is norm-bounded as $\|\mathbf{X}\|_2 \leq R$. Then, any function in $\text{conv}(\mathcal{F}_{nn})$ can be uniformly approximated over the ball $\|\mathbf{x}\|_2 \leq R$ within ϵ -error by a uniform combination $\hat{f}(\mathbf{x}) = \frac{1}{m} \sum_{i=1}^m f_{\mathbf{w}_i}(\mathbf{x})$ of $m = \mathcal{O}\left(\frac{M^2 k \log(LR/\epsilon)}{\epsilon^2}\right)$ functions $(f_{\mathbf{w}_i})_{i=1}^m \in \mathcal{F}_{nn}$.

Proof. We defer the proof to the Appendix. \square

The above theorem suggests using a uniform combination of multiple discriminator nets to find a better approximation of the solution to the divergence minimization problem in Theorem 1 solved over $\text{conv}(\mathcal{F}_{nn})$. Note that this approach is different from MIX-GAN [10] proposed for achieving equilibrium in GAN minimax game. While our approach considers a uniform combination of multiple neural nets as the discriminator, MIX-GAN considers a randomized combination of the minimax game over multiple neural net discriminators and generators.

6 Minimum-sum hybrid of f-divergence and Wasserstein distance: GAN with Lipschitz or adversarially-trained discriminator

Here we apply the convex duality framework to a novel class of divergence measures. For each f-divergence d_f we define divergence d_{f,W_1} , which is the minimum sum hybrid of d_f and W_1 divergences, as follows

$$d_{f,W_1}(P_1, P_2) := \inf_Q W_1(P_1, Q) + d_f(Q, P_2). \quad (19)$$

The above infimum is taken over all distributions on random \mathbf{X} , searching for distribution Q minimizing the sum of the Wasserstein distance between P_1 and Q and the f-divergence from Q to P_2 . Note that the hybrid of JS-divergence and W_1 -distance defined earlier in (4) is a special case of the above definition. While f-divergence in f-GAN does not change continuously with the generator parameters, the following theorem shows that similar to the continuous behavior of W_1 -distance shown in [19, 4] the proposed hybrid divergence changes continuously with the generative model. We defer the proofs of this section's results to the Appendix.

Theorem 5. *Suppose $G_\theta \in \mathcal{G}$ is continuously changing with parameters θ . Then, for any Q and \mathbf{Z} , $d_{f,W_1}(P_{G_\theta(\mathbf{Z})}, Q)$ will behave continuously as a function of θ . Moreover, if G_θ is assumed to be locally Lipschitz, then $d_{f,W_1}(P_{G_\theta(\mathbf{Z})}, Q)$ will be differentiable w.r.t. θ almost everywhere.*

Our next result reveals the minimax problem dual to minimizing this hybrid divergence with symmetric f-divergence component. We note that this symmetricity condition is met by the JS-divergence and the squared Hellinger divergence among the f-divergence examples discussed in [3].

Theorem 6. *Consider d_{f,W_1} with a symmetric f-divergence d_f , i.e. $d_f(P, Q) = d_f(Q, P)$, satisfying the assumptions in Theorem 2. If the composition $f^* \circ D$ is 1-Lipschitz for all $D \in \mathcal{F}$, the minimax problem in Theorem 1 for the hybrid d_{f,W_1} reduces to the f-GAN problem, i.e.*

$$\min_{G \in \mathcal{G}} \max_{D \in \mathcal{F}} \mathbb{E}[D(\mathbf{X})] - \mathbb{E}[f^*(D(G(\mathbf{Z})))]. \quad (20)$$

The above theorem reveals that when the Lipschitz constant of discriminator D in f-GAN is properly regularized, then solving the f-GAN problem over the regularized discriminator also minimizes the continuous divergence measure d_{f,W_1} . As a special case, in the vanilla GAN problem (16) we only need to constrain discriminator \tilde{D} to be 1-Lipschitz, which can be done via the gradient penalty [13] or spectral normalization of \tilde{D} 's weight matrices [12], and then we minimize the continuously-behaving d_{JSD,W_1} . This result is also consistent with [12]'s empirical observations that regularizing the Lipschitz constant of the discriminator improves the training performance in vanilla GAN.

Our discussion has so far focused on the mixture of f-divergence and the first order Wasserstein distance, which suggests training f-GAN over Lipschitz-bounded discriminators. As a second solution, we prove that the desired continuity property can also be achieved through the following hybrid using the second-order Wasserstein (W_2) distance-squared:

$$d_{f,W_2}(P_1, P_2) := \inf_Q W_2^2(P_1, Q) + d_f(Q, P_2). \quad (21)$$

Theorem 7. *Suppose $G_\theta \in \mathcal{G}$ continuously changes with parameters $\theta \in \mathbb{R}^k$. Then, for any distribution Q and random vector \mathbf{Z} , $d_{f,W_2}(P_{G_\theta(\mathbf{Z})}, Q)$ will be continuous in θ . Also, if we further assume G_θ is bounded and locally-Lipschitz w.r.t. θ , then the hybrid divergence $d_{f,W_2}(P_{G_\theta(\mathbf{Z})}, Q)$ is almost everywhere differentiable w.r.t. θ .*

The following result shows that minimizing d_{f,W_2} reduces to f-GAN problem where the discriminator is being adversarially trained.

Theorem 8. Assume d_f and \mathcal{F} satisfy the assumptions in Theorem 6. Then, the minimax problem in Theorem 1 corresponding to the hybrid d_{f,W_2} divergence reduces to

$$\min_{G \in \mathcal{G}} \max_{D \in \mathcal{F}} \mathbb{E}[D(\mathbf{X})] + \mathbb{E} \left[\min_{\mathbf{u}} -f^*(D(G(\mathbf{Z}) + \mathbf{u})) + \|\mathbf{u}\|^2 \right]. \quad (22)$$

The above result reduces minimizing the hybrid d_{f,W_2} divergence to an f-GAN minimax game with a new third player. Here the third player assists the generator by perturbing the generated fake samples in order to make them harder to be distinguished from the real samples by the discriminator. The cost for perturbing a fake sample $G(\mathbf{Z})$ to $G(\mathbf{Z}) + \mathbf{u}$ will be $\|\mathbf{u}\|^2$, which constrains the power of the third player who can be interpreted as an adversary to the discriminator. To implement the game between these three players, we can adversarially learn the discriminator while we are training GAN, using the Wasserstein risk minimization (WRM) adversarial learning scheme discussed in [14].

7 Numerical Experiments

To evaluate our theoretical results, we used the CelebA [20] and LSUN-bedroom [21] datasets. Furthermore, in the Appendix we include the results of our experiments over the MNIST [22] dataset. We considered vanilla GAN [1] with the minimax formulation in (16) and DCGAN [23] convolutional architecture for discriminator and generator. We used the code provided by [13] and trained DCGAN via Adam optimizer [24] for 200,000 generator iterations. We applied 5 discriminator updates for each generator update.

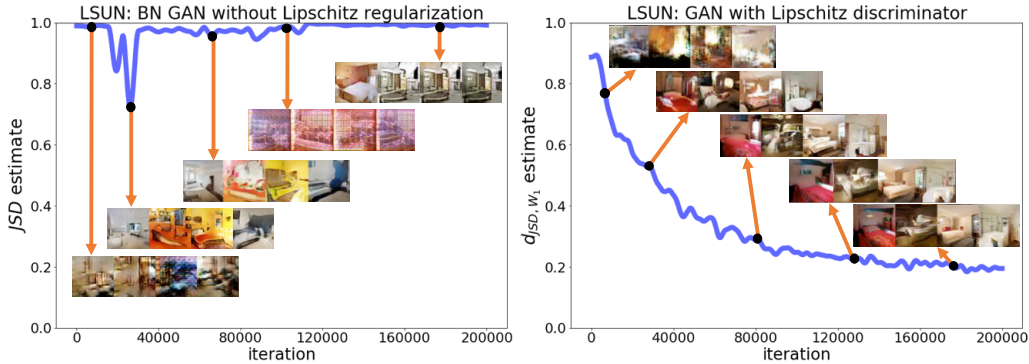


Figure 2: Divergence estimate in DCGAN trained over LSUN samples, (left) JS-divergence in standard DCGAN regularized with batch normalization, (right) hybrid d_{JS, W_1} in DCGAN with 1-Lipschitz discriminator regularized via spectral normalization.

Figure 2 shows how the discriminator loss evaluated over 2000 validation samples, which is an estimate of the divergence measure, changes as we train the DCGAN over LSUN samples. Using standard DCGAN regularized by only batch normalization (BN) [25], we observed (Figure 2-left) that the JS-divergence estimate always remains close to its maximum value 1 and also poorly correlates with the visual quality of generated samples. In this experiment, the GAN training failed and led to mode collapse starting at about the 110,000th iteration. On the other hand, after replacing BN with spectral normalization (SN) [12] to ensure the discriminator’s Lipschitzness, the discriminator loss decreased in a desired monotonic fashion (Figure 2-right). This observation is consistent with Theorems 5 and 6 showing that the discriminator loss becomes an estimate for the hybrid d_{JS, W_1} divergence changing continuously with the generator parameters. Also, the samples generated by the Lipschitz-regularized DCGAN looked qualitatively better and correlated well with the estimate of d_{JS, W_1} divergence.

Figure 3 shows the results of similar experiments over the CelebA dataset. Again, we observed (Figure 3-top left) that the JS-divergence estimate remains close to 1 while training DCGAN with BN. However, after applying two different Lipschitz regularization methods, SN and the gradient penalty (GP) [13] in Figures 3-top right and bottom left, we observed that the hybrid d_{JS, W_1} changed nicely and monotonically, and correlated properly with the sharpness of samples generated. Figure 3-bottom right shows that a similar desired behavior can also be achieved using the second-order hybrid d_{JS, W_2} divergence. In this case, we trained the DCGAN discriminator via the WRM adversarial learning scheme [14].

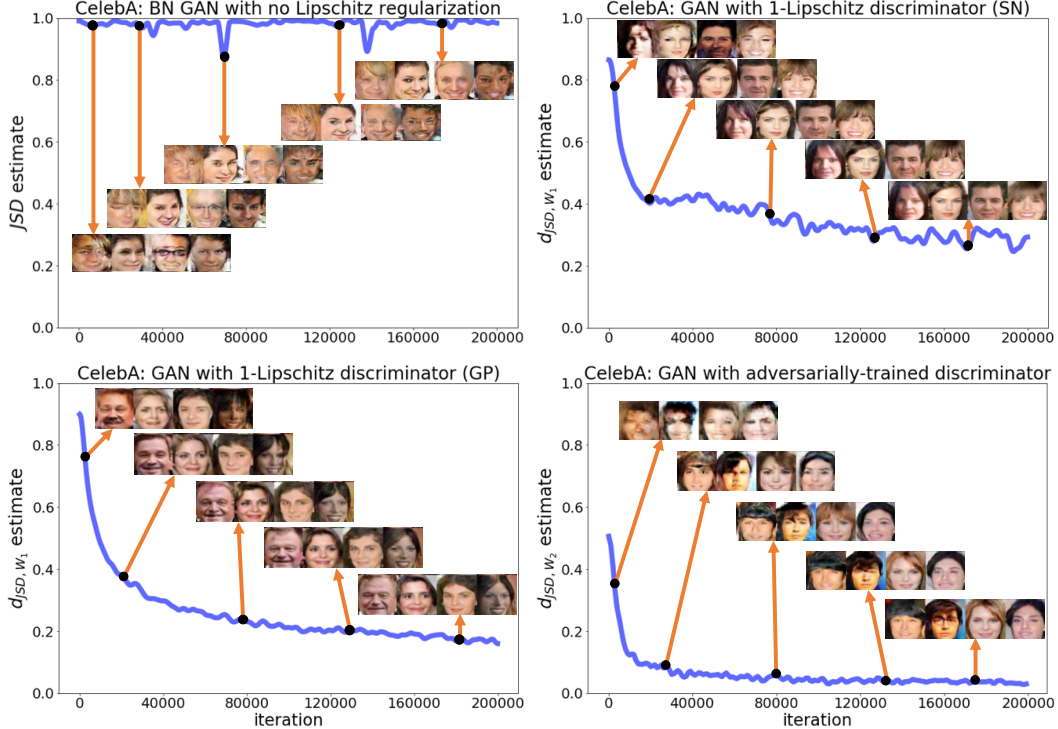


Figure 3: Divergence estimate in DCGAN trained over CelebA samples, (top-left) JS-divergence in DCGAN regularized with batch normalization, (top-right) hybrid d_{JSD, W_1} in DCGAN with spectrally-normalized discriminator, (bottom-left) hybrid d_{JSD, W_1} in DCGAN with 1-Lipschitz discriminator regularized via the gradient penalty, (bottom-right) hybrid d_{JSD, W_2} in DCGAN with discriminator being adversarially-trained using WRM.

8 Related Work

Theoretical studies of GAN have focused on three different aspects: approximation, generalization, and optimization. On the approximation properties of GAN, [11] studies GAN’s approximation power using a moment-matching approach. The authors view the maximized discriminator objective as an \mathcal{F} -based adversarial divergence, showing that the adversarial divergence between two distributions takes its minimum value if and only if the two distributions share the same moments over \mathcal{F} . Our convex duality framework interprets their result and further draws the connection to the original divergence measure. [26] studies the f-GAN problem through an information geometric approach based on the Bregman divergence and its connection to f-divergence.

Analyzing GAN’s generalization performance is another problem of interest in several recent works. [10] proves generalization guarantees for GANs in terms of \mathcal{F} -based distance measures. [27] uses an elegant approach based on the Birthday Paradox to empirically study the generalizability of GAN’s learned models. [28] develops a quantitative approach for examining diversity and generalization in GAN’s learned distribution. [29] studies approximation-generalization trade-offs in GAN by analyzing the discriminative power of \mathcal{F} -based distances. Regarding optimization properties of GAN, [30, 31] propose duality-based methods for improving the optimization performance in training deep generative models. [32] suggests applying noise convolution with input data for boosting the training performance in f-GAN. Moreover, several other works including [33, 34, 35, 9, 36] explore the optimization and stability properties of training GANs. Finally, we note that the same convex analysis approach used in this paper has further provided a powerful theoretical framework to analyze various supervised and unsupervised learning problems [37, 38, 39, 40, 41].

Acknowledgments: We are grateful for support under a Stanford Graduate Fellowship, the National Science Foundation grant under CCF-1563098, and the Center for Science of Information (CSOI), an NSF Science and Technology Center under grant agreement CCF-0939370.

References

- [1] Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. Generative adversarial nets. In *Advances in neural information processing systems*, pages 2672–2680, 2014.
- [2] Ian Goodfellow. Nips 2016 tutorial: Generative adversarial networks. *arXiv preprint arXiv:1701.00160*, 2016.
- [3] Sebastian Nowozin, Botond Cseke, and Ryota Tomioka. f-gan: Training generative neural samplers using variational divergence minimization. In *Advances in Neural Information Processing Systems*, pages 271–279, 2016.
- [4] Martin Arjovsky, Soumith Chintala, and Léon Bottou. Wasserstein generative adversarial networks. *International Conference on Machine Learning*, 2017.
- [5] Gintare Karolina Dziugaite, Daniel M Roy, and Zoubin Ghahramani. Training generative neural networks via maximum mean discrepancy optimization. *arXiv preprint arXiv:1505.03906*, 2015.
- [6] Yujia Li, Kevin Swersky, and Rich Zemel. Generative moment matching networks. In *International Conference on Machine Learning*, pages 1718–1727, 2015.
- [7] Chun-Liang Li, Wei-Cheng Chang, Yu Cheng, Yiming Yang, and Barnabás Póczos. Mmd gan: Towards deeper understanding of moment matching network. In *Advances in Neural Information Processing Systems*, pages 2200–2210, 2017.
- [8] Junbo Zhao, Michael Mathieu, and Yann LeCun. Energy-based generative adversarial network. *arXiv preprint arXiv:1609.03126*, 2016.
- [9] Soheil Feizi, Farzan Farnia, Tony Ginart, and David Tse. Understanding gans: the lqg setting. *arXiv preprint arXiv:1710.10793*, 2017.
- [10] Sanjeev Arora, Rong Ge, Yingyu Liang, Tengyu Ma, and Yi Zhang. Generalization and equilibrium in generative adversarial nets (gans). *arXiv preprint arXiv:1703.00573*, 2017.
- [11] Shuang Liu, Olivier Bousquet, and Kamalika Chaudhuri. Approximation and convergence properties of generative adversarial learning. In *Advances in Neural Information Processing Systems*, pages 5551–5559, 2017.
- [12] Takeru Miyato, Toshiki Kataoka, Masanori Koyama, and Yuichi Yoshida. Spectral normalization for generative adversarial networks. *International Conference on Learning Representations*, 2018.
- [13] Ishaan Gulrajani, Faruk Ahmed, Martin Arjovsky, Vincent Dumoulin, and Aaron C Courville. Improved training of wasserstein gans. In *Advances in Neural Information Processing Systems*, pages 5769–5779, 2017.
- [14] Aman Sinha, Hongseok Namkoong, and John Duchi. Certifiable distributional robustness with principled adversarial training. In *International Conference on Learning Representations*, 2018.
- [15] Imre Csiszár, Paul C Shields, et al. Information theory and statistics: A tutorial. *Foundations and Trends® in Communications and Information Theory*, 1(4):417–528, 2004.
- [16] Cédric Villani. *Optimal transport: old and new*, volume 338. Springer Science & Business Media, 2008.
- [17] Yasemin Altun and Alex Smola. Unifying divergence minimization and statistical inference via convex duality. In *International Conference on Computational Learning Theory*, pages 139–153, 2006.
- [18] Stephen Boyd and Lieven Vandenberghe. *Convex optimization*. Cambridge university press, 2004.
- [19] Martin Arjovsky and Léon Bottou. Towards principled methods for training generative adversarial networks. *arXiv preprint arXiv:1701.04862*, 2017.
- [20] Ziwei Liu, Ping Luo, Xiaogang Wang, and Xiaoou Tang. Deep learning face attributes in the wild. In *Proceedings of International Conference on Computer Vision (ICCV)*, 2015.
- [21] Fisher Yu, Ari Seff, Yinda Zhang, Shuran Song, Thomas Funkhouser, and Jianxiong Xiao. Lsun: Construction of a large-scale image dataset using deep learning with humans in the loop. *arXiv preprint arXiv:1506.03365*, 2015.
- [22] Yann LeCun. The mnist database of handwritten digits. <http://yann.lecun.com/exdb/mnist/>, 1998.
- [23] Alec Radford, Luke Metz, and Soumith Chintala. Unsupervised representation learning with deep convolutional generative adversarial networks. *arXiv preprint arXiv:1511.06434*, 2015.
- [24] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014.
- [25] Sergey Ioffe and Christian Szegedy. Batch normalization: Accelerating deep network training by reducing internal covariate shift. *arXiv preprint arXiv:1502.03167*, 2015.

- [26] Richard Nock, Zac Cranko, Aditya K Menon, Lizhen Qu, and Robert C Williamson. f-gans in an information geometric nutshell. In *Advances in Neural Information Processing Systems*, pages 456–464, 2017.
- [27] Sanjeev Arora and Yi Zhang. Do gans actually learn the distribution? an empirical study. *arXiv preprint arXiv:1706.08224*, 2017.
- [28] Shibani Santurkar, Ludwig Schmidt, and Aleksander Madry. A classification-based perspective on gan distributions. *arXiv preprint arXiv:1711.00970*, 2017.
- [29] Pengchuan Zhang, Qiang Liu, Dengyong Zhou, Tao Xu, and Xiaodong He. On the discrimination-generalization tradeoff in GANs. *International Conference on Learning Representations*, 2018.
- [30] Xu Chen, Jiang Wang, and Hao Ge. Training generative adversarial networks via primal-dual subgradient methods: a Lagrangian perspective on GAN. In *International Conference on Learning Representations*, 2018.
- [31] Shengjia Zhao, Jiaming Song, and Stefano Ermon. The information autoencoding family: A lagrangian perspective on latent variable generative models. *arXiv preprint arXiv:1806.06514*, 2018.
- [32] Kevin Roth, Aurelien Lucchi, Sebastian Nowozin, and Thomas Hofmann. Stabilizing training of generative adversarial networks through regularization. In *Advances in Neural Information Processing Systems*, pages 2015–2025, 2017.
- [33] Vaishnavh Nagarajan and J Zico Kolter. Gradient descent gan optimization is locally stable. In *Advances in Neural Information Processing Systems*, pages 5591–5600, 2017.
- [34] Lars Mescheder, Sebastian Nowozin, and Andreas Geiger. The numerics of gans. In *Advances in Neural Information Processing Systems*, pages 1823–1833, 2017.
- [35] Constantinos Daskalakis, Andrew Ilyas, Vasilis Syrgkanis, and Haoyang Zeng. Training gans with optimism. *arXiv preprint arXiv:1711.00141*, 2017.
- [36] Maziar Sanjabi, Jimmy Ba, Meisam Razaviyayn, and Jason D Lee. Solving approximate wasserstein gans to stationarity. *arXiv preprint arXiv:1802.08249*, 2018.
- [37] Miroslav Dudík, Steven J Phillips, and Robert E Schapire. Maximum entropy density estimation with generalized regularization and an application to species distribution modeling. *Journal of Machine Learning Research*, 8(Jun):1217–1260, 2007.
- [38] Meisam Razaviyayn, Farzan Farnia, and David Tse. Discrete rényi classifiers. In *Advances in Neural Information Processing Systems*, pages 3276–3284, 2015.
- [39] Farzan Farnia and David Tse. A minimax approach to supervised learning. In *Advances in Neural Information Processing Systems*, pages 4240–4248, 2016.
- [40] Rizal Fathony, Anqi Liu, Kaiser Asif, and Brian Ziebart. Adversarial multiclass classification: A risk minimization perspective. In *Advances in Neural Information Processing Systems*, pages 559–567, 2016.
- [41] Rizal Fathony, Mohammad Ali Bashiri, and Brian Ziebart. Adversarial surrogate losses for ordinal regression. In *Advances in Neural Information Processing Systems*, pages 563–573, 2017.
- [42] Yusuke Tsuzuku, Issei Sato, and Masashi Sugiyama. Lipschitz-margin training: Scalable certification of perturbation invariance for deep neural networks. *arXiv preprint arXiv:1802.04034*, 2018.
- [43] Jonathan M Borwein. A very complicated proof of the minimax theorem. *Minimax Theory and its Applications*, 1(1):21–27, 2016.
- [44] Patrick Billingsley. *Convergence of probability measures*. John Wiley & Sons, 2013.

9 Appendix

9.1 Additional numerical results

9.1.1 LSUN divergence estimates for different training schemes

Figure 4 shows the complete divergence estimates over LSUN dataset for the GAN training schemes described in the main text. While the hybrid divergence measures d_{JSD,W_1} , d_{JSD,W_2} decreased smoothly as the DCGAN was being trained, the JS-divergence always remained close to its maximum value 1 which led to lower-quality produced samples.

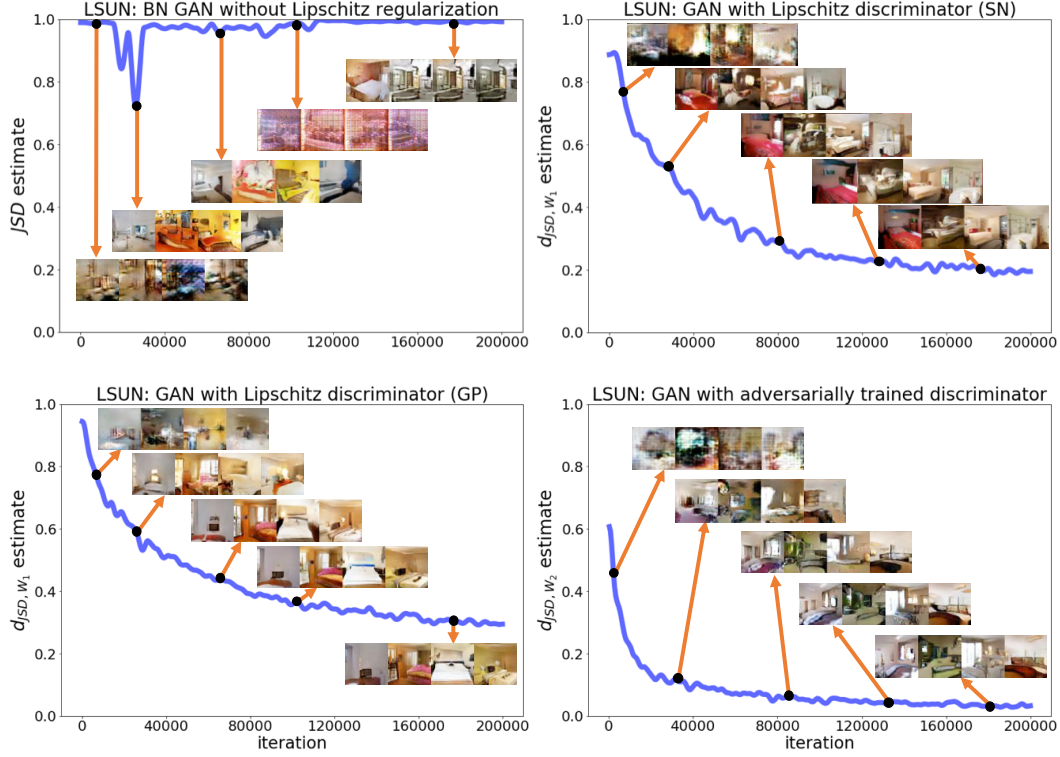


Figure 4: Divergence estimate in DCGAN trained over LSUN samples, (top-left) JS-divergence in DCGAN regularized with batch normalization, (top-right) hybrid d_{JSD, W_1} in DCGAN with spectrally-normalized discriminator, (bottom-left) hybrid d_{JSD, W_1} in DCGAN with 1-Lipschitz discriminator regularized via the gradient penalty, (bottom-right) hybrid d_{JSD, W_2} in DCGAN with discriminator being adversarially-trained using WRM.

9.1.2 CelebA, LSUN, MNIST images generated by different trainings of DCGAN

Figures 5, 6, and 7 show the CelebA, LSUN, and MNIST samples generated by vanilla DCGAN trained via the different methods described in the main text. Observe that applying Lipschitz regularization and adversarial training to the discriminator consistently result in the highest quality generator output samples. We note that tight SN in these figures refers to [42]’s spectral normalization method for convolutional layers, which precisely normalizes a conv layer’s spectral norm and hence guarantees the 1-Lipschitzness of the discriminator neural net. Note that for non-tight SN we use the original heuristic for normalizing convolutional layers’ operator norm introduced in [12].

9.2 Proof of Theorem 1

Theorem 1 and Corollary 1 directly result from the following two lemmas.

Lemma 1. *Suppose divergence $d(P, Q)$ is non-negative, lower semicontinuous and convex in distribution Q . Consider a convex subset of continuous functions \mathcal{F} and assume support set \mathcal{X} is compact. Then, the following duality holds for any pair of distributions P_1, P_2 :*

$$\max_{D \in \mathcal{F}} \mathbb{E}_{P_2}[D(\mathbf{X})] - d_{P_1}^*(D) = \min_Q \{ d(P_1, Q) + \max_{D \in \mathcal{F}} \{ \mathbb{E}_{P_2}[D(\mathbf{X})] - \mathbb{E}_Q[D(\mathbf{X})] \} \}. \quad (23)$$

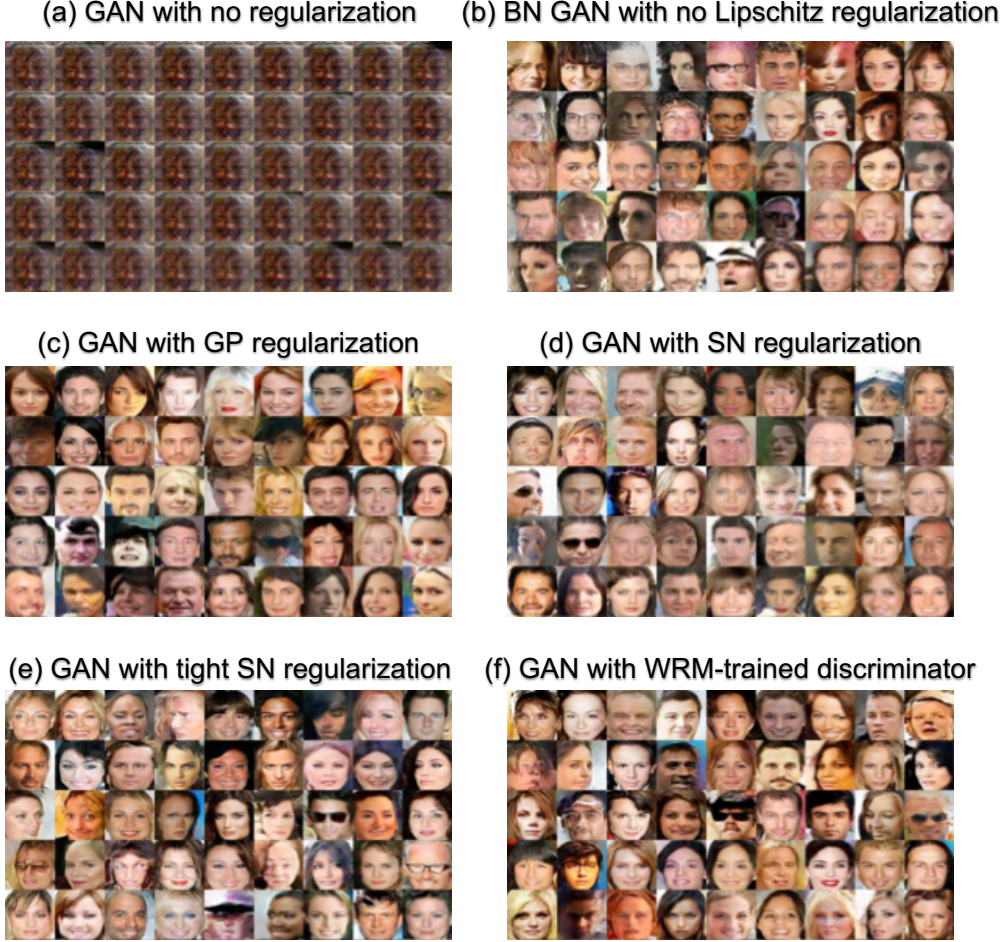


Figure 5: Samples generated by DCGAN trained over CelebA samples

Proof. Note that

$$\begin{aligned}
& \min_Q \{ d(P_1, Q) + \max_{D \in \mathcal{F}} \{ \mathbb{E}_{P_2}[D(\mathbf{X})] - \mathbb{E}_Q[D(\mathbf{X})] \} \} \\
&= \min_Q \max_{D \in \mathcal{F}} \{ d(P_1, Q) + \mathbb{E}_{P_2}[D(\mathbf{X})] - \mathbb{E}_Q[D(\mathbf{X})] \} \\
&\stackrel{(a)}{=} \max_{D \in \mathcal{F}} \min_Q \{ d(P_1, Q) + \mathbb{E}_{P_2}[D(\mathbf{X})] - \mathbb{E}_Q[D(\mathbf{X})] \} \\
&= \max_{D \in \mathcal{F}} \{ \mathbb{E}_{P_2}[D(\mathbf{X})] + \min_Q \{ d(P_1, Q) - \mathbb{E}_Q[D(\mathbf{X})] \} \} \\
&= \max_{D \in \mathcal{F}} \{ \mathbb{E}_{P_2}[D(\mathbf{X})] - \max_Q \{ \mathbb{E}_Q[D(\mathbf{X})] - d(P_1, Q) \} \} \\
&\stackrel{(b)}{=} \max_{D \in \mathcal{F}} \mathbb{E}_{P_2}[D(\mathbf{X})] - d_{P_1}^*(D).
\end{aligned} \tag{24}$$

Here (a) is a consequence of the generalized Sion's minimax theorem [43], because the space of probability measures on compact \mathcal{X} is convex and weakly compact [44], \mathcal{F} is assumed to be convex, the minimax objective is lower semicontinuous and convex in Q and linear in D . (b) holds according to the conjugate d_P^* 's definition. \square

Lemma 2. Assume divergence $d(P, Q)$ is non-negative, lower semicontinuous and convex in distribution Q over compact \mathcal{X} . Consider a linear space subset of continuous functions \mathcal{F} . Then, the following duality holds for any pair of distributions P_1, P_2 :

$$\min_{Q \in \mathcal{P}_{\mathcal{F}}(P_2)} d(P_1, Q) = \max_{D \in \mathcal{F}} \mathbb{E}_{P_2}[D(\mathbf{X})] - d_{P_1}^*(D). \tag{25}$$



Figure 6: Samples generated by DCGAN trained over LSUN-bedroom samples

Proof. This lemma is a consequence of Lemma 1. Note that a linear space \mathcal{F} is a convex set. Therefore, Lemma 1 applies to \mathcal{F} . However, since \mathcal{F} is a linear space i.e. for any $D \in \mathcal{F}$ and $\lambda \in \mathbb{R}$ it includes λD we have

$$\max_{D \in \mathcal{F}} \{ \mathbb{E}_{P_2}[D(\mathbf{X})] - \mathbb{E}_Q[D(\mathbf{X})] \} = \begin{cases} 0 & \text{if } Q \in \mathcal{P}_{\mathcal{F}}(P_2) \\ +\infty & \text{otherwise.} \end{cases} \quad (26)$$

As a result, the minimizing Q^* precisely matches the moments over \mathcal{F} to P_2 's moments, which completes the proof. \square

9.3 Proof of Theorem 2

We first prove the following lemma.

Lemma 3. Consider f -divergence d_f corresponding to function f which has a non-decreasing convex-conjugate f^* . Then, for any continuous D

$$d_{fP}^*(D) = \mathbb{E}_P[f^*(D(\mathbf{X}) + \lambda_0)] - \lambda_0 \quad (27)$$

where $\lambda_0 \in \mathbb{R}$ satisfies $\mathbb{E}_P[f^{*'}(D(\mathbf{X}) + \lambda_0)] = 1$. Here $f^{*'}$ stands for the derivative of conjugate function f^* which is supposed to be non-negative everywhere.

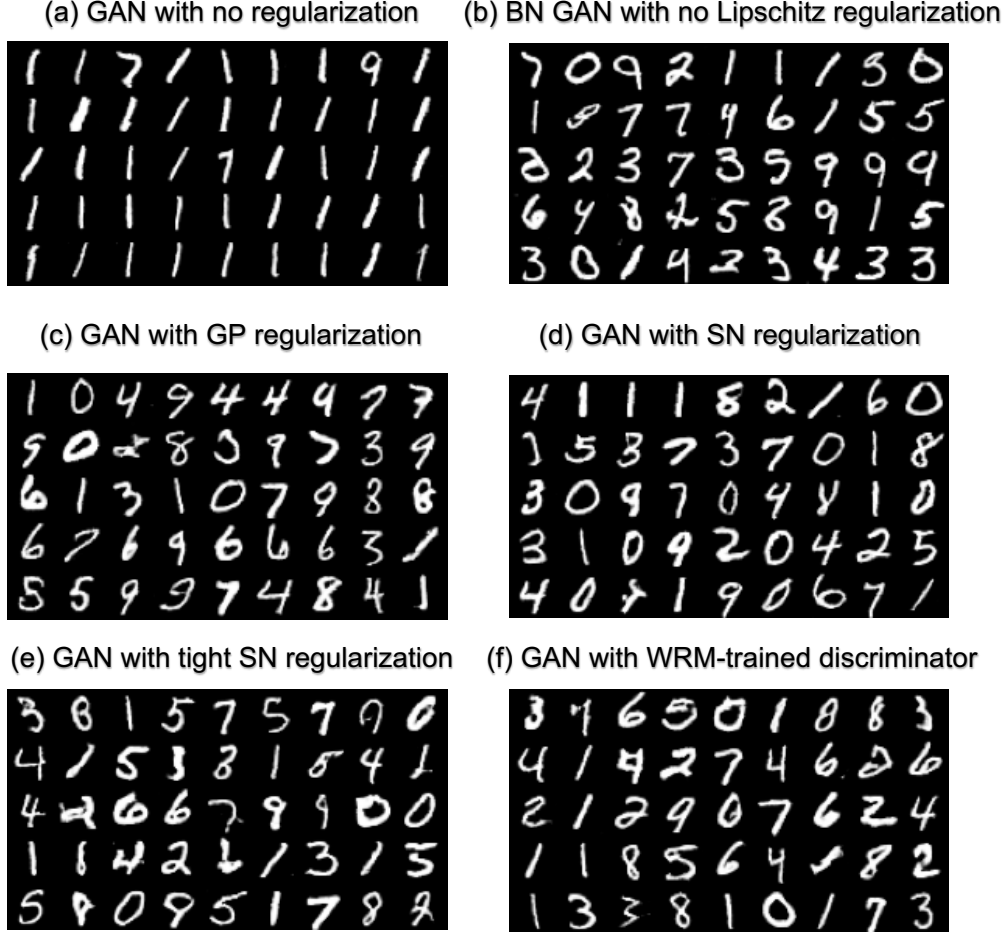


Figure 7: Samples generated by DCGAN trained over MNIST samples

Proof. Note that

$$\begin{aligned}
 d_{fP}^*(D) &\stackrel{(a)}{=} \sup_Q \mathbb{E}_Q[D(\mathbf{X})] - d_f(P, Q) \\
 &\stackrel{(b)}{=} \sup_Q \mathbb{E}_Q[D(\mathbf{X})] - \mathbb{E}_P\left[f\left(\frac{q(\mathbf{X})}{p(\mathbf{X})}\right)\right] \\
 &\stackrel{(c)}{=} \max_{q(\mathbf{x}) \geq 0, \int q(\mathbf{x}) d\mathbf{x} = 1} \int q(\mathbf{x}) D(\mathbf{x}) d\mathbf{x} - \mathbb{E}_P\left[f\left(\frac{q(\mathbf{X})}{p(\mathbf{X})}\right)\right] \\
 &\stackrel{(d)}{=} \min_{\lambda \in \mathbb{R}} -\lambda + \max_{q(\mathbf{x}) \geq 0} \int q(\mathbf{x}) (D(\mathbf{x}) + \lambda) d\mathbf{x} - \mathbb{E}_P\left[f\left(\frac{q(\mathbf{X})}{p(\mathbf{X})}\right)\right] \\
 &\stackrel{(e)}{=} \min_{\lambda \in \mathbb{R}} -\lambda + \max_{r(\mathbf{x}) \geq 0} \mathbb{E}_P[r(\mathbf{X})(D(\mathbf{X}) + \lambda) - f\left(\frac{r(\mathbf{X})}{p(\mathbf{X})}\right)] \\
 &\stackrel{(f)}{=} \min_{\lambda \in \mathbb{R}} -\lambda + \mathbb{E}_P\left[\max_{r(\mathbf{X}) \geq 0} r(\mathbf{X})(D(\mathbf{X}) + \lambda) - f\left(\frac{r(\mathbf{X})}{p(\mathbf{X})}\right)\right] \\
 &\stackrel{(g)}{=} \min_{\lambda \in \mathbb{R}} -\lambda + \mathbb{E}_P[f^*(D(\mathbf{X}) + \lambda)] \\
 &= -\max_{\lambda \in \mathbb{R}} \lambda - \mathbb{E}_P[f^*(D(\mathbf{X}) + \lambda)] \tag{28}
 \end{aligned}$$

$$\stackrel{(h)}{=} -\lambda_0 + \mathbb{E}_P[f^*(D(\mathbf{X}) + \lambda_0)]. \tag{29}$$

Here (a) and (b) follow from the conjugate d_P^* and f-divergence d_f definitions. (c) rewrites the optimization problem in terms of the density function q corresponding to distribution Q . (d) uses the strong convex duality to move the density constraint $\int q(\mathbf{x}) d\mathbf{x} = 1$ to the objective. Note that strong duality holds, since we have a convex optimization problem with affine constraints. (e) rewrites the problem after a change of variable $r(\mathbf{x}) = q(\mathbf{x})/p(\mathbf{x})$. (f) holds since f and D are assumed to be continuous. (g) follows from the assumption that the derivative of f^* takes non-negative values, and hence the minimizing $r(\mathbf{x}) \geq 0$ also minimizes the unconstrained optimization for the convex conjugate f^*

$$f^*(D(\mathbf{X}) + \lambda) := \max_{r(\mathbf{X})} r(\mathbf{X})(D(\mathbf{X}) + \lambda) - f(r(\mathbf{X})).$$

Taking the derivative of the concave objective, the λ value maximizing the objective solves the equation $\mathbb{E}_P[f^{*'}(D(\mathbf{X}) + \lambda)] = 1$ which is assumed to be λ_0 . Therefore, (h) holds and the proof is complete. \square

Now we prove Theorem 2 which can be broken into two parts as follows.

Theorem (Theorem 2). *Consider f-divergence d_f where f has a non-decreasing conjugate f^* .*

(a) *Suppose \mathcal{F} is a convex set closed to a constant addition, i.e. for any $D \in \mathcal{F}$, $\lambda \in \mathbb{R}$ we have $D + \lambda \in \mathcal{F}$. Then,*

$$\begin{aligned} & \min_{P_{G(\mathbf{Z})} \in \mathcal{P}_G} \min_{Q_{\mathbf{X}}} d_f(P_{G(\mathbf{Z})}, Q) + \max_{D \in \mathcal{F}} \{\mathbb{E}_{P_{\mathbf{X}}}[D(\mathbf{X})] - \mathbb{E}_Q[D(\mathbf{X})]\} \\ &= \min_{G \in \mathcal{G}} \max_{D \in \mathcal{F}} \mathbb{E}_{P_{\mathbf{X}}}[D(\mathbf{X})] - \mathbb{E}[f^*(D(G(\mathbf{Z})))]. \end{aligned} \quad (30)$$

(b) *Suppose \mathcal{F} is a linear space including the constant function $D_0(\mathbf{x}) = 1$. Then,*

$$\min_{P_{G(\mathbf{Z})} \in \mathcal{P}_G} \min_{Q_{\mathbf{X}} \in \mathcal{P}_{\mathcal{F}}(P_{\mathbf{X}})} d_f(P_{G(\mathbf{Z})}, Q) = \min_{G \in \mathcal{G}} \max_{D \in \mathcal{F}} \mathbb{E}_{P_{\mathbf{X}}}[D(\mathbf{X})] - \mathbb{E}[f^*(D(G(\mathbf{Z})))]. \quad (31)$$

Proof. This theorem is an application of Theorem 1 and Corollary 1. For part (a) we have

$$\begin{aligned} & \min_{P_{G(\mathbf{Z})} \in \mathcal{P}_G} \min_{Q_{\mathbf{X}}} d_f(P_{G(\mathbf{Z})}, Q) + \max_{D \in \mathcal{F}} \{\mathbb{E}_{P_{\mathbf{X}}}[D(\mathbf{X})] - \mathbb{E}_Q[D(\mathbf{X})]\} \\ & \stackrel{(c)}{=} \min_{G \in \mathcal{G}} \max_{D \in \mathcal{F}} \mathbb{E}_{P_{\mathbf{X}}}[D(\mathbf{X})] - d_{f_{P_{G(\mathbf{Z})}}}^*(D) \\ & \stackrel{(d)}{=} \min_{G \in \mathcal{G}} \max_{D \in \mathcal{F}} \mathbb{E}_{P_{\mathbf{X}}}[D(\mathbf{X})] + \max_{\lambda \in \mathbb{R}} \lambda - \mathbb{E}[f^*(D(G(\mathbf{Z})) + \lambda)] \\ &= \min_{G \in \mathcal{G}} \max_{D \in \mathcal{F}, \lambda \in \mathbb{R}} \mathbb{E}_{P_{\mathbf{X}}}[D(\mathbf{X}) + \lambda] - \mathbb{E}[f^*(D(G(\mathbf{Z})) + \lambda)] \\ & \stackrel{(e)}{=} \min_{G \in \mathcal{G}} \max_{D \in \mathcal{F}} \mathbb{E}_{P_{\mathbf{X}}}[D(\mathbf{X})] - \mathbb{E}[f^*(D(G(\mathbf{Z})))]. \end{aligned}$$

Here (c) is a direct result of Theorem 1. (d) uses the simplified version (28) for $d_{f_P}^*$. (e) follows from the assumption that \mathcal{F} is closed to constant additions.

For part (b) note that since \mathcal{F} is a linear space and includes $D_0(\mathbf{x}) = 1$, it is closed to constant additions. Hence, an application of Corollary 1 reveals

$$\begin{aligned} \min_{P_{G(\mathbf{Z})} \in \mathcal{P}_G} \min_{Q_{\mathbf{X}} \in \mathcal{P}_{\mathcal{F}}(P_{\mathbf{X}})} d_f(P_{G(\mathbf{Z})}, Q) &= \min_{G \in \mathcal{G}} \max_{D \in \mathcal{F}} \mathbb{E}_{P_{\mathbf{X}}}[D(\mathbf{X})] - d_{f_{P_{G(\mathbf{Z})}}}^*(D) \\ &= \min_{G \in \mathcal{G}} \max_{D \in \mathcal{F}} \mathbb{E}_{P_{\mathbf{X}}}[D(\mathbf{X})] + \max_{\lambda \in \mathbb{R}} \lambda - \mathbb{E}[f^*(D(G(\mathbf{Z})) + \lambda)] \\ &= \min_{G \in \mathcal{G}} \max_{D \in \mathcal{F}, \lambda \in \mathbb{R}} \mathbb{E}_{P_{\mathbf{X}}}[D(\mathbf{X}) + \lambda] - \mathbb{E}[f^*(D(G(\mathbf{Z})) + \lambda)] \\ &= \min_{G \in \mathcal{G}} \max_{D \in \mathcal{F}} \mathbb{E}_{P_{\mathbf{X}}}[D(\mathbf{X})] - \mathbb{E}[f^*(D(G(\mathbf{Z})))], \end{aligned}$$

which makes the proof complete. \square

9.4 Proof of Theorem 3

Theorem 3 is a direct application of the following lemma to Theorem 1 and Corollary 1.

Lemma 4. Let c be a lower semicontinuous non-negative cost function. Considering the c -transform operation D^c defined in the text, the following holds for any continuous D

$$OT_{cP}^*(D) = \mathbb{E}_P[D^c(\mathbf{X})]. \quad (32)$$

Proof. We have

$$\begin{aligned} OT_{cP}^*(D) &\stackrel{(a)}{=} \sup_Q \mathbb{E}_Q[D(\mathbf{X}')] - OT_c(P, Q) \\ &\stackrel{(b)}{=} - \inf_Q \inf_{M \in \Pi(P, Q)} \mathbb{E}_M [c(\mathbf{X}, \mathbf{X}') - D(\mathbf{X}')] \\ &= - \inf_{Q, M \in \Pi(P, Q)} \mathbb{E}_M [c(\mathbf{X}, \mathbf{X}') - D(\mathbf{X}')] \\ &\stackrel{(c)}{\geq} - \mathbb{E}_P [\inf_{\mathbf{x}'} c(\mathbf{X}, \mathbf{x}') - D(\mathbf{x}')] \\ &= \mathbb{E}_P [\sup_{\mathbf{x}'} D(\mathbf{x}') - c(\mathbf{X}, \mathbf{x}')] \\ &\stackrel{(d)}{=} \mathbb{E}_P [D^c(\mathbf{X})]. \end{aligned}$$

Here (a), (b), (d) hold according to the definitions. Moreover, we show (c) will hold with equality under the lemma's assumptions. $c(\mathbf{x}, \mathbf{x}') - D(\mathbf{x}')$ is lower semicontinuous, and hence for every $\epsilon > 0$ there exists a measurable function $v(\mathbf{x})$ such that for the coupling $M = \pi_{\mathbf{X}, v(\mathbf{X})}$ the absolute difference $|\mathbb{E}_M [c(\mathbf{X}, \mathbf{X}') - D(\mathbf{X}')] - \mathbb{E}_P [\inf_{\mathbf{x}'} c(\mathbf{X}, \mathbf{x}') - D(\mathbf{x}')]| < \epsilon$ is ϵ -bounded. Therefore, (c) holds with equality and the proof is complete. \square

9.5 Proof of Theorem 4

Consider a convex combination of functions from \mathcal{F}_{nm} as $f_\alpha(\mathbf{x}) = \int \alpha(\mathbf{w}) f_{\mathbf{w}}(\mathbf{x}) d\mathbf{w}$ where α can be considered as a probability density function over feasible set \mathcal{W} . Consider m samples $(\mathbf{W}_i)_{i=1}^m$ taken i.i.d. from α . Since any $f_{\mathbf{w}}$ is M -bounded, according to Hoeffding's inequality for a fixed \mathbf{x} we have

$$\Pr \left(\left| \frac{1}{m} \sum_{i=1}^m f_{\mathbf{W}_i}(\mathbf{x}) - \mathbb{E}_{\mathbf{W} \sim \alpha} [f_{\mathbf{W}}(\mathbf{x})] \right| \geq \frac{\epsilon}{2} \right) \leq 2 \exp\left(-\frac{m\epsilon^2}{8M^2}\right). \quad (33)$$

Next we consider a δ -covering for the ball $\{\mathbf{x} : \|\mathbf{x}\|_2 \leq R\}$, where we choose $\delta = \frac{\epsilon}{4L}$. We know a δ -covering $\{\mathbf{x}_j : 1 \leq j \leq N\}$ exists with a bounded size $N \leq (12LR/\epsilon)^k$ [15]. Then, an application of the union bound implies

$$\begin{aligned} \Pr \left(\max_{1 \leq j \leq N} \left| \frac{1}{m} \sum_{i=1}^m f_{\mathbf{W}_i}(\mathbf{x}_j) - \mathbb{E}_{\mathbf{W} \sim \alpha} [f_{\mathbf{W}}(\mathbf{x}_j)] \right| \geq \frac{\epsilon}{2} \right) &\leq 2N \exp\left(-\frac{m\epsilon^2}{8M^2}\right) \\ &\leq \exp\left(-\frac{m\epsilon^2}{8M^2} + k \log\left(\frac{12LR}{\epsilon}\right) + \log 2\right) \end{aligned}$$

Hence if we have $-\frac{m\epsilon^2}{8M^2} + k \log\left(\frac{12LR}{\epsilon}\right) + \log 2 < 0$ the above upper-bound is strictly less than 1, showing there exists at least one outcome $(\mathbf{w}_i)_{i=1}^m$ satisfying

$$\max_{1 \leq j \leq N} \left| \frac{1}{m} \sum_{i=1}^m f_{\mathbf{w}_i}(\mathbf{x}_j) - \mathbb{E}_{\mathbf{W} \sim \alpha} [f_{\mathbf{W}}(\mathbf{x}_j)] \right| < \frac{\epsilon}{2}. \quad (34)$$

Then, we claim the following holds over the norm-bounded $\{\mathbf{x} : \|\mathbf{x}\|_2 \leq R\}$:

$$\sup_{\|\mathbf{x}\|_2 \leq R} \left| \frac{1}{m} \sum_{i=1}^m f_{\mathbf{w}_i}(\mathbf{x}_j) - \mathbb{E}_{\mathbf{W} \sim \alpha} [f_{\mathbf{W}}(\mathbf{x}_j)] \right| < \epsilon. \quad (35)$$

This is because due to the definition of a δ -covering for any $\|\mathbf{x}\|_2 \leq R$ there exists \mathbf{x}_j for which $\|\mathbf{x}_j - \mathbf{x}\| \leq \frac{\epsilon}{4L}$. Then, since any $f_{\mathbf{w}}$ is supposed to be L -Lipschitz we have

$$\left| \frac{1}{m} \sum_{i=1}^m f_{\mathbf{w}_i}(\mathbf{x}_j) - \frac{1}{m} \sum_{i=1}^m f_{\mathbf{w}_i}(\mathbf{x}) \right| \leq \frac{\epsilon}{4}, \quad \left| \mathbb{E}_{\mathbf{W} \sim \alpha} [f_{\mathbf{W}}(\mathbf{x}_j)] - \mathbb{E}_{\mathbf{W} \sim \alpha} [f_{\mathbf{W}}(\mathbf{x})] \right| \leq \frac{\epsilon}{4} \quad (36)$$

which together with (34) shows (35). Hence, if we choose

$$m = \frac{8M^2}{\epsilon^2} (k \log(12LR/\epsilon) + \log 2) = \mathcal{O}\left(\frac{M^2 k \log(LR/\epsilon)}{\epsilon^2}\right) \quad (37)$$

there will be some weight assignments $(\mathbf{w}_i)_{i=1}^m$ such that their uniform combination $\frac{1}{m} \sum_{i=1}^m f_{\mathbf{w}_i}(\mathbf{x})$ ϵ -approximates the convex combination f_α uniformly over $\{\mathbf{x} : \|\mathbf{x}\|_2 \leq R\}$.

9.6 Proof of Theorem 5

We show that for any distributions P_0, P_1, P_2 the following holds

$$|d_{f, W_1}(P_0, P_2) - d_{f, W_1}(P_1, P_2)| \leq W_1(P_0, P_1). \quad (38)$$

The above inequality holds since if Q_0 and Q_1 solve the minimum sum optimization problems for $d_{f, W_1}(P_0, P_2), d_{f, W_1}(P_1, P_2)$, we have

$$\begin{aligned} d_{f, W_1}(P_0, P_2) - d_{f, W_1}(P_1, P_2) &\leq W_1(P_0, Q_1) - W_1(P_1, Q_1) \leq W_1(P_0, P_1), \\ d_{f, W_1}(P_1, P_2) - d_{f, W_1}(P_0, P_2) &\leq W_1(P_1, Q_0) - W_1(P_0, Q_0) \leq W_1(P_0, P_1) \end{aligned}$$

where the second inequalities in both these lines follow from the symmetricity and triangle inequality property of the W_1 -distance. Therefore, the following holds for any Q :

$$|d_{f, W_1}(P_{G_\theta(\mathbf{Z})}, Q) - d_{f, W_1}(P_{G_{\theta'}(\mathbf{Z})}, Q)| \leq W_1(P_{G_\theta(\mathbf{Z})}, P_{G_{\theta'}(\mathbf{Z})}).$$

Hence, we only need to show $W_1(P_{G_\theta(\mathbf{Z})}, Q)$ is changing continuously with θ and is almost everywhere differentiable. We prove these things using a similar proof to [4]'s proof for the continuity of the first-order Wasserstein distance.

Consider two functions $G_\theta, G_{\theta'}$. The joint distribution M for $(G_\theta(\mathbf{Z}), G_{\theta'}(\mathbf{Z}))$ is contained in $\Pi(P_{G_\theta(\mathbf{Z})}, P_{G_{\theta'}(\mathbf{Z})})$, which results in

$$\begin{aligned} W_1(P_{G_\theta(\mathbf{Z})}, P_{G_{\theta'}(\mathbf{Z})}) &\leq \mathbb{E}_M[\|\mathbf{X} - \mathbf{X}'\|] \\ &= \mathbb{E}[\|G_\theta(\mathbf{Z}) - G_{\theta'}(\mathbf{Z})\|]. \end{aligned} \quad (39)$$

If we let $\theta' \rightarrow \theta$ then $G_{\theta'}(\mathbf{z}) \rightarrow G_\theta(\mathbf{z})$ and hence $\|G_{\theta'}(\mathbf{z}) - G_\theta(\mathbf{z})\| \rightarrow 0$ hold pointwise. Since \mathcal{X} is assumed to be compact, there exists some finite R for which $0 \leq \|\mathbf{x} - \mathbf{x}'\| \leq R$ holds over the compact $\mathcal{X} \times \mathcal{X}$. Then the bounded convergence theorem implies $\mathbb{E}[\|G_\theta(\mathbf{Z}) - G_{\theta'}(\mathbf{Z})\|]$ converges to 0 as $\theta' \rightarrow \theta$. Then, since W_1 -distance always takes non-negative values

$$W_1(P_{G_\theta(\mathbf{Z})}, P_{G_{\theta'}(\mathbf{Z})}) \xrightarrow{\theta' \rightarrow \theta} 0.$$

Thus, W_1 satisfies the discussed continuity property and as a result $d_{f, W_1}(P_{G_\theta(\mathbf{Z})}, Q)$ changes continuously with θ . Furthermore, if G_θ is locally-Lipschitz and its Lipschitz constant w.r.t. parameters θ is bounded above by L ,

$$\begin{aligned} d_{f, W_1}(P_{G_\theta(\mathbf{Z})}, P_{G_{\theta'}(\mathbf{Z})}) &\leq W_1(P_{G_\theta(\mathbf{Z})}, P_{G_{\theta'}(\mathbf{Z})}) \\ &\leq \mathbb{E}[\|G_\theta(\mathbf{Z}) - G_{\theta'}(\mathbf{Z})\|] \\ &\leq L\|\theta - \theta'\|, \end{aligned} \quad (40)$$

which implies both $W_1(P_{G_\theta(\mathbf{Z})}, Q)$ and $d_{f, W_1}(P_{G_\theta(\mathbf{Z})}, Q)$ are everywhere continuous and almost everywhere differentiable w.r.t. θ .

9.7 Proof of Theorem 6

We first generalize the definition of the hybrid divergence to a general minimum-sum hybrid of an f-divergence and an optimal transport cost. For f-divergence d_f and optimal transport cost OT_c corresponding to convex function f and cost c respectively, we define the following hybrid $d_{f,c}$ of the two divergence measures:

$$d_{f,c}(P_1, P_2) := \inf_Q OT_c(P_1, Q) + d_f(Q, P_2). \quad (41)$$

Lemma 5. Given a symmetric f -divergence d_f with convex lower semicontinuous f and a non-negative lower semicontinuous c , $d_{f,c}(P_1, P_2)$ will be a convex function of P_1 and P_2 , and further satisfies the following generalization of the Kantorovich duality [16]:

$$d_{f,c}(P_1, P_2) = \sup_{D \text{ c-concave}} \mathbb{E}_{P_1}[D(\mathbf{X})] - \mathbb{E}_{P_2}[f^*(D^c(\mathbf{X}))]. \quad (42)$$

Proof. According to the Kantorovich duality [16] we have

$$\begin{aligned} d_{f,c}(P_1, P_2) &\stackrel{(a)}{=} \inf_Q OT_c(P_1, Q) + d_f(Q, P_2) \\ &\stackrel{(b)}{=} \inf_Q \sup_{D \text{ c-concave}} \mathbb{E}_{P_1}[D(\mathbf{X})] - \mathbb{E}_Q[D^c(\mathbf{X})] + d_f(Q, P_2) \\ &\stackrel{(c)}{=} \inf_Q \sup_{D \text{ c-concave}} \mathbb{E}_{P_1}[D(\mathbf{X})] - \mathbb{E}_Q[D^c(\mathbf{X})] + d_f(P_2, Q) \\ &\stackrel{(d)}{=} \sup_{D \text{ c-concave}} \inf_Q \mathbb{E}_{P_1}[D(\mathbf{X})] - \mathbb{E}_Q[D^c(\mathbf{X})] + d_f(P_2, Q) \\ &= \sup_{D \text{ c-concave}} \mathbb{E}_{P_1}[D(\mathbf{X})] + \inf_Q d_f(P_2, Q) - \mathbb{E}_Q[D^c(\mathbf{X})] \\ &\stackrel{(e)}{=} \sup_{D \text{ c-concave}} \mathbb{E}_{P_1}[D(\mathbf{X})] - d_{f_{P_2}}^*(D^c) \\ &\stackrel{(f)}{=} \sup_{D \text{ c-concave}} \mathbb{E}_{P_1}[D(\mathbf{X})] + \max_{\lambda \in \mathbb{R}} \lambda - \mathbb{E}_{P_2}[f^*(D^c(\mathbf{X}) + \lambda)] \\ &= \sup_{D \text{ c-concave}, \lambda \in \mathbb{R}} \mathbb{E}_{P_1}[D(\mathbf{X}) + \lambda] - \mathbb{E}_{P_2}[f^*(D^c(\mathbf{X}) + \lambda)]. \\ &= \sup_{D \text{ c-concave}} \mathbb{E}_{P_1}[D(\mathbf{X})] - \mathbb{E}_{P_2}[f^*(D^c(\mathbf{X}))]. \end{aligned}$$

Here (a) holds according to the definition. (b) is a consequence of the Kantorovich duality ([16], Theorem 5.10). (c) holds because d_f is assumed to be symmetric. (d) holds due to the generalized minimax theorem [43], since the space of distributions over compact \mathcal{X} is convex and weakly compact, the set of c -concave functions is convex, the minimax objective is concave in D and convex in Q . (e) holds according to the conjugate d_P^* 's definition, and (f) is based on our earlier result in (28). Note that the final expression is maximizing an objective linear in P_2 , which is convex in P_2 . The last equality holds since for any constant $\lambda \in \mathbb{R}$ if D^c is the c -transform of D , $D^c + \lambda$ will be the c -transform of $D + \lambda$. Finally, note that $d_{f,c}(P_1, P_2)$ is the supremum of some linear functions of P_1 and P_2 with compact support sets. Hence $d_{f,c}$ will be a convex function of P_1 and P_2 . \square

Now we prove the following generalization of Theorem 6, which directly results in Theorem 6 for the difference norm cost $c_1(\mathbf{x}, \mathbf{x}') = \|\mathbf{x} - \mathbf{x}'\|$. Here note that for cost c_1 the c -transform of a 1-Lipschitz function D will be D itself, which implies if $f^* \circ D$ is 1-Lipschitz then

$$-f^*(D(G(\mathbf{Z}))) = \inf_{\mathbf{x}'} -f^*(D(\mathbf{x}')) + c_1(G(\mathbf{Z}), \mathbf{x}').$$

Theorem (Generalization of Theorem 6). Assume d_f is a symmetric f -divergence, i.e. $d_f(P, Q) = d_f(Q, P)$, satisfying the assumptions in Lemma 2. Suppose \mathcal{F} is a convex set of continuous functions closed to constant additions and cost function c is non-negative and continuous. Then, the minimax problem in Theorem 1 and Corollary 1 for the mixed divergence $d_{f,c}$ reduces to

$$\min_{G \in \mathcal{G}} \max_{D \in \mathcal{F}} \mathbb{E}_{P_{\mathbf{X}}}[D(\mathbf{X})] + \mathbb{E}[\inf_{\mathbf{x}'} -f^*(D(\mathbf{x}')) + c(G(\mathbf{Z}), \mathbf{x}')]. \quad (43)$$

Proof. According to Lemma 5, $d_{f,c}(P, Q)$ satisfies the convexity property in Q . Hence, the assumptions of Theorem 1 and Corollary 1 hold and we only need to plug in the conjugate $d_{f,c_{P_1}}^*$ into

Corollary 1. According to the definition,

$$\begin{aligned}
d_{f,c}^*(D) &= \sup_{P_2} \mathbb{E}_{P_2}[D(\mathbf{X})] - d_{f,c}(P_1, P_2) \\
&= \sup_{P_2} \sup_Q -OT_c(P_1, Q) - d_f(Q, P_2) + \mathbb{E}_{P_2}[D(\mathbf{X})] \\
&= \sup_Q \sup_{P_2} -OT_c(P_1, Q) - d_f(Q, P_2) + \mathbb{E}_{P_2}[D(\mathbf{X})] \\
&= \sup_Q -OT_c(P_1, Q) + \sup_{P_2} \mathbb{E}_{P_2}[D(\mathbf{X})] - d_f(Q, P_2) \\
&= \sup_Q -OT_c(P_1, Q) + d_{f,Q}^*(D) \\
&\stackrel{(g)}{=} \sup_Q -OT_c(P_1, Q) + \min_{\lambda \in \mathbb{R}} -\lambda + \mathbb{E}_Q[f^*(D(\mathbf{X}) + \lambda)] \\
&= \sup_Q \min_{\lambda \in \mathbb{R}} -OT_c(P_1, Q) - \lambda + \mathbb{E}_Q[f^*(D(\mathbf{X}) + \lambda)] \\
&\stackrel{(h)}{=} \min_{\lambda \in \mathbb{R}} \sup_Q -OT_c(P_1, Q) - \lambda + \mathbb{E}_Q[f^*(D(\mathbf{X}) + \lambda)] \\
&\stackrel{(i)}{=} \inf_{\lambda \in \mathbb{R}} -\lambda + \mathbb{E}_{P_1}[(f^* \circ (D + \lambda))^c(\mathbf{X})].
\end{aligned}$$

Here (g) holds based on our earlier result in (28). (h) is a consequence of the minimax theorem, since the space of distributions over compact \mathcal{X} is convex and compact, and the objective is concave in λ and lower semicontinuous and convex in Q . (i) is implied by Lemma 3. Therefore, according to Corollary 1

$$\begin{aligned}
&\min_{P_{G(\mathbf{Z})} \in \mathcal{P}_{\mathcal{G}}} \min_{Q_{\mathbf{X}}} d_{f,c}(P_{G(\mathbf{Z})}, Q) + \max_{D \in \mathcal{F}} \{\mathbb{E}_{P_{\mathbf{X}}}[D(\mathbf{X})] - \mathbb{E}_Q[D(\mathbf{X})]\} \\
&= \min_{G \in \mathcal{G}} \max_{D \in \mathcal{F}} \mathbb{E}_{P_{\mathbf{X}}}[D(\mathbf{X})] - d_{f,c}^*(P_{G(\mathbf{Z})}, D) \\
&= \min_{G \in \mathcal{G}} \max_{D \in \mathcal{F}} \mathbb{E}_{P_{\mathbf{X}}}[D(\mathbf{X})] + \max_{\lambda \in \mathbb{R}} \lambda - \mathbb{E}[(f^* \circ (D + \lambda))^c(G(\mathbf{Z}))] \\
&= \min_{G \in \mathcal{G}} \max_{D \in \mathcal{F}, \lambda \in \mathbb{R}} \mathbb{E}_{P_{\mathbf{X}}}[D(\mathbf{X}) + \lambda] - \mathbb{E}[(f^* \circ (D + \lambda))^c(G(\mathbf{Z}))] \\
&\stackrel{(j)}{=} \min_{G \in \mathcal{G}} \max_{D \in \mathcal{F}} \mathbb{E}_{P_{\mathbf{X}}}[D(\mathbf{X})] - \mathbb{E}[(f^* \circ D)^c(G(\mathbf{Z}))] \\
&= \min_{G \in \mathcal{G}} \max_{D \in \mathcal{F}} \mathbb{E}_{P_{\mathbf{X}}}[D(\mathbf{X})] - \mathbb{E}[\sup_{\mathbf{x}'} f^*(D(\mathbf{x}')) - c(G(\mathbf{Z}), \mathbf{x}')] \\
&= \min_{G \in \mathcal{G}} \max_{D \in \mathcal{F}} \mathbb{E}_{P_{\mathbf{X}}}[D(\mathbf{X})] + \mathbb{E}[\inf_{\mathbf{x}'} -f^*(D(\mathbf{x}')) + c(G(\mathbf{Z}), \mathbf{x}')].
\end{aligned}$$

Here (j) holds since \mathcal{F} is assumed to be closed to constant additions. Hence, the proof is complete. \square

9.8 Proof of Theorem 7

Consider distributions P_0, P_1, P_2 . Let Q_0, Q_1 be the optimal solutions to the minimum sum optimization problems for $d_{f,W_2}(P_0, P_2)$ and $d_{f,W_2}(P_1, P_2)$, respectively. Then, according to the definition

$$\begin{aligned}
d_{f,W_2}(P_0, P_2) - d_{f,W_2}(P_1, P_2) &\leq W_2^2(P_0, Q_1) - W_2^2(P_1, Q_1), \\
d_{f,W_2}(P_1, P_2) - d_{f,W_2}(P_0, P_2) &\leq W_2^2(P_1, Q_0) - W_2^2(P_0, Q_0)
\end{aligned}$$

which implies

$$|d_{f,W_2}(P_0, P_2) - d_{f,W_2}(P_1, P_2)| \leq \sup_Q |W_2^2(P_0, Q) - W_2^2(P_1, Q)|.$$

Hence, for G_{θ} , $G_{\theta'}$ and any distribution P_2 we have

$$|d_{f,W_2}(P_{G_{\theta}(\mathbf{Z})}, P_2) - d_{f,W_2}(P_{G_{\theta'}(\mathbf{Z})}, P_2)| \leq \sup_Q |W_2^2(P_{G_{\theta}(\mathbf{Z})}, Q) - W_2^2(P_{G_{\theta'}(\mathbf{Z})}, Q)|. \quad (44)$$

Fix a distribution Q over the compact \mathcal{X} . Then, for any $(G_{\theta}(\mathbf{Z}), \mathbf{X}')$ whose joint distribution is in $\Pi(P_{G_{\theta}(\mathbf{Z})}, Q)$, $(G_{\theta'}(\mathbf{Z}), \mathbf{X}')$ has a joint distribution in $\Pi(P_{G_{\theta'}(\mathbf{Z})}, Q)$. Moreover, since \mathcal{X} is a

compact set in a Hilbert space, any $\mathbf{x} \in \mathcal{X}$ is norm-bounded for some finite R as $\|\mathbf{x}\| \leq R$, which implies

$$\begin{aligned}
& \left| W_2^2(P_{G_\theta(\mathbf{Z})}, Q) - W_2^2(P_{G_{\theta'}(\mathbf{Z})}, Q) \right| \\
& \leq \sup_{M_{\mathbf{Z}, \mathbf{X}' \in \Pi(P_{\mathbf{Z}}, Q)}} \left| \mathbb{E}_M \left[\|G_\theta(\mathbf{Z}) - \mathbf{X}'\|^2 - \|G_{\theta'}(\mathbf{Z}) - \mathbf{X}'\|^2 \right] \right| \\
& \leq \sup_{M_{\mathbf{Z}, \mathbf{X}' \in \Pi(P_{\mathbf{Z}}, Q)}} \mathbb{E}_M \left[\left| \|G_\theta(\mathbf{Z})\|^2 - \|G_{\theta'}(\mathbf{Z})\|^2 \right| + 2\|\mathbf{X}'\| \|G_{\theta'}(\mathbf{Z}) - G_\theta(\mathbf{Z})\| \right] \\
& \leq \mathbb{E}_{P_{\mathbf{Z}}} \left[\left| \|G_\theta(\mathbf{Z})\|^2 - \|G_{\theta'}(\mathbf{Z})\|^2 \right| + 2R \|G_{\theta'}(\mathbf{Z}) - G_\theta(\mathbf{Z})\| \right].
\end{aligned}$$

Taking a supremum over Q from both sides of the above inequality shows

$$\begin{aligned}
& \sup_Q \left| W_2^2(P_{G_\theta(\mathbf{Z})}, Q) - W_2^2(P_{G_{\theta'}(\mathbf{Z})}, Q) \right| \\
& \leq \mathbb{E}_{P_{\mathbf{Z}}} \left[\left| \|G_\theta(\mathbf{Z})\|^2 - \|G_{\theta'}(\mathbf{Z})\|^2 \right| + 2R \|G_{\theta'}(\mathbf{Z}) - G_\theta(\mathbf{Z})\| \right]. \tag{45}
\end{aligned}$$

Since G_θ changes continuously with θ , $\left| \|G_\theta(\mathbf{z})\|^2 - \|G_{\theta'}(\mathbf{z})\|^2 \right| + 2R \|G_{\theta'}(\mathbf{z}) - G_\theta(\mathbf{z})\| \rightarrow 0$ as $\theta' \rightarrow \theta$ holds pointwise. Therefore, since \mathcal{X} is compact and hence bounded, the bounded convergence theorem together with (45) implies

$$\sup_Q \left| W_2^2(P_{G_\theta(\mathbf{Z})}, Q) - W_2^2(P_{G_{\theta'}(\mathbf{Z})}, Q) \right| \xrightarrow{\theta' \rightarrow \theta} 0. \tag{46}$$

Now, combining (44) and (46) shows for any distribution P_2

$$\left| d_{f, W_2}(P_{G_\theta(\mathbf{Z})}, P_2) - d_{f, W_2}(P_{G_{\theta'}(\mathbf{Z})}, P_2) \right| \xrightarrow{\theta' \rightarrow \theta} 0. \tag{47}$$

Also, if we further assume G_θ is bounded by T locally-Lipschitz w.r.t. θ with Lipschitz constant L , then

$$\begin{aligned}
& \sup_Q \left| W_2^2(P_{G_\theta(\mathbf{Z})}, Q) - W_2^2(P_{G_{\theta'}(\mathbf{Z})}, Q) \right| \\
& \leq \mathbb{E}_{P_{\mathbf{Z}}} \left[\left| \|G_\theta(\mathbf{Z})\|^2 - \|G_{\theta'}(\mathbf{Z})\|^2 \right| + 2R \|G_{\theta'}(\mathbf{Z}) - G_\theta(\mathbf{Z})\| \right] \tag{48} \\
& \leq \mathbb{E}_{P_{\mathbf{Z}}} \left[\left(\|G_\theta(\mathbf{Z})\| + \|G_{\theta'}(\mathbf{Z})\| \right) \left(\|G_\theta(\mathbf{Z})\| - \|G_{\theta'}(\mathbf{Z})\| \right) + 2R \|G_{\theta'}(\mathbf{Z}) - G_\theta(\mathbf{Z})\| \right] \\
& \leq \mathbb{E}_{P_{\mathbf{Z}}} \left[2T \left| \|G_\theta(\mathbf{Z})\| - \|G_{\theta'}(\mathbf{Z})\| \right| + 2R \|G_{\theta'}(\mathbf{Z}) - G_\theta(\mathbf{Z})\| \right] \\
& \leq \mathbb{E}_{P_{\mathbf{Z}}} \left[2(T + R) \|G_{\theta'}(\mathbf{Z}) - G_\theta(\mathbf{Z})\| \right] \\
& \leq 2(T + R)L \|\theta' - \theta\|,
\end{aligned}$$

implying $d_{f, W_2}(P_{G_\theta(\mathbf{Z})}, Q)$ is continuous everywhere and differentiable almost everywhere as a function of θ .

10 Proof of Theorem 8

Note that applying the generalized version of Theorem 6 proved in the Appendix to difference norm-squared cost $c_2(\mathbf{x}, \mathbf{x}') = \|\mathbf{x} - \mathbf{x}'\|^2$ reveals that for a symmetric f-divergence d_f and convex set \mathcal{F} closed to constant additions the minimax problem in Theorem 1 and Corollary 1 for the mixed

divergence d_{f,c_2} reduces to

$$\begin{aligned}
& \min_{G \in \mathcal{G}} \max_{D \in \mathcal{F}} \mathbb{E}_{P_{\mathbf{X}}} [D(\mathbf{X})] + \mathbb{E} \left[\min_{\mathbf{x}'} -f^*(D(\mathbf{x}')) + c_2(G(\mathbf{Z}), \mathbf{x}') \right] \\
&= \min_{G \in \mathcal{G}} \max_{D \in \mathcal{F}} \mathbb{E}_{P_{\mathbf{X}}} [D(\mathbf{X})] + \mathbb{E} \left[\min_{\mathbf{x}'} -f^*(D(\mathbf{x}')) + \|G(\mathbf{Z}) - \mathbf{x}'\|^2 \right] \\
&= \min_{G \in \mathcal{G}} \max_{D \in \mathcal{F}} \mathbb{E}_{P_{\mathbf{X}}} [D(\mathbf{X})] + \mathbb{E} \left[\min_{\mathbf{u}} -f^*(D(G(\mathbf{Z}) + \mathbf{u})) + \|\mathbf{u}\|^2 \right].
\end{aligned} \tag{49}$$

Here the last equality follows the change of variable $\mathbf{u} = \mathbf{x}' - G(\mathbf{Z})$. Also, note that d_{f,W_2} defined in the main text is the same as the special case of the generalized hybrid divergence $d_{f,c}$ with cost c_2 . Hence, the proof is complete.

10.1 Two additional examples for convex duality framework applied to Wasserstein distances

10.1.1 Total variation distance: Energy-based GAN

Consider the total variation distance $\delta(P, Q)$ which is defined as

$$\delta(P, Q) := \sup_{A \in \Sigma} |P(A) - Q(A)|, \tag{50}$$

where Σ is the set all Borel subsets of support set \mathcal{X} . More generally we consider $\delta_m(P, Q) = m\delta(P, Q)$ for any positive $m > 0$. Under mild assumptions, the total variation distance can be cast as a Wasserstein distance for the indicator cost $c_{m,I}(\mathbf{x}, \mathbf{x}') = m \mathbb{I}(\mathbf{x} \neq \mathbf{x}')$ [16], i.e. $\delta_m(P, Q) = OT_{c_{m,I}}(P, Q)$. Note that $c_{m,I}$ is a lower semicontinuous distance function, and hence Lemma 3 applies to $c_{m,I}$ indicating

$$\begin{aligned}
\delta_{mP}^*(D) &= OT_{c_{I,m}^*}^*(D) \\
&= \mathbb{E}_P [D^{c_{I,m}^*}(\mathbf{X})] \\
&= \mathbb{E}_P \left[\sup_{\mathbf{x}'} D(\mathbf{x}') - m c_I(\mathbf{X}, \mathbf{x}') \right] \\
&= \mathbb{E}_P \left[\max \left\{ D(\mathbf{X}), \max_{\mathbf{x}'} D(\mathbf{x}') - m \right\} \right] \\
&= \mathbb{E}_P \left[\max \left\{ m + D(\mathbf{X}) - \max_{\mathbf{x}'} D(\mathbf{x}'), 0 \right\} \right] + \max_{\mathbf{x}'} D(\mathbf{x}') - m
\end{aligned}$$

Without loss of generality, we can assume that the maximum discriminator output is always 0 which results in

$$\delta_{mP}^*(D) = \mathbb{E}_P \left[\max \left\{ m + D(\mathbf{X}), 0 \right\} \right] - m$$

Therefore, the minimax problem in Corollaries 1,2 for the total variation distance will be

$$\begin{aligned}
& \min_{G \in \mathcal{G}} \max_{D \in \mathcal{F}} \mathbb{E}_P [D(\mathbf{X})] - \delta_{mP}^*(D) \\
&= \min_{G \in \mathcal{G}} \max_{D \in \mathcal{F}} \mathbb{E}_P [D(\mathbf{X})] - \mathbb{E}_P \left[\max \left\{ m + D(G(\mathbf{Z})), 0 \right\} \right] + m \\
&= \min_{G \in \mathcal{G}} \max_{-D \in \mathcal{F}} -\mathbb{E}_P [D(\mathbf{X})] - \mathbb{E}_P \left[\max \left\{ m - D(G(\mathbf{Z})), 0 \right\} \right] + m \\
&= \min_{G \in \mathcal{G}} \max_{\tilde{D} \in \mathcal{F}} -\mathbb{E}_P [\tilde{D}(\mathbf{X})] - \mathbb{E}_P \left[\max \left\{ m - \tilde{D}(G(\mathbf{Z})), 0 \right\} \right] + m
\end{aligned}$$

where the last equality follows from the assumption that for any $D \in \mathcal{F}$ we have $-D \in \mathcal{F}$. Since D is assumed to be non-positive, \tilde{D} takes non-negative values. Note that this problem is equivalent to a minimax game where discriminator D is *minimizing* the following cost over \mathcal{F} :

$$L_D(G, D) = \mathbb{E}_P [D(\mathbf{X})] + \mathbb{E}_P \left[\max \left\{ m - D(G(\mathbf{Z})), 0 \right\} \right] \tag{51}$$

which is also the discriminator cost function in the energy-based GAN [8]. Hence, for any fixed $G \in \mathcal{G}$, the optimal discriminator $D \in \mathcal{F}$ for the total variation's minimax problem is the same as the energy-based GAN's optimal discriminator.

10.1.2 Second-order Wasserstein distance: the LQG setting

Consider the second-order Wasserstein distance $W_2(P, Q)$, and suppose \mathcal{F} is the set of quadratic functions over \mathbf{X} , which is a linear space. Also assume the generator G is a linear function and the r -dimensional noise \mathbf{Z} is Gaussianly-distributed with zero-mean and identity covariance matrix $I_{r \times r}$. According to the interpretation provided in Corollary 2, the second-order Wasserstein GAN finds the multivariate Gaussian distribution with rank r covariance matrix minimizing the W_2 distance to the set of distributions with their second-order moments matched to $P_{\mathbf{X}}$'s moments.

Since the value of $\mathbb{E}[\|\mathbf{X} - G(\mathbf{Z})\|^2]$ depends only on the second-order moments of the vector $[\mathbf{X}, G(\mathbf{Z})]$, we can minimize the W_2 -distance between the two sets by minimizing this expectation over Gaussianly-distributed vectors $[\mathbf{X}, G(\mathbf{Z})]$ subject to a rank r covariance matrix for $G(\mathbf{Z})$ and a pre-determined covariance matrix for $[\mathbf{X}]$. Hence, the optimal G^* simply corresponds to the r -PCA solution for $P_{\mathbf{X}}$.

This example shows Theorem 3 provides another way to recover [9]'s main result under the linear generator, quadratic discriminator and Gaussianly-distributed data assumptions.