

TRAINING GENERATIVE NETWORKS USING RANDOM DISCRIMINATORS

**Babak Barazandeh* †*Meisam Razaviyayn* ‡*Maziar Sanjabi*

*†‡University of Southern California
{barazand, razaviya, sanjabi}@usc.edu

ABSTRACT

In recent years, Generative Adversarial Networks (GANs) have drawn a lot of attentions for learning the underlying distribution of data in various applications. Despite their wide applicability, training GANs is notoriously difficult. This difficulty is due to the min-max nature of the resulting optimization problem and the lack of proper tools of solving general (non-convex, non-concave) min-max optimization problems. In this paper, we try to alleviate this problem by proposing a new generative network that relies on the use of random discriminators instead of adversarial design. This design helps us to avoid the min-max formulation and leads to an optimization problem that is stable and could be solved efficiently. The performance of the proposed method is evaluated using handwritten digits (MNIST) and Fashion products (Fashion-MNIST) data sets. While the resulting images are not as sharp as adversarial training, the use of random discriminator leads to a much faster algorithm as compared to the adversarial counterpart. This observation, at the minimum, illustrates the potential of the random discriminator approach for warm-start in training GANs.

Index Terms— Generative Adversarial Networks, Deep Neural Network, Randomized Learning, Non-convex Min-Max Optimization

1. INTRODUCTION

Generative Adversarial Networks (GANs) [1] have been relatively successful in learning underlying distribution of data, especially in application such as image generation. GANs aims to find the mapping that matches a known distribution to the underlying distribution of the data. The way they perform this task is by projecting the inputs to a higher dimension using Neural Networks [2] and then minimizing the distance between the mapped distribution and the unknown distribution in the projected space. To find the optimal network, [1] proposed using Jensen-Shannon divergence [3] for measuring the distance between projected distribution and the data distribution. Later on, [4] generalized the idea by using the f-divergence as the measure. [5] and [6] proposed using least square and absolute deviation as the measure.

The most recent works proposed using Wasserstein distance and Maximum Mean Discrepancy (MMD) as the distance measure [7–9]. Unlike Jensen-Shannon divergence, the recent measures are continuous and almost everywhere differentiable. The common thread between all these approaches is that the problem is usually formulated as a game between two agents, i.e. generator and discriminator. Generator’s role is to generate samples as close as possible to real data and discriminator is responsible for distinguishing between real data and the generated samples. The result is a non-convex min-max game which is difficult to solve. The learning process, which should solve the resulting non-convex min-max game, is hard to tackle, due to many factors such as using discontinuous [7] or non-smooth [2] measure. In addition to these factors, the fact that all of these models try to learn the mapping transformation adversarially makes the training unstable. Adding regularization or starting from a good initial point is one approach to overcome these problems [2]. However, for most problems finding a good initial point might be as hard as solving the problem itself.

Randomization has shown promising improvement in machine learning algorithms [10, 11]. As the result, to prevent over-mentioned issues, we propose learning underlying distribution of data not through adversarial player but through a random projection. This random projection not only decreases the computation time by removing the optimization steps needed for most of the discriminator’s role, but also leads to a more stable optimization problem. The proposed method has the state of the art performance for simple datasets such as MNIST and Fashion-MNIST.

2. PROBLEM FORMULATION

Let $x \in \mathbb{R}^d$ be a random variable with distribution P_x representing the real data; and z be a random variable representing a known distribution such as standard Gaussian. Our goal is to find a function or a neural network $G(\cdot)$ such that $G(z)$ has a similar distribution to the real data distribution P_x . Therefore, our objective is to solve the following optimization problem

$$\min_G \text{dist}(P_{G(z)}, P_x), \quad (1)$$

where $P_{G(z)}$ is the distribution of $G(z)$ and $\text{dist}(\cdot, \cdot)$ is a distance measure between the two distributions.

A natural question to ask is about what distance metric to use. The original paper of Goodfellow [1] suggests the use of Jensen–Shannon divergence. However, as mentioned in [7], this divergence is not continuous. Therefore, [2, 7] suggest to use the optimal transport distance. In what follows, we first review this distance and then discuss our methodology for solving (1).

3. OPTIMAL TRANSPORT DISTANCE

Let p and q be two discrete distributions taking m different values/states. Thus the distributions p and q can be represented by m -dimensional vectors (p_1, \dots, p_m) and (q_1, \dots, q_m) . The optimal transport distance is defined as the minimum amount of work needs to be done for transporting distribution p to q (and vice versa). Let $\pi_{i,j}$ be the amount of mass moved from state i to state j ; and c_{ij} represent the per-unit cost of this move. Then the optimal transport distance between the two distributions p and q is defined as [12]:

$$\begin{aligned} \text{dist}(p, q) = \min_{\pi \geq 0} & \sum_{i=1}^m \sum_{j=1}^m c_{ij} \pi_{ij} \\ \text{s.t.} & \sum_{j=1}^m \pi_{ij} = p_i, \quad \forall i = 1, \dots, m \\ & \sum_{i=1}^m \pi_{ij} = q_j, \quad \forall j = 1, \dots, m, \end{aligned} \quad (2)$$

where the constrains guarantee that the mapping π is a valid transport. In practice, a popular approach is to solve the dual problem. It is not hard to see that the dual of the optimization problem (2) can be written as

$$\begin{aligned} \text{dist}(p, q) = \max_{\lambda, \gamma} & \sum_{i=1}^m \gamma_i p_i + \sum_{j=1}^m \lambda_j q_j \\ \text{s.t.} & \lambda_j + \gamma_i \leq c_{ij}, \quad \forall i, j = 1, \dots, m. \end{aligned} \quad (3)$$

When c is a proper distance, this dual variable should satisfy $\lambda = -\gamma$ [12]. In practice, since the dimension m is large and estimating p and q accurately is not possible, we parameterize the dual variable with a neural network and solve the dual optimization problem by training two neural networks simultaneously [7]. However, this approach leads to a non-convex min-max optimization problem. Unlike special cases such as convex-concave set-up [13], there is no algorithm to date in the literature which can find even an ϵ -stationary point in the general non-convex setting; see [14] and the references therein. Therefore, training generative adversarial networks (GANs) can become notoriously difficult in practice and may require significant tuning of training parameters. A natural solution is to not parameterize the dual function and instead solve (2) or (3) directly which leads to a convex reformulation. However, as mentioned earlier, since the dimension m is large, approximating p and q is statistically not possible. Moreover, the distance in the original feature domain may not reflect the

actual distance between the distributions. Thus, we suggest an alternative formulation in the next section.

4. TRAINING IN DIFFERENT FEATURE DOMAIN

In many applications, the closeness of samples in the original feature domain does not reflect the actual similarity between the samples. For example, two images of the same object may have a large difference when the distance is computed in the pixel domain. Therefore, other mappings of the features, such as features obtained by Convolutional Neural Network (CNN) may be used to extract meaningful features from samples [15].

Let $\mathcal{D} = \{D_1, D_2, \dots, D_K\}$ be a collection of meaningful features we are interested in. In other words, each function $D \in \mathcal{D}$ is a mapping from our original feature domain to the domain of interest, i.e., $D_k(\cdot) : \mathbb{R}^d \mapsto \mathbb{R}^{d'}, \forall k = 1, \dots, K$. Then, instead of solving (1), one might be interested in solving the following optimization problem

$$\min_G \sum_{k=1}^K w_k \text{dist}(P_{D_k(G(z))}, P_{D_k(x)}), \quad (4)$$

where $P_{D_k(G(z))}$ represents the distribution of the random variable $D_k(G(z))$; $P_{D_k(x)}$ is the distribution of $D_k(x)$; and w_k is a weight coefficient indicating the importance of the k -th feature D_k .

In the general setting, we may have uncountable number of mappings D_k . Thus, by defining a measure on the set \mathcal{D} , we can generalize (4) to the following optimization problem

$$\min_G \mathbb{E}_D \left[\text{dist}(P_{D(G(z))}, P_{D(x)}) \right]. \quad (5)$$

Remark 1. We use the notation D since the function D plays the role of a discriminator in the Generative Adversarial Learning (GANs) context.

Plugging (2) in the equation (5) leads to the optimization problem

$$\min_G \mathbb{E}_D \left[\begin{array}{l} \max_{(\lambda, \gamma) \in \mathcal{C}} \sum_{i=1}^m \gamma_i P_{D(G(z))}^i + \sum_{j=1}^m \lambda_j P_{D(x)}^j \\ \text{s.t.} \quad \lambda_j + \gamma_i \leq c_{ij}, \quad \forall i, j. \end{array} \right], \quad (6)$$

where $\mathcal{C} = \{(\lambda, \gamma) \mid \lambda_i + \gamma_j \leq c_{ij}, \forall i, j\}$.

Unfortunately, in practice, we do not have access to the actual values of the distributions $P_{D(x)}$ and $P_{D(G(z))}$. However, we can estimate them using a batch of generated and real samples. The following simple lemma motivates the use of a natural surrogate function.

Lemma 1. Let p and q be two discrete distributions with $p = (p_1, \dots, p_m)$ and $q = (q_1, \dots, q_m)$. Let $x \in \mathbb{R}^m$

and $y \in \mathbb{R}^m$ be the corresponding one-hot encoded random variables, i.e., $P(x = e_i) = p_i, \forall i = 1, \dots, m$ and $P(y = e_i) = q_i, \forall i = 1, \dots, m$, where e_i is the i -th standard basis. Assume further that $\text{dist}(p, q)$ is the optimal transport distance between p and q defined in (2). Let \hat{p}^n and \hat{q}^n be the natural unbiased estimator of p and q based on n i.i.d. samples. In other words, $\hat{p}^n = \frac{1}{n} \sum_{\ell=1}^n x_\ell$ and $\hat{q}^n = \frac{1}{n} \sum_{\ell=1}^n y_\ell$, where x_ℓ and $y_\ell, \ell = 1, \dots, n$, are i.i.d samples obtained from distributions p and q , respectively. Then,

$$\mathbb{E} [\text{dist}(\hat{p}^{n+1}, \hat{q}^{n+1})] \leq \mathbb{E} [\text{dist}(\hat{p}^n, \hat{q}^n)].$$

Moreover,

$$\lim_{n \rightarrow \infty} \text{dist}(\hat{p}^n, \hat{q}^n) = \text{dist}(p, q), \text{ almost surely.}$$

Proof. The proof is similar to the standard proof in sample average approximation method; see [16, Proposition 5.6]. Notice that,

$$\begin{aligned} & \mathbb{E} [\text{dist}(\hat{p}^{n+1}, \hat{q}^{n+1})] \\ &= \mathbb{E} \left[\max_{(\lambda, \gamma) \in \mathcal{C}} \sum_{i=1}^m \gamma_i \hat{p}_i^{n+1} + \sum_{j=1}^m \lambda_j \hat{q}_j^{n+1} \right] \\ &= \mathbb{E} \left[\max_{(\lambda, \gamma) \in \mathcal{C}} \langle \hat{p}^{n+1}, \gamma \rangle + \langle \hat{q}^{n+1}, \lambda \rangle \right] \\ &= \frac{1}{n+1} \mathbb{E} \left[\max_{(\lambda, \gamma) \in \mathcal{C}} \langle \gamma, \sum_{\ell} x_\ell \rangle + \langle \lambda, \sum_{\ell} y_\ell \rangle \right] \\ &= \frac{1}{n(n+1)} \mathbb{E} \left[\max_{(\lambda, \gamma) \in \mathcal{C}} \sum_{t=1}^{n+1} \langle \gamma, \sum_{\ell \neq t} x_\ell \rangle + \sum_{t=1}^{n+1} \langle \lambda, \sum_{\ell \neq t} y_\ell \rangle \right] \\ &\leq \frac{1}{n(n+1)} \mathbb{E} \left[\sum_{t=1}^{n+1} \max_{(\lambda, \gamma) \in \mathcal{C}} \langle \gamma, \sum_{\ell \neq t} x_\ell \rangle + \langle \lambda, \sum_{\ell \neq t} y_\ell \rangle \right] \\ &= \frac{1}{(n+1)} \sum_{t=1}^{n+1} \mathbb{E} \left[\frac{1}{n} \max_{(\lambda, \gamma) \in \mathcal{C}} \langle \gamma, \sum_{\ell \neq t} x_\ell \rangle + \langle \lambda, \sum_{\ell \neq t} y_\ell \rangle \right] \\ &= \mathbb{E} [\text{dist}(\hat{p}^n, \hat{q}^n)]. \end{aligned}$$

The proof of the almost sure convergence follows directly from the facts that $\lim_{n \rightarrow \infty} \hat{p}^n = p$, $\lim_{n \rightarrow \infty} \hat{q}^n = q$, and the continuity of the distance function. \square

The above lemma suggests a natural upper-bound for the objective function in (6). More precisely, instead of solving (6), we can solve

$$\min_G \mathbb{E} \left[\begin{array}{l} \max_{(\lambda, \gamma) \in \mathcal{C}} \sum_{i=1}^m \gamma_i \hat{P}_{D(G(z))}^i + \sum_{j=1}^m \lambda_j \hat{P}_{D(x)}^j \\ \text{s.t.} \quad \lambda_j + \gamma_i \leq c_{ij}, \quad \forall i, j \end{array} \right], \quad (7)$$

where $\hat{P}_{D(G(z))}$ and $\hat{P}_{D(x)}$ are the unbiased estimators of $P_{D(G(z))}$ and $P_{D(x)}$ based on our i.i.d samples. Moreover,

the expectation is taken with respect to both, the function D as well as the batch of samples which is drawn for estimating the distributions. As we will see later, in practice it is easier to use the primal form for solving the inner problem in (7), i.e.,

$$\min_G \mathbb{E} \left[\begin{array}{l} \min_{\pi \geq 0} \sum_{i=1}^m \sum_{j=1}^m c_{ij} \pi_{ij} \\ \text{s.t.} \quad \sum_{j=1}^m \pi_{ij} = \hat{P}_{D(G(z))}^i, \sum_{i=1}^m \pi_{ij} = \hat{P}_{D(x)}^j, \quad \forall i, j \end{array} \right],$$

To show the dependence of c_{ij} to G , let us assume that our generator G is generating the output $h(w, z)$ from the input z . Here w represents the weights of the network needed to be learned. Moreover, in practice, the value of $\hat{P}_{D(G(z))}^i$ is estimated by taking the average over all batch of data. Hence, by duplicating variables if necessary, we can re-write the above optimization problem as

$$\min_w \mathbb{E}_{z, x, D} \left[\begin{array}{l} \min_{\pi \geq 0} \sum_{i=1}^n \sum_{j=1}^n \|D(h(w, z_i)) - D(x_j)\| \pi_{ij} \\ \text{s.t.} \quad \pi \vec{1} = \frac{1}{n}, \quad \pi^T \vec{1} = \frac{1}{n} \end{array} \right]. \quad (8)$$

Here, n is the batch size and we ignored the entries of $\hat{P}_{D(G(z))}$ and $\hat{P}_{D(x)}$ that are zero. Notice that to obtain an algorithm with convergence guarantee for solving this optimization problem, one can properly regularize the inner optimization problem to obtain unbiased estimates of the gradient of the objective function [2, 14]. However, in this work, due to practical considerations, we suggest to *approximately* solve the inner problem and use the approximate solution for solving (8).

Solving the inner-problem approximately. In order to solve the inner problem in (8), we need to solve

$$\begin{aligned} & \min_{\pi \geq 0} \sum_{i=1}^n \sum_{j=1}^n \|D(h(w, z_i)) - D(x_j)\| \pi_{ij} \\ & \text{s.t.} \quad \pi \vec{1} = \frac{1}{n}, \quad \pi^T \vec{1} = \frac{1}{n}. \end{aligned} \quad (9)$$

Notice that this problem is the classical *optimal assignment problem* which can be solved using Hungarian method [17], Auction algorithm [18], or many other methods proposed in the literature. Based on our observations, even the greedy method of assigning each column to the lowest unassigned row worked in our numerical experiments. The benefit of the greedy method is that it can be performed almost linearly in m by the use of a proper hash function.

Algorithm 1 summarizes our proposed Generative Networks using Random Discriminator (GN-RD) algorithm for solving (8).

Remark 2. The training approach in Algorithm 1 relies on two neural networks: the generative and the discriminator. Hence, Algorithm 1 can be viewed as a GANs training approach where

Algorithm 1: Generative Networks using Random Discriminator (GN-RD)

Input : w_0 : Initialization for generator’s parameter, α : Learning rate, n : Batch size, N_{Itr} : Maximum iteration number

```
1 for  $t = 1 : N_{\max}$  do
2   Sample an i.i.d. batch of real data  $(x_1, \dots, x_n)$ 
3   Sample an i.i.d. batch of noise  $(z_1, \dots, z_n)$ 
4   Create a random discriminator neural network  $D$ 
   with random weights
5   Solve (9) by finding the optimal assignment value
   between real data and generated sample
6   Update generator’s parameter,
    $w_{t+1} = w_t - \alpha \nabla_w G(w_t)$ 
```

```
7 end
Output :  $G(w_{N_{\max}})$ 
```

we use a random discriminator at each iteration of updating the generator.

Remark 3. *The recent works [19, 20] have similarities in terms of learning generative models through min-min formulation instead of min-max formulation. However, unlike their method, 1) our algorithm is based on mapping images via randomly generated discriminators; 2) In our analysis, we establish that this formulation leads to an upper-bound of the distance measure; 3) our algorithm is based on the use of optimal assignment, while the works [19, 20] suggests a greedy matching, which is more difficult to understand and analyze.*

5. NUMERICAL EXPERIMENTS

In this section, we evaluate the performance of the proposed GN-RD algorithm for learning generative networks to create samples from MNIST [21] and Fashion-MNIST [22] datasets. As mentioned previously, the proposed algorithm does not require any optimization on the discriminator network and only needs randomly generated discriminator to learn the underlying distribution of the data¹.

5.1. Learning handwritten digits and fashion products

In this section, we use GN-RD for generating samples from handwritten digits and Fashion-MNIST datasets. Each of these datasets contains 50K training samples.

Architecture of the Neural Networks: The generator’s Neural Network consists of two fully connected layer with 1024 and 6272 neurons. The output of the second fully connected layer is followed by two deconvolutional layers to generate the final 28×28 image.

¹All the experiments have been run on a machine with single GeForce GTX 1050 Ti GPU.

The discriminator neural network has two convolution layers each followed by a max pool. The size of the both convolutional layers are 64. The last layer has been flattened to create the output. The design of both neural networks is summarized below:

- Generator: [FC(100, 1024), Leaky ReLU($\alpha = 0.2$), FC(1024, 6272), Leaky ReLU($\alpha = 0.2$), DECONV(64, kernel size = 4, stride = 2), Leaky ReLU(alpha = 0.2), DECONV(1, kernel size = 4, stride = 2), Sigmoid].
- Discriminator: [CONV(64, filter size = 5, stride = 1), Leaky ReLU(alpha = 0.2), Max Pool (kernel size = 2, stride = 2), COVN(64, filter size = 5, stride = 1), Max Pool (kernel size = 2, stride = 2), Flating].

We have used originally proposed adversarial discriminator for Wasserstein GAN (WGAN) [7], Wasserstein GAN with gradient penalty (WGAN-GP) [8]² and Cramér GAN [23]³.

As mentioned in Algorithm 1, it is important to notice that unlike benchmark methods, the proposed method only optimizes the generator’s parameters. However, at each iteration, weights in the convolutional layers of the discriminator are randomly generated from normal distribution.

Hyper parameters: We have used Adam with step size 0.001 and $\beta_1 = 0.5$ and $\beta_2 = 0.9$ as the optimizer for our generator. The batch size is set to 100.

Fig.1 shows the result of the generated digits and the corresponding inception score [24] using different benchmark methods. As seen from the figure, the proposed GN-RD is able to quickly learn the underlying distribution of the data and generate promising samples.

Fig. 2 shows the result of using the proposed method for generating samples from fashion MNIST dataset. The sample is generated only after 600 iterations (~ 10 minutes) of the proposed method which shows that the GN-RD quickly converges and generates promising samples.

6. CONCLUSION

Generative Adversarial Networks (GANs) have been able to learn the underlying distribution of the data and generate samples from it. Training GANs is notoriously unstable due to their non-convex min-max formulation. In this work, we propose the use of randomized discriminator to avoid facing the complexity of solving non-convex min-max problems. Evaluating the performance of the proposed method on real data set of MNIST and Fashion-MNIST shows the ability of the proposed method in generating promising samples without adversarial learning.

²For WGAN and WGAN-GP implementation visit https://github.com/igul222/improved_wgan_training

³For Cramér GAN implementation visit <https://github.com/jiamings/cramer-gan>

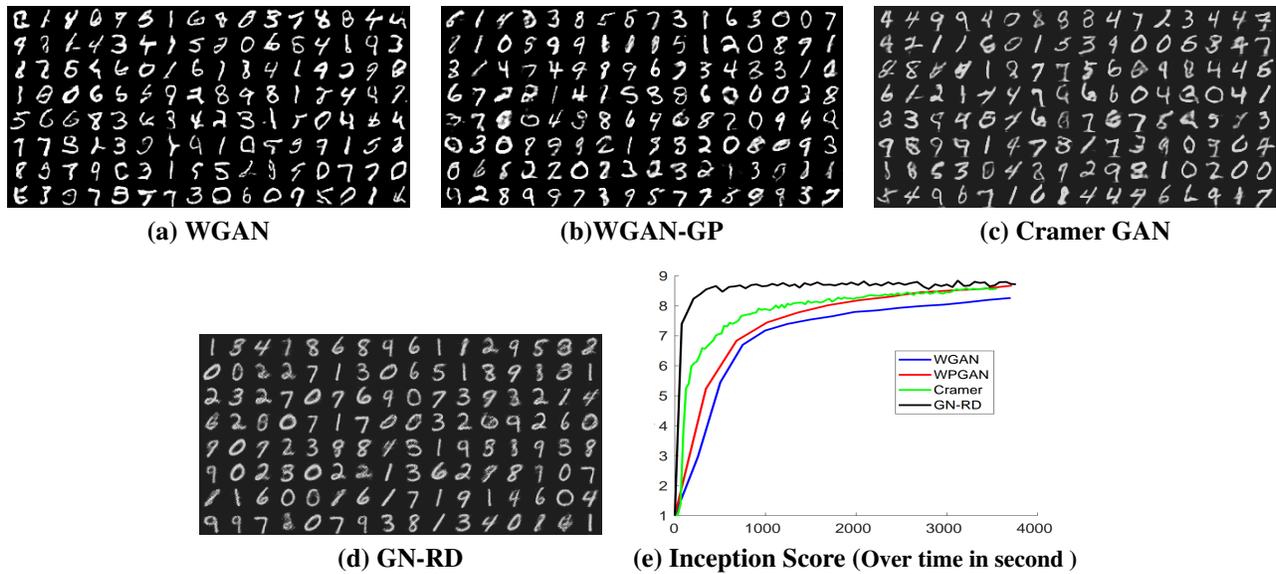


Fig. 1. Generating hand-written digits using MNIST dataset

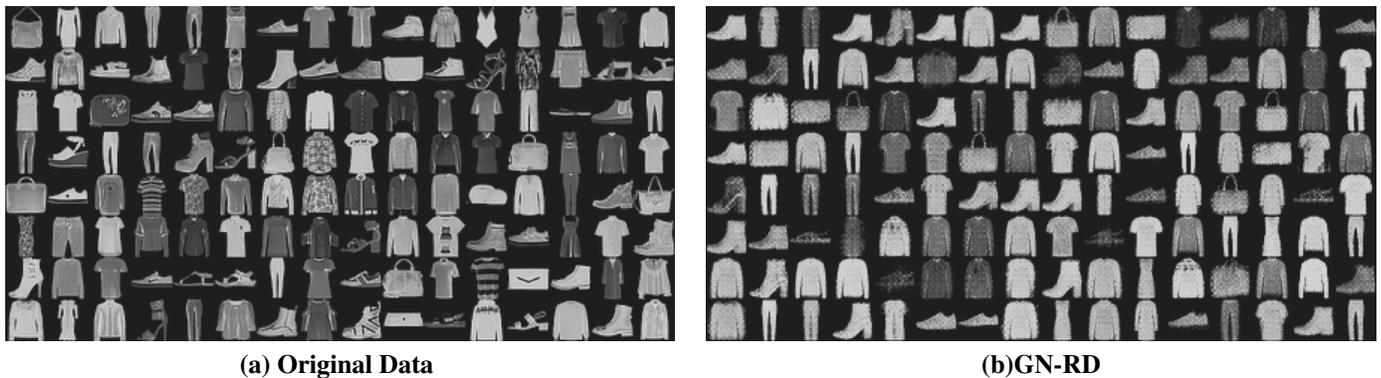


Fig. 2. Generating fashion products using Fashion-MNIST dataset

Acknowledgement

The authors would like to thank Mohammad Norouzi for his insightful feedback.

7. REFERENCES

- [1] I. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. Courville, and Y. Bengio, "Generative adversarial nets," in *Advances in neural information processing systems*, 2014, pp. 2672–2680.
- [2] M. Sanjabi, J. Ba, M. Razaviyayn, and J. D. Lee, "On the convergence and robustness of training gans with regularized optimal transport," in *Advances in Neural Information Processing Systems*, 2018, pp. 7091–7101.
- [3] J. Lin, "Divergence measures based on the shannon entropy," *IEEE Transactions on Information theory*, vol. 37, no. 1, pp. 145–151, 1991.
- [4] S. Nowozin, B. Cseke, and R. Tomioka, "f-gan: Training generative neural samplers using variational divergence minimization," in *Advances in neural information processing systems*, 2016, pp. 271–279.
- [5] X. Mao, Q. Li, H. Xie, R. Y. Lau, Z. Wang, and S. Paul Smolley, "Least squares generative adversarial networks," in *Proceedings of the IEEE International Conference on Computer Vision*, 2017, pp. 2794–2802.
- [6] J. Zhao, M. Mathieu, and Y. LeCun, "Energy-based generative adversarial network," *arXiv preprint arXiv:1609.03126*, 2016.

- [7] M. Arjovsky, S. Chintala, and L. Bottou, "Wasserstein gan," *arXiv preprint arXiv:1701.07875*, 2017.
- [8] I. Gulrajani, F. Ahmed, M. Arjovsky, V. Dumoulin, and A. C. Courville, "Improved training of wasserstein gans," in *Advances in Neural Information Processing Systems*, 2017, pp. 5767–5777.
- [9] M. Bińkowski, D. J. Sutherland, M. Arbel, and A. Gretton, "Demystifying mmd gans," *arXiv preprint arXiv:1801.01401*, 2018.
- [10] B. Barazandeh and M. Razaviyayn, "On the behavior of the expectation-maximization algorithm for mixture models," in *2018 IEEE Global Conference on Signal and Information Processing (GlobalSIP)*. IEEE, 2018, pp. 61–65.
- [11] Y. Sun, A. Gilbert, and A. Tewari, "Random relu features: Universality, approximation, and composition," *arXiv preprint arXiv:1810.04374*, 2018.
- [12] C. Villani, "Optimal transport—old and new, volume 338 of a series of comprehensive studies in mathematics," 2009.
- [13] A. Juditsky and A. Nemirovski, "Solving variational inequalities with monotone operators on domains given by linear minimization oracles," *Mathematical Programming*, vol. 156, no. 1-2, pp. 221–256, 2016.
- [14] M. Nouiehed, M. Sanjabi, J. D. Lee, and M. Razaviyayn, "Solving a class of non-convex min-max games using iterative first order methods," *arXiv preprint arXiv:1902.08297*, 2019.
- [15] K. O'Shea and R. Nash, "An introduction to convolutional neural networks," *arXiv preprint arXiv:1511.08458*, 2015.
- [16] A. Shapiro, D. Dentcheva, and A. Ruszczyński, *Lectures on stochastic programming: modeling and theory*. SIAM, 2009.
- [17] H. W. Kuhn, "The hungarian method for the assignment problem," *Naval research logistics quarterly*, vol. 2, no. 1-2, pp. 83–97, 1955.
- [18] D. P. Bertsekas, "The auction algorithm: A distributed relaxation method for the assignment problem," *Annals of operations research*, vol. 14, no. 1, pp. 105–123, 1988.
- [19] K. Li and J. Malik, "On the implicit assumptions of gans," *arXiv preprint arXiv:1811.12402*, 2018.
- [20] K. Li and J. Malik, "Implicit maximum likelihood estimation," *arXiv preprint arXiv:1809.09087*, 2018.
- [21] Y. LeCun, C. Cortes, and C. Burges, "Mnist handwritten digit database, 1998," URL <http://www.research.att.com/~yann/ocr/mnist>, 1998.
- [22] H. Xiao, K. Rasul, and R. Vollgraf, "Fashion-mnist: a novel image dataset for benchmarking machine learning algorithms," *arXiv preprint arXiv:1708.07747*, 2017.
- [23] M. G. Bellemare, I. Danihelka, W. Dabney, S. Mohamed, B. Lakshminarayanan, S. Hoyer, and R. Munos, "The cramer distance as a solution to biased wasserstein gradients," *arXiv preprint arXiv:1705.10743*, 2017.
- [24] S. Barratt and R. Sharma, "A note on the inception score," *arXiv preprint arXiv:1801.01973*, 2018.