

Wasserstein GAN

Martin Arjovsky¹, Soumith Chintala², and Léon Bottou^{1,2}

¹Courant Institute of Mathematical Sciences

²Facebook AI Research

1 Introduction

The problem this paper is concerned with is that of unsupervised learning. Mainly, what does it mean to learn a probability distribution? The classical answer to this is to learn a probability density. This is often done by defining a parametric family of densities $(P_\theta)_{\theta \in \mathbb{R}^d}$ and finding the one that maximized the likelihood on our data: if we have real data examples $\{x^{(i)}\}_{i=1}^m$, we would solve the problem

$$\max_{\theta \in \mathbb{R}^d} \frac{1}{m} \sum_{i=1}^m \log P_\theta(x^{(i)})$$

If the real data distribution \mathbb{P}_r admits a density and \mathbb{P}_θ is the distribution of the parametrized density P_θ , then, asymptotically, this amounts to minimizing the Kullback-Leibler divergence $KL(\mathbb{P}_r \parallel \mathbb{P}_\theta)$.

For this to make sense, we need the model density P_θ to exist. This is not the case in the rather common situation where we are dealing with distributions supported by low dimensional manifolds. It is then unlikely that the model manifold and the true distribution's support have a non-negligible intersection (see [1]), and this means that the KL distance is not defined (or simply infinite).

The typical remedy is to add a noise term to the model distribution. This is why virtually all generative models described in the classical machine learning literature include a noise component. In the simplest case, one assumes a Gaussian noise with relatively high bandwidth in order to cover all the examples. It is well known, for instance, that in the case of image generation models, this noise degrades the quality of the samples and makes them blurry. For example, we can see in the recent paper [22] that the optimal standard deviation of the noise added to the model when maximizing likelihood is around 0.1 to each pixel in a generated image, when the pixels were already normalized to be in the range $[0, 1]$. This is a very high amount of noise, so much that when papers report the samples of their models, they don't add the noise term on which they report likelihood numbers. In other words, the added noise term is clearly incorrect for the problem, but is needed to make the maximum likelihood approach work.

Rather than estimating the density of \mathbb{P}_r which may not exist, we can define a random variable Z with a fixed distribution $p(z)$ and pass it through a parametric function $g_\theta : \mathcal{Z} \rightarrow \mathcal{X}$ (typically a neural network of some kind) that directly generates samples following a certain distribution \mathbb{P}_θ . By varying θ , we can change this distribution and make it close to the real data distribution \mathbb{P}_r . This is useful in two ways. First of all, unlike densities, this approach can represent distributions confined to a low dimensional manifold. Second, the ability to easily generate samples is often more useful than knowing the numerical value of the density (for example in image superresolution or semantic segmentation when considering the conditional distribution of the output image given the input image). In general, it is computationally difficult to generate samples given an arbitrary high dimensional density [15].

Variational Auto-Encoders (VAEs) [9] and Generative Adversarial Networks (GANs) [4] are well known examples of this approach. Because VAEs focus on the approximate likelihood of the examples, they share the limitation of the standard models and need to fiddle with additional noise terms. GANs offer much more flexibility in the definition of the objective function, including Jensen-Shannon [4], and all f -divergences [16] as well as some exotic combinations [6]. On the other hand, training GANs is well known for being delicate and unstable, for reasons theoretically investigated in [1].

In this paper, we direct our attention on the various ways to measure how close the model distribution and the real distribution are, or equivalently, on the various ways to define a distance or divergence $\rho(\mathbb{P}_\theta, \mathbb{P}_r)$. The most fundamental difference between such distances is their impact on the convergence of sequences of probability distributions. A sequence of distributions $(\mathbb{P}_t)_{t \in \mathbb{N}}$ converges if and only if there is a distribution \mathbb{P}_∞ such that $\rho(\mathbb{P}_t, \mathbb{P}_\infty)$ tends to zero, something that depends on how exactly the distance ρ is defined. Informally, a distance ρ induces a weaker topology when it makes it easier for a sequence of distribution to converge.¹ Section 2 clarifies how popular probability distances differ in that respect.

In order to optimize the parameter θ , it is of course desirable to define our model distribution \mathbb{P}_θ in a manner that makes the mapping $\theta \mapsto \mathbb{P}_\theta$ continuous. Continuity means that when a sequence of parameters θ_t converges to θ , the distributions \mathbb{P}_{θ_t} also converge to \mathbb{P}_θ . However, it is essential to remember that the notion of the convergence of the distributions \mathbb{P}_{θ_t} depends on the way we compute the distance between distributions. The weaker this distance, the easier it is to define a continuous mapping from θ -space to \mathbb{P}_θ -space, since it's easier for the distributions to converge. The main reason we care about the mapping $\theta \mapsto \mathbb{P}_\theta$ to be continuous is as follows. If ρ is our notion of distance between two distributions, we would like to have a loss function $\theta \mapsto \rho(\mathbb{P}_\theta, \mathbb{P}_r)$ that is continuous, and this is equivalent to having the mapping $\theta \mapsto \mathbb{P}_\theta$ be continuous when using the distance between distributions ρ .

¹More exactly, the topology induced by ρ is weaker than that induced by ρ' when the set of convergent sequences under ρ is a superset of that under ρ' .

The contributions of this paper are:

- In Section 2, we provide a comprehensive theoretical analysis of how the Earth Mover (EM) distance behaves in comparison to popular probability distances and divergences used in the context of learning distributions.
- In Section 3, we define a form of GAN called Wasserstein-GAN that minimizes a reasonable and efficient approximation of the EM distance, and we theoretically show that the corresponding optimization problem is sound.
- In Section 4, we empirically show that WGANs cure the main training problems of GANs. In particular, training WGANs does not require maintaining a careful balance in training of the discriminator and the generator, and does not require a careful design of the network architecture either. The mode dropping phenomenon that is typical in GANs is also drastically reduced. One of the most compelling practical benefits of WGANs is the ability to continuously estimate of the EM distance by training the discriminator to optimality. Plotting these learning curves is not only useful for debugging and hyperparameter searches, but also correlate remarkably well with the observed sample quality.

2 Different Distances

We now introduce our notation. Let \mathcal{X} be a compact metric set (such as the space of images $[0, 1]^d$) and let Σ denote the set of all the Borel subsets of \mathcal{X} . Let $\text{Prob}(\mathcal{X})$ denote the space of probability measures defined on \mathcal{X} . We can now define elementary distances and divergences between two distributions $\mathbb{P}_r, \mathbb{P}_g \in \text{Prob}(\mathcal{X})$:

- The *Total Variation* (TV) distance

$$\delta(\mathbb{P}_r, \mathbb{P}_g) = \sup_{A \in \Sigma} |\mathbb{P}_r(A) - \mathbb{P}_g(A)| .$$

- The *Kullback-Leibler* (KL) divergence

$$KL(\mathbb{P}_r \parallel \mathbb{P}_g) = \int \log \left(\frac{P_r(x)}{P_g(x)} \right) P_r(x) d\mu(x) ,$$

where both \mathbb{P}_r and \mathbb{P}_g are assumed to be absolutely continuous, and therefore admit densities, with respect to a same measure μ defined on \mathcal{X} .² The KL divergence is famously assymmetric and possibly infinite when there are points such that $P_g(x) = 0$ and $P_r(x) > 0$.

²Recall that a probability distribution $\mathbb{P}_r \in \text{Prob}(\mathcal{X})$ admits a density $p_r(x)$ with respect to μ , that is, $\forall A \in \Sigma, \mathbb{P}_r(A) = \int_A P_r(x) d\mu(x)$, if and only if it is absolutely continuous with respect to μ , that is, $\forall A \in \Sigma, \mu(A) = 0 \Rightarrow \mathbb{P}_r(A) = 0$.

- The *Jensen-Shannon* (JS) divergence

$$JS(\mathbb{P}_r, \mathbb{P}_g) = KL(\mathbb{P}_r \| \mathbb{P}_m) + KL(\mathbb{P}_g \| \mathbb{P}_m) ,$$

where \mathbb{P}_m is the mixture $(\mathbb{P}_r + \mathbb{P}_g)/2$. This divergence is symmetrical and always defined because we can choose $\mu = \mathbb{P}_m$.

- The *Earth-Mover* (EM) distance or Wasserstein-1

$$W(\mathbb{P}_r, \mathbb{P}_g) = \inf_{\gamma \in \Pi(\mathbb{P}_r, \mathbb{P}_g)} \mathbb{E}_{(x,y) \sim \gamma} [\|x - y\|] , \quad (1)$$

where $\Pi(\mathbb{P}_r, \mathbb{P}_g)$ denotes the set of all joint distributions $\gamma(x, y)$ whose marginals are respectively \mathbb{P}_r and \mathbb{P}_g . Intuitively, $\gamma(x, y)$ indicates how much “mass” must be transported from x to y in order to transform the distributions \mathbb{P}_r into the distribution \mathbb{P}_g . The EM distance then is the “cost” of the optimal transport plan.

The following example illustrates how apparently simple sequences of probability distributions converge under the EM distance but do not converge under the other distances and divergences defined above.

Example 1 (Learning parallel lines). Let $Z \sim U[0, 1]$ the uniform distribution on the unit interval. Let \mathbb{P}_0 be the distribution of $(0, Z) \in \mathbb{R}^2$ (a 0 on the x-axis and the random variable Z on the y-axis), uniform on a straight vertical line passing through the origin. Now let $g_\theta(z) = (\theta, z)$ with θ a single real parameter. It is easy to see that in this case,

- $W(\mathbb{P}_0, \mathbb{P}_\theta) = |\theta|$,
- $JS(\mathbb{P}_0, \mathbb{P}_\theta) = \begin{cases} \log 2 & \text{if } \theta \neq 0 , \\ 0 & \text{if } \theta = 0 , \end{cases}$
- $KL(\mathbb{P}_\theta \| \mathbb{P}_0) = KL(\mathbb{P}_0 \| \mathbb{P}_\theta) = \begin{cases} +\infty & \text{if } \theta \neq 0 , \\ 0 & \text{if } \theta = 0 , \end{cases}$
- and $\delta(\mathbb{P}_0, \mathbb{P}_\theta) = \begin{cases} 1 & \text{if } \theta \neq 0 , \\ 0 & \text{if } \theta = 0 . \end{cases}$

When $\theta_t \rightarrow 0$, the sequence $(\mathbb{P}_{\theta_t})_{t \in \mathbb{N}}$ converges to \mathbb{P}_0 under the EM distance, but does not converge at all under either the JS, KL, reverse KL, or TV divergences. Figure 1 illustrates this for the case of the EM and JS distances.

Example 1 gives us a case where we can learn a probability distribution over a low dimensional manifold by doing gradient descent on the EM distance. This cannot be done with the other distances and divergences because the resulting loss function is not even continuous. Although this simple example features distributions with disjoint supports, the same conclusion holds when the supports have a non empty

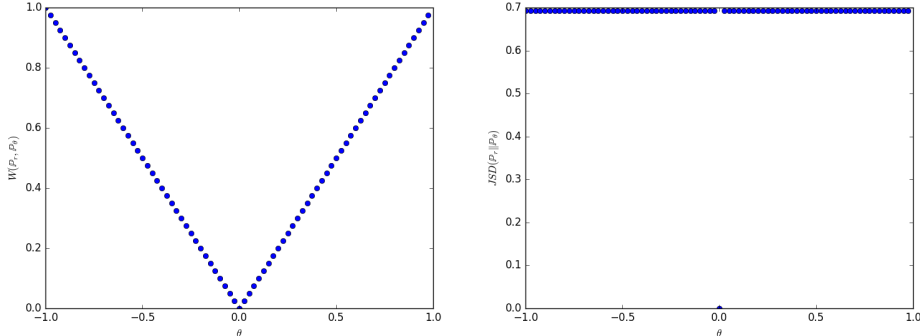


Figure 1: These plots show $\rho(\mathbb{P}_\theta, \mathbb{P}_0)$ as a function of θ when ρ is the EM distance (left plot) or the JS divergence (right plot). The EM plot is continuous and provides a usable gradient everywhere. The JS plot is not continuous and does not provide a usable gradient.

intersection contained in a set of measure zero. This happens to be the case when two low dimensional manifolds intersect in general position [1].

Since the Wasserstein distance is much weaker than the JS distance³, we can now ask whether $W(\mathbb{P}_r, \mathbb{P}_\theta)$ is a continuous loss function on θ under mild assumptions. This, and more, is true, as we now state and prove.

Theorem 1. *Let \mathbb{P}_r be a fixed distribution over \mathcal{X} . Let Z be a random variable (e.g Gaussian) over another space \mathcal{Z} . Let $g : \mathcal{Z} \times \mathbb{R}^d \rightarrow \mathcal{X}$ be a function, that will be denoted $g_\theta(z)$ with z the first coordinate and θ the second. Let \mathbb{P}_θ denote the distribution of $g_\theta(Z)$. Then,*

1. *If g is continuous in θ , so is $W(\mathbb{P}_r, \mathbb{P}_\theta)$.*
2. *If g is locally Lipschitz and satisfies regularity assumption 1, then $W(\mathbb{P}_r, \mathbb{P}_\theta)$ is continuous everywhere, and differentiable almost everywhere.*
3. *Statements 1-2 are false for the Jensen-Shannon divergence $JS(\mathbb{P}_r, \mathbb{P}_\theta)$ and all the KLS.*

Proof. See Appendix C □

The following corollary tells us that learning by minimizing the EM distance makes sense (at least in theory) with neural networks.

Corollary 1. *Let g_θ be any feedforward neural network⁴ parameterized by θ , and $p(z)$ a prior over z such that $\mathbb{E}_{z \sim p(z)}[\|z\|] < \infty$ (e.g. Gaussian, uniform, etc.).*

³ The argument for *why* this happens, and indeed how we arrived to the idea that Wasserstein is what we should really be optimizing is displayed in Appendix A. We strongly encourage the interested reader who is not afraid of the mathematics to go through it.

⁴By a feedforward neural network we mean a function composed by affine transformations and pointwise nonlinearities which are smooth Lipschitz functions (such as the sigmoid, tanh, elu, softplus, etc). **Note:** the statement is also true for rectifier nonlinearities but the proof is more technical (even though very similar) so we omit it.

Then assumption 1 is satisfied and therefore $W(\mathbb{P}_r, \mathbb{P}_\theta)$ is continuous everywhere and differentiable almost everywhere.

Proof. See Appendix C □

All this shows that EM is a much more sensible cost function for our problem than at least the Jensen-Shannon divergence. The following theorem describes the relative strength of the topologies induced by these distances and divergences, with KL the strongest, followed by JS and TV, and EM the weakest.

Theorem 2. *Let \mathbb{P} be a distribution on a compact space \mathcal{X} and $(\mathbb{P}_n)_{n \in \mathbb{N}}$ be a sequence of distributions on \mathcal{X} . Then, considering all limits as $n \rightarrow \infty$,*

1. *The following statements are equivalent*
 - $\delta(\mathbb{P}_n, \mathbb{P}) \rightarrow 0$ with δ the total variation distance.
 - $JS(\mathbb{P}_n, \mathbb{P}) \rightarrow 0$ with JS the Jensen-Shannon divergence.
2. *The following statements are equivalent*
 - $W(\mathbb{P}_n, \mathbb{P}) \rightarrow 0$.
 - $\mathbb{P}_n \xrightarrow{\mathcal{D}} \mathbb{P}$ where $\xrightarrow{\mathcal{D}}$ represents convergence in distribution for random variables.
3. *KL($\mathbb{P}_n \parallel \mathbb{P}$) $\rightarrow 0$ or KL($\mathbb{P} \parallel \mathbb{P}_n$) $\rightarrow 0$ imply the statements in (1).*
4. *The statements in (1) imply the statements in (2).*

Proof. See Appendix C □

This highlights the fact that the KL, JS, and TV distances are not sensible cost functions when learning distributions supported by low dimensional manifolds. However the EM distance is sensible in that setup. This obviously leads us to the next section where we introduce a practical approximation of optimizing the EM distance.

3 Wasserstein GAN

Again, Theorem 2 points to the fact that $W(\mathbb{P}_r, \mathbb{P}_\theta)$ might have nicer properties when optimized than $JS(\mathbb{P}_r, \mathbb{P}_\theta)$. However, the infimum in (1) is highly intractable. On the other hand, the Kantorovich-Rubinstein duality [21] tells us that

$$W(\mathbb{P}_r, \mathbb{P}_\theta) = \sup_{\|f\|_L \leq 1} \mathbb{E}_{x \sim \mathbb{P}_r}[f(x)] - \mathbb{E}_{x \sim \mathbb{P}_\theta}[f(x)] \quad (2)$$

where the supremum is over all the 1-Lipschitz functions $f : \mathcal{X} \rightarrow \mathbb{R}$. Note that if we replace $\|f\|_L \leq 1$ for $\|f\|_L \leq K$ (consider K -Lipschitz for some constant K), then we end up with $K \cdot W(\mathbb{P}_r, \mathbb{P}_\theta)$. Therefore, if we have a parameterized family of

functions $\{f_w\}_{w \in \mathcal{W}}$ that are all K -Lipschitz for some K , we could consider solving the problem

$$\max_{w \in \mathcal{W}} \mathbb{E}_{x \sim \mathbb{P}_r}[f_w(x)] - \mathbb{E}_{z \sim p(z)}[f_w(g_\theta(z))] \quad (3)$$

and if the supremum in (2) is attained for some $w \in \mathcal{W}$ (a pretty strong assumption akin to what’s assumed when proving consistency of an estimator), this process would yield a calculation of $W(\mathbb{P}_r, \mathbb{P}_\theta)$ up to a multiplicative constant. Furthermore, we could consider differentiating $W(\mathbb{P}_r, \mathbb{P}_\theta)$ (again, up to a constant) by back-proping through equation (2) via estimating $\mathbb{E}_{z \sim p(z)}[\nabla_\theta f_w(g_\theta(z))]$. While this is all intuition, we now prove that this process is principled under the optimality assumption.

Theorem 3. *Let \mathbb{P}_r be any distribution. Let \mathbb{P}_θ be the distribution of $g_\theta(Z)$ with Z a random variable with density p and g_θ a function satisfying assumption 1. Then, there is a solution $f : \mathcal{X} \rightarrow \mathbb{R}$ to the problem*

$$\max_{\|f\|_L \leq 1} \mathbb{E}_{x \sim \mathbb{P}_r}[f(x)] - \mathbb{E}_{x \sim \mathbb{P}_\theta}[f(x)]$$

and we have

$$\nabla_\theta W(\mathbb{P}_r, \mathbb{P}_\theta) = -\mathbb{E}_{z \sim p(z)}[\nabla_\theta f(g_\theta(z))]$$

when both terms are well-defined.

Proof. See Appendix C □

Now comes the question of finding the function f that solves the maximization problem in equation (2). To roughly approximate this, something that we can do is train a neural network parameterized with weights w lying in a compact space \mathcal{W} and then backprop through $\mathbb{E}_{z \sim p(z)}[\nabla_\theta f_w(g_\theta(z))]$, as we would do with a typical GAN. Note that the fact that \mathcal{W} is compact implies that all the functions f_w will be K -Lipschitz for some K that only depends on \mathcal{W} and not the individual weights, therefore approximating (2) up to an irrelevant scaling factor and the capacity of the ‘critic’ f_w . In order to have parameters w lie in a compact space, something simple we can do is clamp the weights to a fixed box (say $\mathcal{W} = [-0.01, 0.01]^l$) after each gradient update. The Wasserstein Generative Adversarial Network (WGAN) procedure is described in Algorithm 1.

The fact that the EM distance is continuous and differentiable a.e. means that we can (and should) train the critic till optimality. The argument is simple, the more we train the critic, the more reliable gradient of the Wasserstein we get, which is actually useful by the fact that Wasserstein is differentiable almost everywhere. For the JS, as the discriminator gets better the gradients get more reliable but the true gradient is 0 since the JS is locally saturated and we get vanishing gradients, as can be seen in Figure 1 of this paper and Theorem 2.4 of [1]. In Figure 2 we show a proof of concept of this, where we train a GAN discriminator and a WGAN critic till optimality. The discriminator learns very quickly to distinguish between fake and real, and as expected provides no reliable gradient information. The critic, however, can’t saturate, and converges to a linear function that gives

Algorithm 1 WGAN, our proposed algorithm. All experiments in the paper used the default values $\alpha = 0.00005$, $c = 0.01$, $m = 64$, $n_{\text{critic}} = 5$.

Require: : α , the learning rate. c , the clipping parameter. m , the batch size. n_{critic} , the number of iterations of the critic per generator iteration.

Require: : w_0 , initial critic parameters. θ_0 , initial generator’s parameters.

```

1: while  $\theta$  has not converged do
2:   for  $t = 0, \dots, n_{\text{critic}}$  do
3:     Sample  $\{x^{(i)}\}_{i=1}^m \sim \mathbb{P}_r$  a batch from the real data.
4:     Sample  $\{z^{(i)}\}_{i=1}^m \sim p(z)$  a batch of prior samples.
5:      $g_w \leftarrow \nabla_w \left[ \frac{1}{m} \sum_{i=1}^m f_w(x^{(i)}) - \frac{1}{m} \sum_{i=1}^m f_w(g_\theta(z^{(i)})) \right]$ 
6:      $w \leftarrow w + \alpha \cdot \text{RMSPProp}(w, g_w)$ 
7:      $w \leftarrow \text{clip}(w, -c, c)$ 
8:   end for
9:   Sample  $\{z^{(i)}\}_{i=1}^m \sim p(z)$  a batch of prior samples.
10:   $g_\theta \leftarrow -\nabla_\theta \frac{1}{m} \sum_{i=1}^m f_w(g_\theta(z^{(i)}))$ 
11:   $\theta \leftarrow \theta - \alpha \cdot \text{RMSPProp}(\theta, g_\theta)$ 
12: end while

```

remarkably clean gradients everywhere. The fact that we constrain the weights limits the possible growth of the function to be at most linear in different parts of the space, forcing the optimal critic to have this behaviour.

Perhaps more importantly, the fact that we can train the critic till optimality makes it impossible to collapse modes when we do. This is due to the fact that mode collapse comes from the fact that the optimal generator for a *fixed* discriminator is a sum of deltas on the points the discriminator assigns the highest values, as brilliantly observed by [11].

In the following section we display the practical benefits of our new algorithm, and we provide an in-depth comparison of its behaviour and that of traditional GANs.

4 Empirical Results

We run experiments on image generation using our Wasserstein-GAN algorithm and show that there are significant practical benefits to using it over the formulation used in standard GANs.

We claim two main benefits:

- a meaningful loss metric that correlates with the generator’s convergence and sample quality
- improved stability of the optimization process

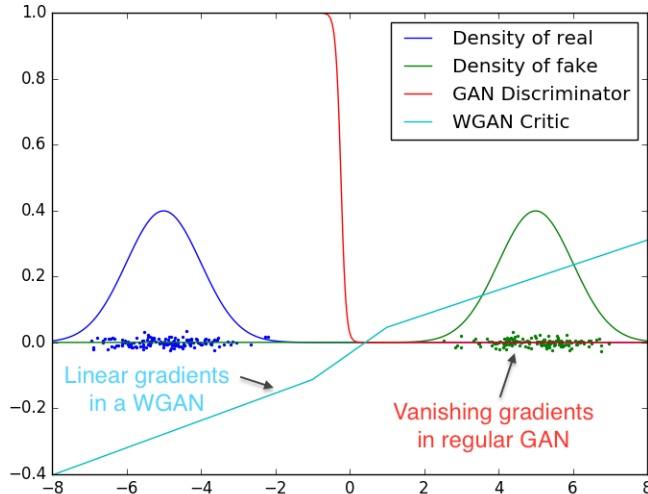


Figure 2: Optimal discriminator and critic when learning to differentiate two Gaussians. As we can see, the traditional GAN discriminator saturates and results in vanishing gradients. Our WGAN critic provides very clean gradients on all parts of the space.

4.1 Experimental Procedure

We run experiments on image generation. The target distribution to learn is the LSUN-Bedrooms dataset [23] – a collection of natural images of indoor bedrooms. Our baseline comparison is DCGAN [17], a GAN with a convolutional architecture trained with the standard GAN procedure using the $-\log D$ trick [4, 1]. The generated samples are 3-channel images of 64x64 pixels in size. We use the hyperparameters specified in Algorithm 1 for all of our experiments.

4.2 Meaningful loss metric

Because the WGAN algorithm attempts to train the critic f (lines 2–8 in Algorithm 1) relatively well before each generator update (line 10 in Algorithm 1), the loss function at this point is an estimate of the EM distance, up to constant factors related to the way we constrain the Lipschitz constant of f .

Our first experiment illustrates how this estimate correlates well with the quality of the generated samples. Besides the convolutional DCGAN architecture, we also ran experiments where we replace the generator or both the generator and the critic by 4-layer ReLU-MLP with 512 hidden units.

Figure 3 plots the evolution of the WGAN estimate (3) of the EM distance during WGAN training for all three architectures. The plots clearly show that these curves correlate well with the visual quality of the generated samples.

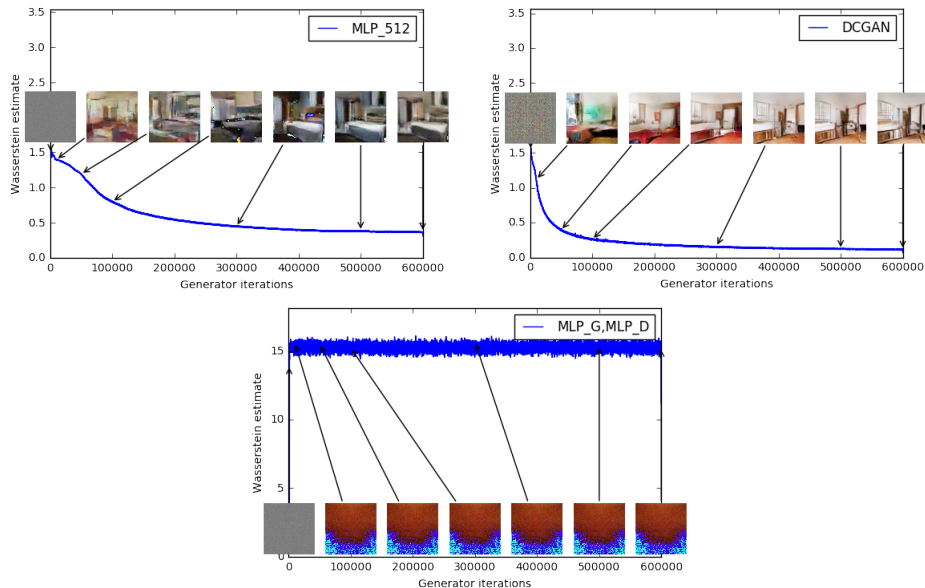


Figure 3: Training curves and samples at different stages of training. We can see a clear correlation between lower error and better sample quality. Upper left: the generator is an MLP with 4 hidden layers and 512 units at each layer. The loss decreases consistently as training progresses and sample quality increases. Upper right: the generator is a standard DCGAN. The loss decreases quickly and sample quality increases as well. In both upper plots the critic is a DCGAN without the sigmoid so losses can be subjected to comparison. Lower half: both the generator and the discriminator are MLPs with substantially high learning rates (so training failed). Loss is constant and samples are constant as well. The training curves were passed through a median filter for visualization purposes.

To our knowledge, this is the first time in GAN literature that such a property is shown, where the loss of the GAN shows properties of convergence. This property is extremely useful when doing research in adversarial networks as one does not need to stare at the generated samples to figure out failure modes and to gain information on which models are doing better over others.

However, we do not claim that this is a new method to quantitatively evaluate generative models yet. The constant scaling factor that depends on the critic’s architecture means it’s hard to compare models with different critics. Even more, in practice the fact that the critic doesn’t have infinite capacity makes it hard to know just how close to the EM distance our estimate really is. This being said, we have successfully used the loss metric to validate our experiments repeatedly and without failure, and we see this as a huge improvement in training GANs which previously had no such facility.

In contrast, Figure 4 plots the evolution of the GAN estimate of the JS distance during GAN training. More precisely, during GAN training, the discriminator is

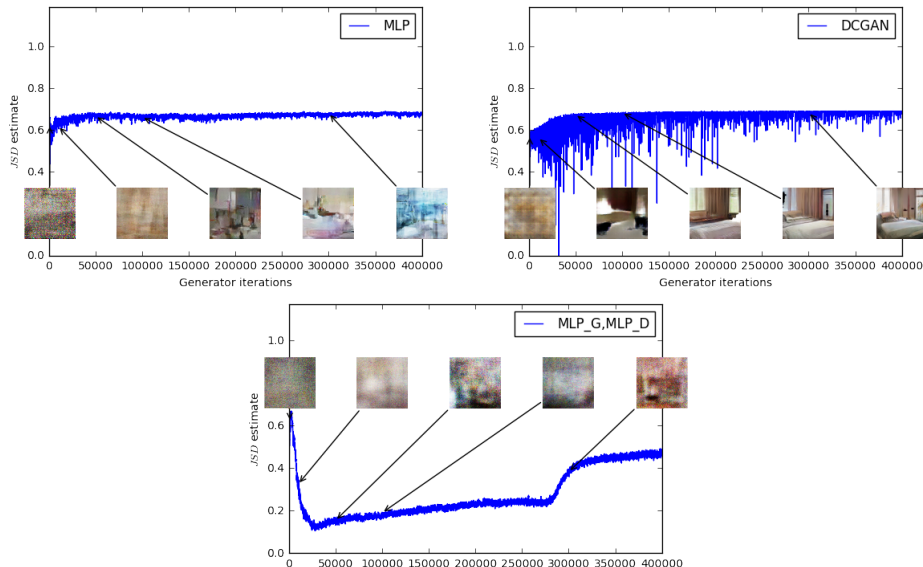


Figure 4: JS estimates for an MLP generator (upper left) and a DCGAN generator (upper right) trained with the standard GAN procedure. Both had a DCGAN discriminator. Both curves have increasing error. Samples get better for the DCGAN but the JS estimate increases or stays constant, pointing towards no significant correlation between sample quality and loss. Bottom: MLP with both generator and discriminator. The curve goes up and down regardless of sample quality. All training curves were passed through the same median filter as in Figure 3.

trained to maximize

$$L(D, g_\theta) = \mathbb{E}_{x \sim \mathbb{P}_r} [\log D(x)] + \mathbb{E}_{x \sim \mathbb{P}_\theta} [\log(1 - D(x))]$$

which is a lower bound of $2JS(\mathbb{P}_r, \mathbb{P}_\theta) - 2 \log 2$. In the figure, we plot the quantity $\frac{1}{2}L(D, g_\theta) + \log 2$, which is a lower bound of the JS distance.

This quantity clearly correlates poorly the sample quality. Note also that the JS estimate usually stays constant or goes up instead of going down. In fact it often remains very close to $\log 2 \approx 0.69$ which is the highest value taken by the JS distance. In other words, the JS distance saturates, the discriminator has zero loss, and the generated samples are in some cases meaningful (DCGAN generator, top right plot) and in other cases collapse to a single nonsensical image. This last phenomenon has been theoretically explained in [1] and highlighted in [11].

When using the $-\log D$ trick [4, 1], the discriminator loss and the generator loss are different. Figure 8 in Appendix E reports the same plots for GAN training, but using the generator loss instead of the discriminator loss. This does not change the conclusions.

Finally, as a negative result, we report that WGAN training becomes unstable at times when one uses a momentum based optimizer such as Adam [8] on the critic,

or when one uses high learning rates. Since the loss for the critic is nonstationary, momentum based methods like Adam seemed to perform worse. We identified momentum as a potential cause because, as the loss blew up and samples got worse, the cosine between the Adam step and the gradient turned negative. The only places where this cosine was negative was in these situations of instability. We therefore switched to RMSProp [20] which is known to perform well even on very nonstationary problems [13].

4.3 Improved stability

One of the benefits of WGAN is that it allows us to train the critic till optimality. When the critic is trained to completion, it simply provides a loss to the generator that we can train as any other neural network. This tells us that we no longer need to balance generator and discriminator’s capacity properly. The better the critic, the higher quality the gradients we use to train the generator.

We observe that WGANs are much more robust than GANs when one varies the architectural choices for the generator. We illustrate this by running experiments on three generator architectures: (1) a convolutional DCGAN generator, (2) a convolutional DCGAN generator without batch normalization and with a constant number of filters, and (3) a 4-layer ReLU-MLP with 512 hidden units. The last two are known to perform very poorly with GANs. We keep the convolutional DCGAN architecture for the WGAN critic or the GAN discriminator.

Figures 5, 6, and 7 show samples generated for these three architectures using both the WGAN and GAN algorithms. We refer the reader to Appendix F for full sheets of generated samples. Samples were not cherry-picked.

In no experiment did we see evidence of mode collapse for the WGAN algorithm.

5 Related Work

There’s been a number of works on the so called Integral Probability Metrics (IPMs) [14]. Given \mathcal{F} a set of functions from \mathcal{X} to \mathbb{R} , we can define

$$d_{\mathcal{F}}(\mathbb{P}_r, \mathbb{P}_\theta) = \sup_{f \in \mathcal{F}} \mathbb{E}_{x \sim \mathbb{P}_r}[f(x)] - \mathbb{E}_{x \sim \mathbb{P}_\theta}[f(x)] \quad (4)$$

as an integral probability metric associated with the function class \mathcal{F} . It is easily verified that if for every $f \in \mathcal{F}$ we have $-f \in \mathcal{F}$ (such as all examples we’ll consider), then $d_{\mathcal{F}}$ is nonnegative, satisfies the triangular inequality, and is symmetric. Thus, $d_{\mathcal{F}}$ is a pseudometric over $\text{Prob}(\mathcal{X})$.

While IPMs might seem to share a similar formula, as we will see different classes of functions can yeald to radically different metrics.

- By the Kantorovich-Rubinstein duality [21], we know that $W(\mathbb{P}_r, \mathbb{P}_\theta) = d_{\mathcal{F}}(\mathbb{P}_r, \mathbb{P}_\theta)$ when \mathcal{F} is the set of 1-Lipschitz functions. Furthermore, if \mathcal{F} is the set of K -Lipschitz functions, we get $K \cdot W(\mathbb{P}_r, \mathbb{P}_\theta) = d_{\mathcal{F}}(\mathbb{P}_r, \mathbb{P}_\theta)$.

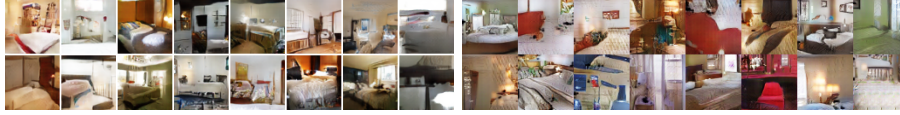


Figure 5: Algorithms trained with a DCGAN generator. Left: WGAN algorithm. Right: standard GAN formulation. Both algorithms produce high quality samples.



Figure 6: Algorithms trained with a generator without batch normalization and constant number of filters at every layer (as opposed to duplicating them every time as in [17]). Aside from taking out batch normalization, the number of parameters is therefore reduced by a bit more than an order of magnitude. Left: WGAN algorithm. Right: standard GAN formulation. As we can see the standard GAN failed to learn while the WGAN still was able to produce samples.

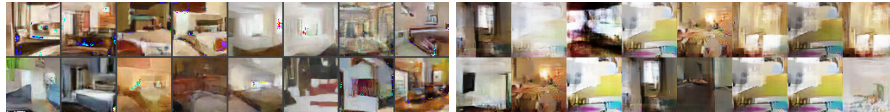


Figure 7: Algorithms trained with an MLP generator with 4 layers and 512 units with ReLU nonlinearities. The number of parameters is similar to that of a DCGAN, but it lacks a strong inductive bias for image generation. Left: WGAN algorithm. Right: standard GAN formulation. The WGAN method still was able to produce samples, lower quality than the DCGAN, and of higher quality than the MLP of the standard GAN. Note the significant degree of mode collapse in the GAN MLP.

- When \mathcal{F} is the set of all measurable functions bounded between -1 and 1 (or all continuous functions between -1 and 1), we retrieve $d_{\mathcal{F}}(\mathbb{P}_r, \mathbb{P}_\theta) = \delta(\mathbb{P}_r, \mathbb{P}_\theta)$ the total variation distance [14]. This already tells us that going from 1-Lipschitz to 1-Bounded functions drastically changes the topology of the space, and the regularity of $d_{\mathcal{F}}(\mathbb{P}_r, \mathbb{P}_\theta)$ as a loss function (as by Theorems 1 and 2).
- Energy-based GANs (EBGANs) [24] can be thought of as the generative approach to the total variation distance. This connection is stated and proven in depth in Appendix D. At the core of the connection is that the discriminator will play the role of f maximizing equation (4) while its only restriction is being between 0 and m for some constant m . This will yield the same behaviour as being restricted to be between -1 and 1 up to a constant scaling factor irrelevant to optimization. Thus, when the discriminator approaches optimality the cost for the generator will approximate the total variation distance $\delta(\mathbb{P}_r, \mathbb{P}_\theta)$.

Since the total variation distance displays the same regularity as the JS, it can be seen that EBGANs will suffer from the same problems of classical GANs regarding not being able to train the discriminator till optimality and thus limiting itself to very imperfect gradients.

- Maximum Mean Discrepancy (MMD) [5] is a specific case of integral probability metrics when $\mathcal{F} = \{f \in \mathcal{H} : \|f\|_\infty \leq 1\}$ for \mathcal{H} some Reproducing Kernel Hilbert Space (RKHS) associated with a given kernel $k : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$. As proved on [5] we know that MMD is a proper metric and not only a pseudometric when the kernel is universal. In the specific case where $\mathcal{H} = L^2(\mathcal{X}, m)$ for m the normalized Lebesgue measure on \mathcal{X} , we know that $\{f \in C_b(\mathcal{X}), \|f\|_\infty \leq 1\}$ will be contained in \mathcal{F} , and therefore $d_{\mathcal{F}}(\mathbb{P}_r, \mathbb{P}_\theta) \leq \delta(\mathbb{P}_r, \mathbb{P}_\theta)$ so the regularity of the MMD distance as a loss function will be at least as bad as the one of the total variation. Nevertheless this is a very extreme case, since we would need a very powerful kernel to approximate the whole L^2 . However, even Gaussian kernels are able to detect tiny noise patterns as recently evidenced by [19]. This points to the fact that especially with low bandwidth kernels, the distance might be close to a saturating regime similar as with total variation or the JS. This obviously doesn't need to be the case for every kernel, and figuring out how and which different MMDs are closer to Wasserstein or total variation distances is an interesting topic of research.

The great aspect of MMD is that via the kernel trick there is no need to train a separate network to maximize equation (4) for the ball of a RKHS. However, this has the disadvantage that evaluating the MMD distance has computational cost that grows quadratically with the amount of samples used to estimate the expectations in (4). This last point makes MMD have limited scalability, and is sometimes inapplicable to many real life applications because of it. There are estimates with linear computational cost for the MMD [5] which in a lot of cases makes MMD very useful, but they also have worse sample complexity.

- Generative Moment Matching Networks (GMMNs) [10, 3] are the generative counterpart of MMD. By backproping through the kernelized formula for equation (4), they directly optimize $d_{MMD}(\mathbb{P}_r, \mathbb{P}_\theta)$ (the IPM when \mathcal{F} is as in the previous item). As mentioned, this has the advantage of not requiring a separate network to approximately maximize equation (4). However, GMMNs have enjoyed limited applicability. Partial explanations for their unsuccess are the quadratic cost as a function of the number of samples and vanishing gradients for low-bandwidth kernels. Furthermore, it may be possible that some kernels used in practice are unsuitable for capturing very complex distances in high dimensional sample spaces such as natural images. This is properly justified by the fact that [18] shows that for the typical Gaussian MMD test to be reliable (as in it's power as a statistical test approaching 1), we need the number of samples to grow linearly with the number of dimensions. Since the MMD computational cost grows quadratically with the number of samples in the batch used to estimate equation (4), this makes the cost of having a

reliable estimator grow quadratically with the number of dimensions, which makes it very inapplicable for high dimensional problems. Indeed, for something as standard as 64x64 images, we would need minibatches of size at least 4096 (without taking into account the constants in the bounds of [18] which would make this number substantially larger) and a total cost per iteration of 4096^2 , over 5 orders of magnitude more than a GAN iteration when using the standard batch size of 64.

That being said, these numbers can be a bit unfair to the MMD, in the sense that we are comparing empirical sample complexity of GANs with the theoretical sample complexity of MMDs, which tends to be worse. However, in the original GMMN paper [10] they indeed used a minibatch of size 1000, much larger than the standard 32 or 64 (even when this incurred in quadratic computational cost). While estimates that have linear computational cost as a function of the number of samples exist [5], they have worse sample complexity, and to the best of our knowledge they haven't been yet applied in a generative context such as in GMMNs.

6 Conclusion

We introduced an algorithm that we deemed WGAN, an alternative to traditional GAN training. In this new model, we showed that we can improve the stability of learning, get rid of problems like mode collapse, and provide meaningful learning curves useful for debugging and hyperparameter searches. Furthermore, we showed that the corresponding optimization problem is sound, and provided extensive theoretical work highlighting the deep connections to other distances between distributions.

References

- [1] Martin Arjovsky and Léon Bottou. Towards principled methods for training generative adversarial networks. In *International Conference on Learning Representations*, 2017. Under review.
- [2] Patrick Billingsley. *Convergence of probability measures*. Wiley Series in Probability and Statistics: Probability and Statistics. John Wiley & Sons Inc., New York, second edition, 1999. A Wiley-Interscience Publication.
- [3] Gintare Karolina Dziugaite, Daniel M. Roy, and Zoubin Ghahramani. Training generative neural networks via maximum mean discrepancy optimization. *CoRR*, abs/1505.03906, 2015.
- [4] Ian J. 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 27*, pages 2672–2680. Curran Associates, Inc., 2014.

- [5] Arthur Gretton, Karsten M. Borgwardt, Malte J. Rasch, Bernhard Schölkopf, and Alexander Smola. A kernel two-sample test. *J. Mach. Learn. Res.*, 13:723–773, 2012.
- [6] Ferenc Huszar. How (not) to train your generative model: Scheduled sampling, likelihood, adversary? *CoRR*, abs/1511.05101, 2015.
- [7] Shizuo Kakutani. Concrete representation of abstract (m)-spaces (a characterization of the space of continuous functions). *Annals of Mathematics*, 42(4):994–1024, 1941.
- [8] Diederik P. Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *CoRR*, abs/1412.6980, 2014.
- [9] Diederik P. Kingma and Max Welling. Auto-encoding variational bayes. *CoRR*, abs/1312.6114, 2013.
- [10] Yujia Li, Kevin Swersky, and Rich Zemel. Generative moment matching networks. In *Proceedings of the 32nd International Conference on Machine Learning (ICML-15)*, pages 1718–1727. JMLR Workshop and Conference Proceedings, 2015.
- [11] Luke Metz, Ben Poole, David Pfau, and Jascha Sohl-Dickstein. Unrolled generative adversarial networks. *Corr*, abs/1611.02163, 2016.
- [12] Paul Milgrom and Ilya Segal. Envelope theorems for arbitrary choice sets. *Econometrica*, 70(2):583–601, 2002.
- [13] Volodymyr Mnih, Adrià Puigdomènech Badia, Mehdi Mirza, Alex Graves, Timothy P. Lillicrap, Tim Harley, David Silver, and Koray Kavukcuoglu. Asynchronous methods for deep reinforcement learning. In *Proceedings of the 33rd International Conference on Machine Learning, ICML 2016, New York City, NY, USA, June 19-24, 2016*, pages 1928–1937, 2016.
- [14] Alfred Mller. Integral probability metrics and their generating classes of functions. *Advances in Applied Probability*, 29(2):429–443, 1997.
- [15] Radford M. Neal. Annealed importance sampling. *Statistics and Computing*, 11(2):125–139, April 2001.
- [16] Sebastian Nowozin, Botond Cseke, and Ryota Tomioka. f-gan: Training generative neural samplers using variational divergence minimization. pages 271–279, 2016.
- [17] Alec Radford, Luke Metz, and Soumith Chintala. Unsupervised representation learning with deep convolutional generative adversarial networks. *CoRR*, abs/1511.06434, 2015.

- [18] Aaditya Ramdas, Sashank J. Reddi, Barnabas Poczos, Aarti Singh, and Larry Wasserman. On the high-dimensional power of linear-time kernel two-sample testing under mean-difference alternatives. *Corr*, abs/1411.6314, 2014.
- [19] Dougal J Sutherland, Hsiao-Yu Tung, Heiko Strathmann, Soumyajit De, Aaditya Ramdas, Alex Smola, and Arthur Gretton. Generative models and model criticism via optimized maximum mean discrepancy. In *International Conference on Learning Representations*, 2017. Under review.
- [20] T. Tieleman and G. Hinton. Lecture 6.5—RmsProp: Divide the gradient by a running average of its recent magnitude. COURSERA: Neural Networks for Machine Learning, 2012.
- [21] Cédric Villani. *Optimal Transport: Old and New*. Grundlehren der mathematischen Wissenschaften. Springer, Berlin, 2009.
- [22] Yuhuai Wu, Yuri Burda, Ruslan Salakhutdinov, and Roger B. Grosse. On the quantitative analysis of decoder-based generative models. *CoRR*, abs/1611.04273, 2016.
- [23] Fisher Yu, Yinda Zhang, Shuran Song, Ari Seff, and Jianxiong Xiao. LSUN: Construction of a large-scale image dataset using deep learning with humans in the loop. *Corr*, abs/1506.03365, 2015.
- [24] Junbo Zhao, Michael Mathieu, and Yann LeCun. Energy-based generative adversarial network. *Corr*, abs/1609.03126, 2016.

A Why Wasserstein is indeed weak

We now introduce our notation. Let $\mathcal{X} \subseteq \mathbb{R}^d$ be a compact set (such as $[0, 1]^d$ the space of images). We define $\text{Prob}(\mathcal{X})$ to be the space of probability measures over \mathcal{X} . We note

$$C_b(\mathcal{X}) = \{f : \mathcal{X} \rightarrow \mathbb{R}, f \text{ is continuous and bounded}\}$$

Note that if $f \in C_b(\mathcal{X})$, we can define $\|f\|_\infty = \max_{x \in \mathcal{X}} |f(x)|$, since f is bounded. With this norm, the space $(C_b(\mathcal{X}), \|\cdot\|_\infty)$ is a normed vector space. As for any normed vector space, we can define its dual

$$C_b(\mathcal{X})^* = \{\phi : C_b(\mathcal{X}) \rightarrow \mathbb{R}, \phi \text{ is linear and continuous}\}$$

and give it the dual norm $\|\phi\| = \sup_{f \in C_b(\mathcal{X}), \|f\|_\infty \leq 1} |\phi(f)|$.

With this definitions, $(C_b(\mathcal{X})^*, \|\cdot\|)$ is another normed space. Now let μ be a signed measure over \mathcal{X} , and let us define the total variation distance

$$\|\mu\|_{TV} = \sup_{A \subseteq \mathcal{X}} |\mu(A)|$$

where the supremum is taken all Borel sets in \mathcal{X} . Since the total variation is a norm, then if we have \mathbb{P}_r and \mathbb{P}_θ two probability distributions over \mathcal{X} ,

$$\delta(\mathbb{P}_r, \mathbb{P}_\theta) := \|\mathbb{P}_r - \mathbb{P}_\theta\|$$

is a distance in $\text{Prob}(\mathcal{X})$ (called the total variation distance).

We can consider

$$\Phi : (\text{Prob}(\mathcal{X}), \delta) \rightarrow (C_b(\mathcal{X})^*, \|\cdot\|)$$

where $\Phi(\mathbb{P})(f) := \mathbb{E}_{x \sim \mathbb{P}}[f(x)]$ is a linear function over $C_b(\mathcal{X})$. The Riesz Representation theorem ([7], Theorem 10) tells us that Φ is an isometric immersion. This tells us that we can effectively consider $\text{Prob}(\mathcal{X})$ with the total variation distance as a subset of $C_b(\mathcal{X})^*$ with the norm distance. Thus, just to accentuate it one more time, the total variation over $\text{Prob}(\mathcal{X})$ is exactly the norm distance over $C_b(\mathcal{X})^*$.

Let us stop for a second and analyze what all this technicality meant. The main thing to carry is that we introduced a distance δ over probability distributions. When looked as a distance over a subset of $C_b(\mathcal{X})^*$, this distance gives the norm topology. The norm topology is very strong. Therefore, we can expect that not many functions $\theta \mapsto \mathbb{P}_\theta$ will be continuous when measuring distances between distributions with δ . As we will show later in Theorem 2, δ gives the same topology as the Jensen-Shannon divergence, pointing to the fact that the JS is a very strong distance, and is thus more propense to give a discontinuous loss function.

Now, all dual spaces (such as $C_b(\mathcal{X})^*$ and thus $\text{Prob}(\mathcal{X})$) have a strong topology (induced by the norm), and a weak* topology. As the name suggests, the weak* topology is much weaker than the strong topology. In the case of $\text{Prob}(\mathcal{X})$, the strong topology is given by the total variation distance, and the weak* topology is given by the Wasserstein distance (among others) [21].

B Assumption definitions

Assumption 1. Let $g : \mathcal{Z} \times \mathbb{R}^d \rightarrow \mathcal{X}$ be locally Lipschitz between finite dimensional vector spaces. We will denote $g_\theta(z)$ it's evaluation on coordinates (z, θ) . We say that g satisfies assumption 1 for a certain probability distribution p over \mathcal{Z} if there are local Lipschitz constants $L(\theta, z)$ such that

$$\mathbb{E}_{z \sim p}[L(\theta, z)] < +\infty$$

C Proofs of things

Proof of Theorem 1. Let θ and θ' be two parameter vectors in \mathbb{R}^d . Then, we will first attempt to bound $W(\mathbb{P}_\theta, \mathbb{P}_{\theta'})$, from where the theorem will come easily. The main element of the proof is the use of the coupling γ , the distribution of the joint $(g_\theta(Z), g_{\theta'}(Z))$, which clearly has $\gamma \in \Pi(\mathbb{P}_\theta, \mathbb{P}_{\theta'})$.

By the definition of the Wasserstein distance, we have

$$\begin{aligned} W(\mathbb{P}_\theta, \mathbb{P}_{\theta'}) &\leq \int_{\mathcal{X} \times \mathcal{X}} \|x - y\| d\gamma \\ &= \mathbb{E}_{(x,y) \sim \gamma} [\|x - y\|] \\ &= \mathbb{E}_z [\|g_\theta(z) - g_{\theta'}(z)\|] \end{aligned}$$

If g is continuous in θ , then $g_\theta(z) \rightarrow_{\theta \rightarrow \theta'} g_{\theta'}(z)$, so $\|g_\theta - g_{\theta'}\| \rightarrow 0$ pointwise as functions of z . Since \mathcal{X} is compact, the distance of any two elements in it has to be uniformly bounded by some constant M , and therefore $\|g_\theta(z) - g_{\theta'}(z)\| \leq M$ for all θ and z uniformly. By the bounded convergence theorem, we therefore have

$$W(\mathbb{P}_\theta, \mathbb{P}_{\theta'}) \leq \mathbb{E}_z [\|g_\theta(z) - g_{\theta'}(z)\|] \rightarrow_{\theta \rightarrow \theta'} 0$$

Finally, we have that

$$|W(\mathbb{P}_r, \mathbb{P}_\theta) - W(\mathbb{P}_r, \mathbb{P}_{\theta'})| \leq W(\mathbb{P}_\theta, \mathbb{P}_{\theta'}) \rightarrow_{\theta \rightarrow \theta'} 0$$

proving the continuity of $W(\mathbb{P}_r, \mathbb{P}_\theta)$.

Now let g be locally Lipschitz. Then, for a given pair (θ, z) there is a constant $L(\theta, z)$ and an open set U such that $(\theta, z) \in U$, such that for every $(\theta', z') \in U$ we have

$$\|g_\theta(z) - g_{\theta'}(z')\| \leq L(\theta, z)(\|\theta - \theta'\| + \|z - z'\|)$$

By taking expectations and $z' = z$ we

$$\mathbb{E}_z [\|g_\theta(z) - g_{\theta'}(z)\|] \leq \|\theta - \theta'\| \mathbb{E}_z [L(\theta, z)]$$

whenever $(\theta', z) \in U$. Therefore, we can define $U_\theta = \{\theta' | (\theta', z) \in U\}$. It's easy to see that since U was open, U_θ is as well. Furthermore, by assumption 1, we can define $L(\theta) = \mathbb{E}_z [L(\theta, z)]$ and achieve

$$|W(\mathbb{P}_r, \mathbb{P}_\theta) - W(\mathbb{P}_r, \mathbb{P}_{\theta'})| \leq W(\mathbb{P}_\theta, \mathbb{P}_{\theta'}) \leq L(\theta) \|\theta - \theta'\|$$

for all $\theta' \in U_\theta$, meaning that $W(\mathbb{P}_r, \mathbb{P}_\theta)$ is locally Lipschitz. This obviously implies that $W(\mathbb{P}_r, \mathbb{P}_\theta)$ is everywhere continuous, and by Radamacher's theorem we know it has to be differentiable almost everywhere.

The counterexample for item 3 of the Theorem is indeed Example 1. \square

Proof of Corollary 1. We begin with the case of smooth nonlinearities. Since g is C^1 as a function of (θ, z) then for any fixed (θ, z) we have $L(\theta, Z) \leq \|\nabla_{\theta, x} g_\theta(z)\| + \epsilon$ is an acceptable local Lipschitz constant for all $\epsilon > 0$. Therefore, it suffices to prove

$$\mathbb{E}_{z \sim p(z)}[\|\nabla_{\theta, z} g_\theta(z)\|] < +\infty$$

If H is the number of layers we know that $\nabla_z g_\theta(z) = \prod_{k=1}^H W_k D_k$ where W_k are the weight matrices and D_k are the diagonal Jacobians of the nonlinearities. Let $f_{i:j}$ be the application of layers i to j inclusively (e.g. $g_\theta = f_{1:H}$). Then, $\nabla_{W_k} g_\theta(z) = \left(\left(\prod_{i=k+1}^H W_i D_i \right) D_k \right) f_{1:k-1}(z)$. We recall that if L is the Lipschitz constant of the nonlinearity, then $\|D_i\| \leq L$ and $\|f_{1:k-1}(z)\| \leq \|z\| L^{k-1} \prod_{i=1}^{k-1} \|W_i\|$. Putting this together,

$$\begin{aligned} \|\nabla_{z, \theta} g_\theta(z)\| &\leq \left\| \prod_{i=1}^H W_i D_i \right\| + \sum_{k=1}^H \left\| \left(\prod_{i=k+1}^H W_i D_i \right) D_k \right\| \|f_{1:k-1}(z)\| \\ &\leq L^H \prod_{i=1}^H \|W_i\| + \sum_{k=1}^H \|z\| L^H \left(\prod_{i=1}^{k-1} \|W_i\| \right) \left(\prod_{i=k+1}^H \|W_i\| \right) \end{aligned}$$

If $C_1(\theta) = L^H \left(\prod_{i=1}^H \|W_i\| \right)$ and $C_2(\theta) = \sum_{k=1}^H L^H \left(\prod_{i=1}^{k-1} \|W_i\| \right) \left(\prod_{i=k+1}^H \|W_i\| \right)$ then

$$\mathbb{E}_{z \sim p(z)}[\|\nabla_{\theta, z} g_\theta(z)\|] \leq C_1(\theta) + C_2(\theta) \mathbb{E}_{z \sim p(z)}[\|z\|] < +\infty$$

finishing the proof \square

Proof of Theorem 2.

1. \bullet $(\delta(\mathbb{P}_n, \mathbb{P}) \rightarrow 0 \Rightarrow JS(\mathbb{P}_n, \mathbb{P}) \rightarrow 0)$ — Let \mathbb{P}_m be the mixture distribution $\mathbb{P}_m = \frac{1}{2}\mathbb{P}_n + \frac{1}{2}\mathbb{P}$ (note that \mathbb{P}_m depends on n). It is easily verified that $\delta(\mathbb{P}_m, \mathbb{P}_n) \leq \delta(\mathbb{P}_n, \mathbb{P})$, and in particular this tends to 0 (as does $\delta(\mathbb{P}_m, \mathbb{P})$). We now show this for completeness. Let μ be a signed measure, we define $\|\mu\|_{TV} = \sup_{A \subseteq \mathcal{X}} |\mu(A)|$. for all Borel sets A . In this case,

$$\begin{aligned} \delta(\mathbb{P}_m, \mathbb{P}_n) &= \|\mathbb{P}_m - \mathbb{P}_n\|_{TV} \\ &= \left\| \frac{1}{2}\mathbb{P} + \frac{1}{2}\mathbb{P}_n - \mathbb{P}_n \right\|_{TV} \\ &= \frac{1}{2} \|\mathbb{P} - \mathbb{P}_n\|_{TV} \\ &= \frac{1}{2} \delta(\mathbb{P}_n, \mathbb{P}) \leq \delta(\mathbb{P}_n, \mathbb{P}) \end{aligned}$$

Let $f_n = \frac{d\mathbb{P}_n}{d\mathbb{P}_m}$ be the Radon-Nykodim derivative between \mathbb{P}_n and the mixture. Note that by construction for every Borel set A we have $\mathbb{P}_n(A) \leq 2\mathbb{P}_m(A)$. If $A = \{f_n > 3\}$ then we get

$$\mathbb{P}_n(A) = \int_A f_n d\mathbb{P}_m \geq 3\mathbb{P}_m(A)$$

which implies $\mathbb{P}_m(A) = 0$. This means that f_n is bounded by 3 \mathbb{P}_m (and therefore \mathbb{P}_n and \mathbb{P})-almost everywhere. We could have done this for any constant larger than 2 but for our purposes 3 will suffice.

Let $\epsilon > 0$ fixed, and $A_n = \{f_n > 1 + \epsilon\}$. Then,

$$\mathbb{P}_n(A_n) = \int_{A_n} f_n d\mathbb{P}_m \geq (1 + \epsilon)\mathbb{P}_m(A_n)$$

Therefore,

$$\begin{aligned} \epsilon\mathbb{P}_m(A_n) &\leq \mathbb{P}_n(A_n) - \mathbb{P}_m(A_n) \\ &\leq |\mathbb{P}_n(A_n) - \mathbb{P}_m(A_n)| \\ &\leq \delta(\mathbb{P}_n, \mathbb{P}_m) \\ &\leq \delta(\mathbb{P}_n, \mathbb{P}). \end{aligned}$$

Which implies $\mathbb{P}_m(A_n) \leq \frac{1}{\epsilon}\delta(\mathbb{P}_n, \mathbb{P})$. Furthermore,

$$\begin{aligned} \mathbb{P}_n(A_n) &\leq \mathbb{P}_m(A_n) + |\mathbb{P}_n(A_n) - \mathbb{P}_m(A_n)| \\ &\leq \frac{1}{\epsilon}\delta(\mathbb{P}_n, \mathbb{P}) + \delta(\mathbb{P}_n, \mathbb{P}_m) \\ &\leq \frac{1}{\epsilon}\delta(\mathbb{P}_n, \mathbb{P}) + \delta(\mathbb{P}_n, \mathbb{P}) \\ &\leq \left(\frac{1}{\epsilon} + 1\right)\delta(\mathbb{P}_n, \mathbb{P}) \end{aligned}$$

We now can see that

$$\begin{aligned} KL(\mathbb{P}_n \|\mathbb{P}_m) &= \int \log(f_n) d\mathbb{P}_n \\ &\leq \log(1 + \epsilon) + \int_{A_n} \log(f_n) d\mathbb{P}_n \\ &\leq \log(1 + \epsilon) + \log(3)\mathbb{P}_n(A_n) \\ &\leq \log(1 + \epsilon) + \log(3) \left(\frac{1}{\epsilon} + 1\right) \delta(\mathbb{P}_n, \mathbb{P}) \end{aligned}$$

Taking limsup we get $0 \leq \limsup KL(\mathbb{P}_n \|\mathbb{P}_m) \leq \log(1 + \epsilon)$ for all $\epsilon > 0$, which means $KL(\mathbb{P}_n \|\mathbb{P}_m) \rightarrow 0$.

In the same way, we can define $g_n = \frac{d\mathbb{P}}{d\mathbb{P}_m}$, and

$$2\mathbb{P}_m(\{g_n > 3\}) \geq \mathbb{P}(\{g_n > 3\}) \geq 3\mathbb{P}_m(\{g_n > 3\})$$

meaning that $\mathbb{P}_m(\{g_n > 3\}) = 0$ and therefore g_n is bounded by 3 almost everywhere for $\mathbb{P}_n, \mathbb{P}_m$ and \mathbb{P} . With the same calculation, $B_n = \{g_n > 1 + \epsilon\}$ and

$$\mathbb{P}(B_n) = \int_{B_n} g_n \, d\mathbb{P}_m \geq (1 + \epsilon)\mathbb{P}_m(B_n)$$

so $\mathbb{P}_m(B_n) \leq \frac{1}{\epsilon}\delta(\mathbb{P}, \mathbb{P}_m) \rightarrow 0$, and therefore $\mathbb{P}(B_n) \rightarrow 0$. We can now show

$$\begin{aligned} KL(\mathbb{P} \parallel \mathbb{P}_m) &= \int \log(g_n) \, d\mathbb{P} \\ &\leq \log(1 + \epsilon) + \int_{B_n} \log(g_n) \, d\mathbb{P} \\ &\leq \log(1 + \epsilon) + \log(3)\mathbb{P}(B_n) \end{aligned}$$

so we achieve $0 \leq \limsup KL(\mathbb{P} \parallel \mathbb{P}_m) \leq \log(1 + \epsilon)$ and then $KL(\mathbb{P} \parallel \mathbb{P}_m) \rightarrow 0$. Finally, we conclude

$$JS(\mathbb{P}_n, \mathbb{P}) = \frac{1}{2}KL(\mathbb{P}_n \parallel \mathbb{P}_m) + \frac{1}{2}KL(\mathbb{P} \parallel \mathbb{P}_m) \rightarrow 0$$

- ($JS(\mathbb{P}_n, \mathbb{P}) \rightarrow 0 \Rightarrow \delta(\mathbb{P}_n, \mathbb{P}) \rightarrow 0$) — by a simple application of the triangular and Pinsker's inequalities we get

$$\begin{aligned} \delta(\mathbb{P}_n, \mathbb{P}) &\leq \delta(\mathbb{P}_n, \mathbb{P}_m) + \delta(\mathbb{P}, \mathbb{P}_m) \\ &\leq \sqrt{\frac{1}{2}KL(\mathbb{P}_n \parallel \mathbb{P}_m)} + \sqrt{\frac{1}{2}KL(\mathbb{P} \parallel \mathbb{P}_m)} \\ &\leq 2\sqrt{JS(\mathbb{P}_n, \mathbb{P})} \rightarrow 0 \end{aligned}$$

2. This is a long known fact that W metrizes the weak* topology of $(C(\mathcal{X}), \|\cdot\|_\infty)$ on $\text{Prob}(\mathcal{X})$, and by definition this is the topology of convergence in distribution. A proof of this can be found (for example) in [21].
3. This is a straightforward application of Pinsker's inequality

$$\begin{aligned} \delta(\mathbb{P}_n, \mathbb{P}) &\leq \sqrt{\frac{1}{2}KL(\mathbb{P}_n \parallel \mathbb{P})} \rightarrow 0 \\ \delta(\mathbb{P}, \mathbb{P}_n) &\leq \sqrt{\frac{1}{2}KL(\mathbb{P} \parallel \mathbb{P}_n)} \rightarrow 0 \end{aligned}$$

4. This is trivial by recalling the fact that δ and W give the strong and weak* topologies on the dual of $(C(\mathcal{X}), \|\cdot\|_\infty)$ when restricted to $\text{Prob}(\mathcal{X})$.

□

Proof of Theorem 3. Let us define

$$\begin{aligned} V(\tilde{f}, \theta) &= \mathbb{E}_{x \sim \mathbb{P}_r}[\tilde{f}(x)] - \mathbb{E}_{x \sim \mathbb{P}_\theta}[\tilde{f}(x)] \\ &= \mathbb{E}_{x \sim \mathbb{P}_r}[\tilde{f}(x)] - \mathbb{E}_{z \sim p(z)}[\tilde{f}(g_\theta(z))] \end{aligned}$$

where \tilde{f} lies in $\mathcal{F} = \{\tilde{f} : \mathcal{X} \rightarrow \mathbb{R}, \tilde{f} \in C_b(\mathcal{X}), \|\tilde{f}\|_L \leq 1\}$ and $\theta \in \mathbb{R}^d$.

Since \mathcal{X} is compact, we know by the Kantorovich-Rubenstein duality [21] that there is an $f \in \mathcal{F}$ that attains the value

$$W(\mathbb{P}_r, \mathbb{P}_\theta) = \sup_{\tilde{f} \in \mathcal{F}} V(\tilde{f}, \theta) = V(f, \theta)$$

Let us define $X^*(\theta) = \{f \in \mathcal{F} : V(f, \theta) = W(\mathbb{P}_r, \mathbb{P}_\theta)\}$. By the above point we know then that $X^*(\theta)$ is non-empty. We know that by a simple envelope theorem ([12], Theorem 1) that

$$\nabla_\theta W(\mathbb{P}_r, \mathbb{P}_\theta) = \nabla_\theta V(f, \theta)$$

for any $f \in X^*(\theta)$ when both terms are well-defined.

Let $f \in X^*(\theta)$, which we know exists since $X^*(\theta)$ is non-empty for all θ . Then, we get

$$\begin{aligned} \nabla_\theta W(\mathbb{P}_r, \mathbb{P}_\theta) &= \nabla_\theta V(f, \theta) \\ &= \nabla_\theta [\mathbb{E}_{x \sim \mathbb{P}_r}[f(x)] - \mathbb{E}_{z \sim p(z)}[f(g_\theta(z))]] \\ &= -\nabla_\theta \mathbb{E}_{z \sim p(z)}[f(g_\theta(z))] \end{aligned}$$

under the condition that the first and last terms are well-defined. The rest of the proof will be dedicated to show that

$$-\nabla_\theta \mathbb{E}_{z \sim p(z)}[f(g_\theta(z))] = -\mathbb{E}_{z \sim p(z)}[\nabla_\theta f(g_\theta(z))] \quad (5)$$

when the right hand side is defined. For the reader who is not interested in such technicalities, he or she can skip the rest of the proof.

Since $f \in \mathcal{F}$, we know that it is 1-Lipschitz. Furthermore, $g_\theta(z)$ is locally Lipschitz as a function of (θ, z) . Therefore, $f(g_\theta(z))$ is locally Lipschitz on (θ, z) with constants $L(\theta, z)$ (the same ones as g). By Radamacher's Theorem, $f(g_\theta(z))$ has to be differentiable almost everywhere for (θ, z) jointly. Rewriting this, the set $A = \{(\theta, z) : f \circ g \text{ is not differentiable}\}$ has measure 0. By Fubini's Theorem, this implies that for almost every θ the section $A_\theta = \{z : (\theta, z) \in A\}$ has measure 0. Let's now fix a θ_0 such that the measure of A_{θ_0} is null (**such as when the right hand side of equation (5) is well defined**). For this θ_0 we have $\nabla_\theta f(g_\theta(z))|_{\theta_0}$ is well-defined for almost any z , and since $p(z)$ has a density, it is defined $p(z)$ -a.e. By assumption 1 we know that

$$\mathbb{E}_{z \sim p(z)}[\|\nabla_\theta f(g_\theta(z))|_{\theta_0}\|] \leq \mathbb{E}_{z \sim p(z)}[L(\theta_0, z)] < +\infty$$

so $\mathbb{E}_{z \sim p(z)}[\nabla_\theta f(g_\theta(z))|_{\theta_0}]$ is well-defined for almost every θ_0 . Now, we can see

$$\frac{\mathbb{E}_{z \sim p(z)}[f(g_\theta(z))] - \mathbb{E}_{z \sim p(z)}[f(g_{\theta_0}(z))] - \langle (\theta - \theta_0), \mathbb{E}_{z \sim p(z)}[\nabla_\theta f(g_\theta(z))|_{\theta_0}] \rangle}{\|\theta - \theta_0\|} \quad (6)$$

$$= \mathbb{E}_{z \sim p(z)} \left[\frac{f(g_\theta(z)) - f(g_{\theta_0}(z)) - \langle (\theta - \theta_0), \nabla_\theta f(g_\theta(z))|_{\theta_0} \rangle}{\|\theta - \theta_0\|} \right]$$

By differentiability, the term inside the integral converges $p(z)$ -a.e. to 0 as $\theta \rightarrow \theta_0$. Furthermore,

$$\begin{aligned} & \left\| \frac{f(g_\theta(z)) - f(g_{\theta_0}(z)) - \langle (\theta - \theta_0), \nabla_\theta f(g_\theta(z))|_{\theta_0} \rangle}{\|\theta - \theta_0\|} \right\| \\ & \leq \frac{\|\theta - \theta_0\| L(\theta_0, z) + \|\theta - \theta_0\| \|\nabla_\theta f(g_\theta(z))|_{\theta_0}\|}{\|\theta - \theta_0\|} \\ & \leq 2L(\theta_0, z) \end{aligned}$$

and since $\mathbb{E}_{z \sim p(z)}[2L(\theta_0, z)] < +\infty$ by assumption 1, we get by dominated convergence that Equation 6 converges to 0 as $\theta \rightarrow \theta_0$ so

$$\nabla_\theta \mathbb{E}_{z \sim p(z)}[f(g_\theta(z))] = \mathbb{E}_{z \sim p(z)}[\nabla_\theta f(g_\theta(z))]$$

for almost every θ , and in particular when the right hand side is well defined. Note that the mere existence of the left hand side (meaning the differentiability a.e. of $\mathbb{E}_{z \sim p(z)}[f(g_\theta(z))]$) had to be proven, which we just did. \square

D Energy-based GANs optimize total variation

In this appendix we show that under an optimal discriminator, energy-based GANs (EBGANs) [24] optimize the total variation distance between the real and generated distributions.

Energy-based GANs are trained in a similar fashion to GANs, only under a different loss function. They have a discriminator D who tries to maximize, and a generator network g_θ that's trained to minimize

$$L(D, g_\theta) = \mathbb{E}_{x \sim \mathbb{P}_r}[D(x)] + \mathbb{E}_{z \sim p(z)}[[m - D(g_\theta(z))]^+]$$

for some $m > 0$ and $[x]^+ = \max(0, x)$. Very importantly, D is constrained to be non-negative, since otherwise the trivial solution for D would be to set everything to arbitrarily low values.

We say that a measurable function $D^* : \mathcal{X} \rightarrow [0, +\infty)$ is optimal for g_θ (or \mathbb{P}_θ) if $L(D^*, g_\theta) \geq L(D, g_\theta)$ for all other measurable functions D . We show that such a discriminator always exists for any two distributions \mathbb{P}_r and \mathbb{P}_θ , and that under such a discriminator, $L(D^*, g_\theta)$ is proportional to $\delta(\mathbb{P}_r, \mathbb{P}_\theta)$. As a simple corollary, we get the fact that $L(D^*, g_\theta)$ attains its minimum value if and only if $\delta(\mathbb{P}_r, \mathbb{P}_\theta)$ is at its minimum value, which is 0, and $\mathbb{P}_r = \mathbb{P}_\theta$ (Theorems 1-2 of [24]).

Theorem 4. *Let \mathbb{P}_r be a the real data distribution over a compact space \mathcal{X} . Let $g_\theta : \mathcal{Z} \rightarrow \mathcal{X}$ be a measurable function (such as any neural network). Then, an optimal discriminator D^* exists for \mathbb{P}_r and \mathbb{P}_θ , and*

$$L(D^*, g_\theta) = m + \frac{m}{2} \delta(\mathbb{P}_r, \mathbb{P}_\theta)$$

Proof. First, we prove that there exists an optimal discriminator. Let $D : \mathcal{X} \rightarrow [0, +\infty)$ be a measurable function, then $D'(x) = \min(D(x), m)$ is also a measurable function, and $L(D', g_\theta) \geq L(D, g_\theta)$. Therefore, a function $D^* : \mathcal{X} \rightarrow [0, +\infty)$ is optimal if and only if $D^{*'}$ is. Furthermore, it is optimal if and only if $L(D^*, g_\theta) \geq L(D, g_\theta)$ for all $D : \mathcal{X} \rightarrow [0, m]$. We are then interested to see if there's an optimal discriminator for the problem $\max_{0 \leq D(x) \leq m} L(D, g_\theta)$.

Note now that if $0 \leq D(x) \leq m$ we have

$$\begin{aligned} L(D^*, g_\theta) &= \mathbb{E}_{x \sim \mathbb{P}_r}[D(x)] + \mathbb{E}_{z \sim p(z)}[[m - D(g_\theta(z))]^+] \\ &= \mathbb{E}_{x \sim \mathbb{P}_r}[D(x)] + \mathbb{E}_{z \sim p(z)}[m - D(g_\theta(z))] \\ &= m + \mathbb{E}_{x \sim \mathbb{P}_r}[D(x)] - \mathbb{E}_{z \sim p(z)}[D(g_\theta(z))] \\ &= m + \mathbb{E}_{x \sim \mathbb{P}_r}[D(x)] - \mathbb{E}_{x \sim \mathbb{P}_\theta}[D(x)] \end{aligned}$$

Therefore, we know that

$$\begin{aligned} \sup_{0 \leq D(x) \leq m} L(D, g_\theta) &= m + \sup_{0 \leq D(x) \leq m} \mathbb{E}_{x \sim \mathbb{P}_r}[D(x)] - \mathbb{E}_{x \sim \mathbb{P}_\theta}[D(x)] \\ &= m + \sup_{-\frac{m}{2} \leq D(x) \leq \frac{m}{2}} \mathbb{E}_{x \sim \mathbb{P}_r}[D(x)] - \mathbb{E}_{x \sim \mathbb{P}_\theta}[D(x)] \\ &= m + \frac{m}{2} \sup_{-1 \leq D(x) \leq 1} \mathbb{E}_{x \sim \mathbb{P}_r}[D(x)] - \mathbb{E}_{x \sim \mathbb{P}_\theta}[D(x)] \end{aligned}$$

The interesting part is that

$$\sup_{-1 \leq f(x) \leq 1} \mathbb{E}_{x \sim \mathbb{P}_r}[f(x)] - \mathbb{E}_{x \sim \mathbb{P}_\theta}[f(x)] = \delta(\mathbb{P}_r, \mathbb{P}_\theta) \quad (7)$$

And furthermore, there is an $f^* : \mathcal{X} \rightarrow [-1, 1]$ such that $\delta(\mathbb{P}_r, \mathbb{P}_\theta) = \mathbb{E}_{x \sim \mathbb{P}_r}[f^*(x)] - \mathbb{E}_{x \sim \mathbb{P}_\theta}[f^*(x)]$. Note that the existence of said f^* then implies that

$$L(D, g_\theta) \leq m + \frac{m}{2} (\mathbb{E}_{x \sim \mathbb{P}_r}[f^*(x)] - \mathbb{E}_{x \sim \mathbb{P}_\theta}[f^*(x)]) = m + \frac{m}{2} \delta(\mathbb{P}_r, \mathbb{P}_\theta)$$

Equation (7) and the existence of said f^* are wide known facts seen for example in [14, 2] but nevertheless we show the proof for completeness. Take $\mu = \mathbb{P}_r - \mathbb{P}_\theta$, which is a signed measure, and (P, Q) its Hahn decomposition. Then, we can define $f^* := \mathbb{1}_P - \mathbb{1}_Q$. By construction, then

$$\mathbb{E}_{x \sim \mathbb{P}_r}[f^*(x)] - \mathbb{E}_{x \sim \mathbb{P}_\theta}[f^*(x)] = \mu(P) - \mu(Q) = \|\mu\|_{TV} = \|\mathbb{P}_r - \mathbb{P}_\theta\|_{TV} = \delta(\mathbb{P}_r, \mathbb{P}_\theta)$$

Furthermore, if f is bounded between -1 and 1, we get

$$\begin{aligned} \mathbb{E}_{x \sim \mathbb{P}_r}[f(x)] - \mathbb{E}_{x \sim \mathbb{P}_\theta}[f(x)] &= \int f \, d\mathbb{P}_r - \int f \, d\mathbb{P}_\theta \\ &= \int f \, d\mu \\ &\leq \int |f| \, d|\mu| \\ &= |\mu|(\mathcal{X}) = \|\mu\|_{TV} = \delta(\mathbb{P}_r, \mathbb{P}_\theta) \end{aligned}$$

Finally, defining $D^* = \frac{m}{2} + \frac{m}{2} f^*$ we end up with

$$\begin{aligned} L(D^*, g_\theta) &= \mathbb{E}_{x \sim \mathbb{P}_r}[D^*(x)] + \mathbb{E}_{x \sim \mathbb{P}_\theta}[m - D^*(x)] \\ &= m + \frac{m}{2} (\mathbb{E}_{x \sim \mathbb{P}_r}[f^*(x)] - \mathbb{E}_{x \sim \mathbb{P}_\theta}[f^*(x)]) \\ &= m + \frac{m}{2} \delta(\mathbb{P}_r, \mathbb{P}_\theta) \\ &\geq L(D, g_\theta) \end{aligned}$$

finishing the proof □

E Generator's cost during normal GAN training

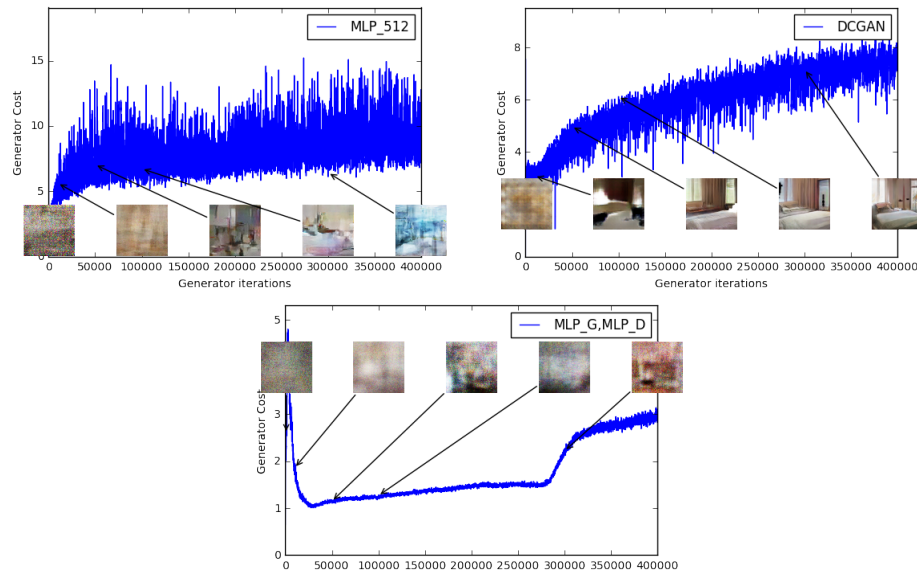


Figure 8: Cost of the generator during normal GAN training, for an MLP generator (upper left) and a DCGAN generator (upper right). Both had a DCGAN discriminator. **Both curves have increasing error.** Samples get better for the DCGAN but the cost of the generator increases, pointing towards no significant correlation between sample quality and loss. Bottom: MLP with both generator and discriminator. The curve goes up and down regardless of sample quality. All training curves were passed through the same median filter as in Figure 3.

F Sheets of samples

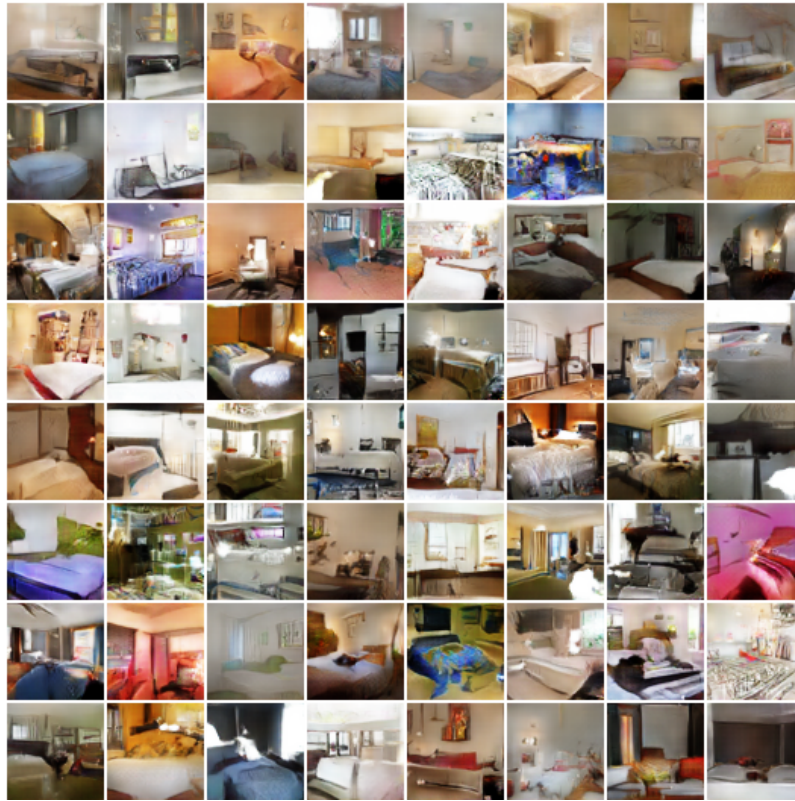


Figure 9: WGAN algorithm: generator and critic are DCGANs.

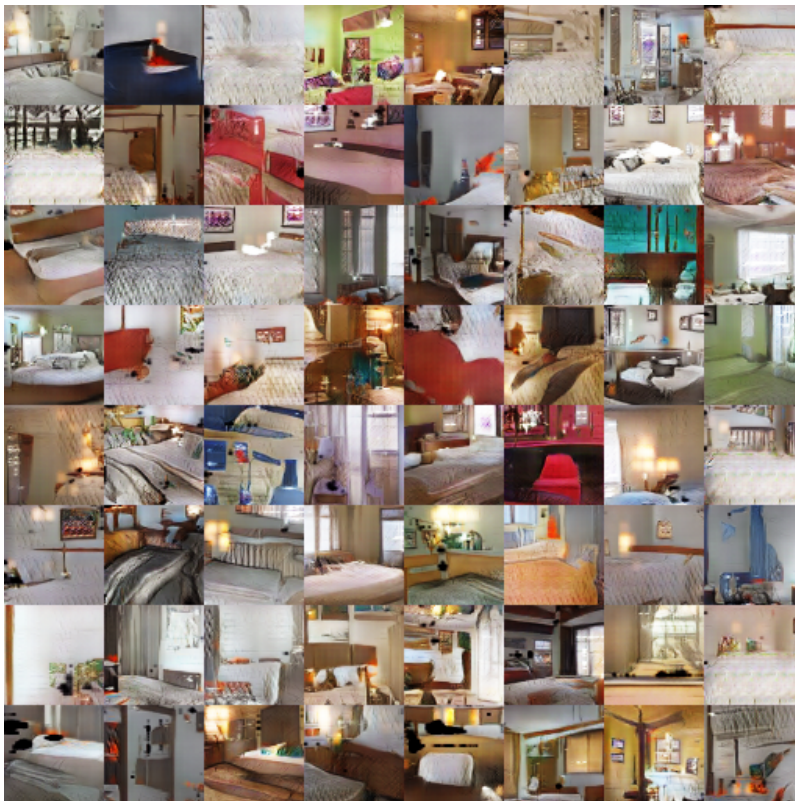


Figure 10: Standard GAN procedure: generator and discriminator are DCGANs.

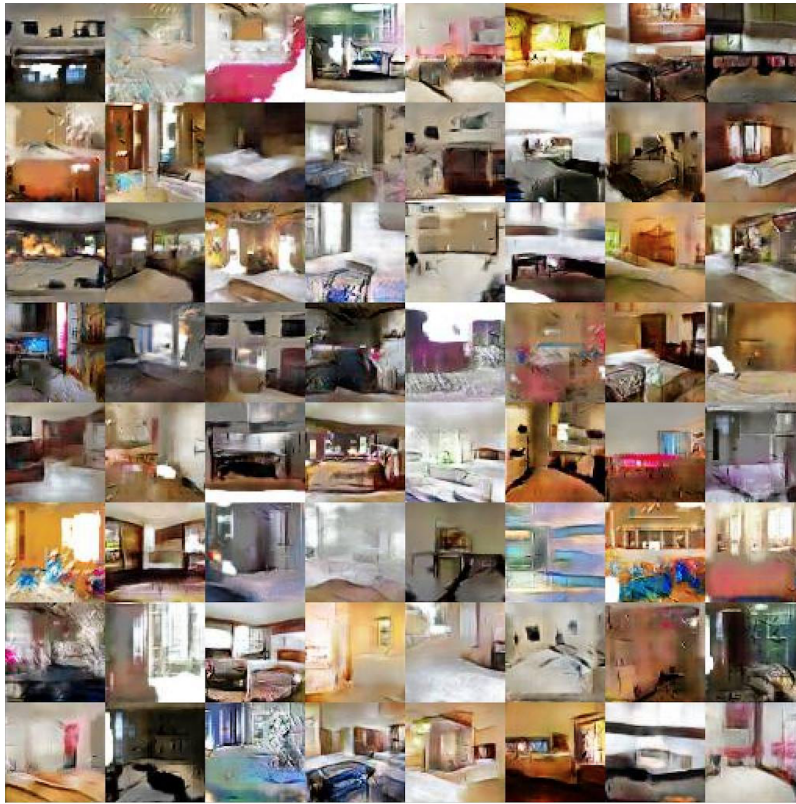


Figure 11: WGAN algorithm: generator is a DCGAN without batchnorm and constant filter size. Critic is a DCGAN.

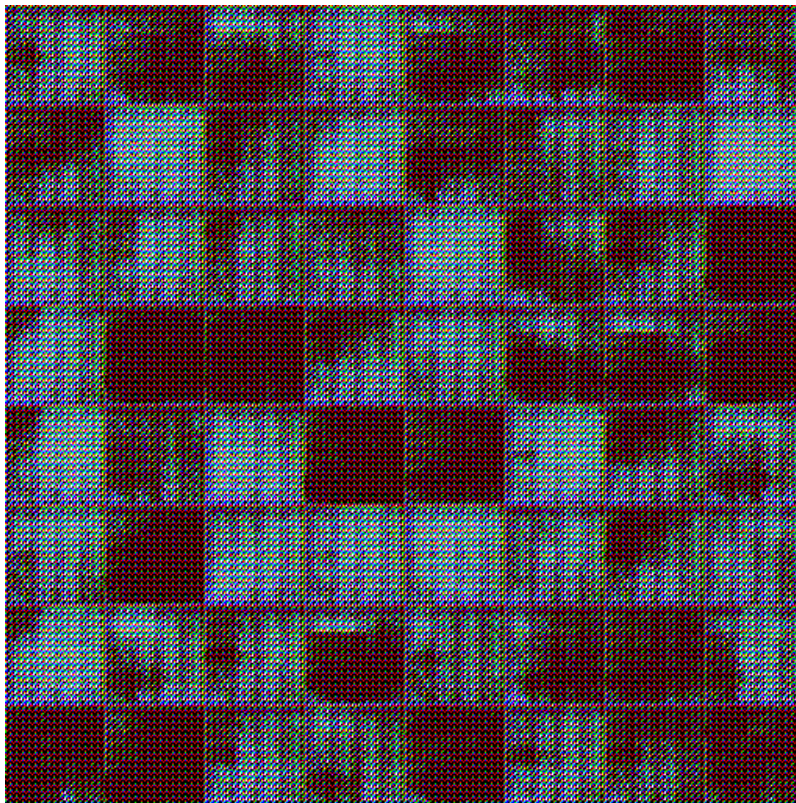


Figure 12: Standard GAN procedure: generator is a DCGAN without batchnorm and constant filter size. Discriminator is a DCGAN.

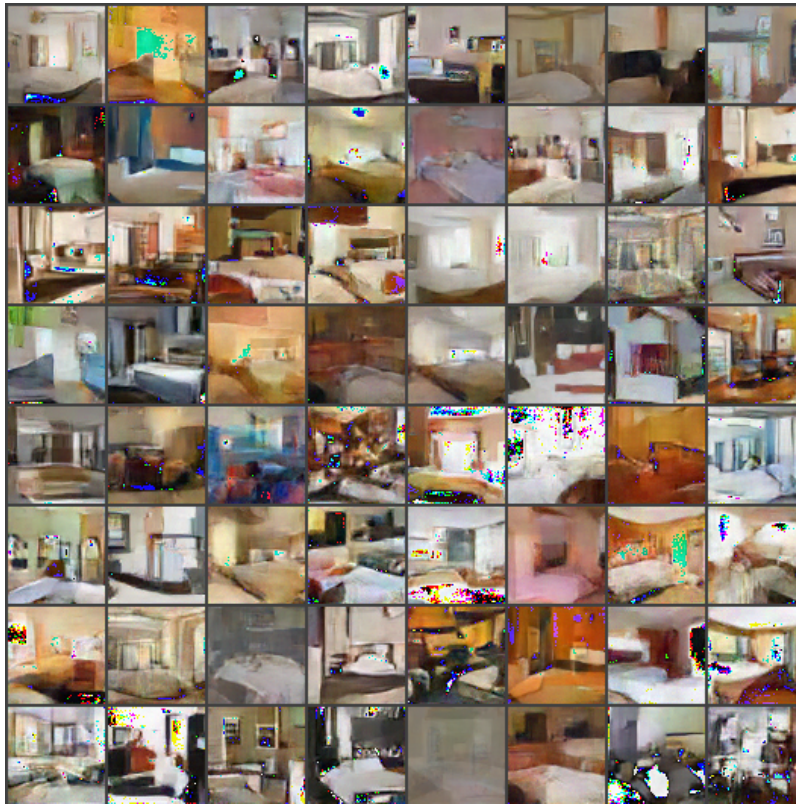


Figure 13: WGAN algorithm: generator is an MLP with 4 hidden layers of 512 units, critic is a DCGAN.



Figure 14: Standard GAN procedure: generator is an MLP with 4 hidden layers of 512 units, discriminator is a DCGAN.