

# A deep introspection on Generative Adversarial Networks

Riccardo Lincetto, Guglielmo Camporese

Department of Information Engineering, University of Padova – Via Gradenigo, 6/b, 35131 Padova, Italy  
emails: riccardo.lincetto, guglielmo.camporese @studenti.unipd.it

**Abstract**—GANs, namely Generative Adversarial Networks, are a hot topic nowadays. These models have the ability of generating good quality data, learning their distributions on a training set. Image generation in particular benefited from this framework, thanks also to the ease of assessment of the results obtained. Here we review the analysis of GAN models, recalling the necessary results from game theory, and propose a parametrization of the problem, which is then assessed analytically and by means of simulation, on image generation with MNIST dataset.

## I. INTRODUCTION

Since its rise, deep learning had a great impact on discriminative models. Generative models instead were not affected by this innovation at first, but this trend changed with the introduction of Generative Adversarial Networks (GAN), a powerful framework first introduced in [1]. Since then, GAN gained more and more momentum because of the ability of training *deep generative models*, avoiding some of the difficulties encountered in other frameworks [2].

GAN is a sub-class of generative models that use, generally two, *neural networks*: given a training set of sample data, distributed according to a probability density function (pdf)  $p_{data}$ , the purpose of GAN is to generate samples according to a distribution  $p_g$ , that mimics  $p_{data}$ , without explicitly defining it. As suggested by the name, this is achieved by putting in competition two entities: a generator (G) and a discriminator (D). The task of G is to generate data that can be regarded as true by D, while D has the purpose of correctly distinguishing real from fake data. The classical real-life analogy with this process involves counterfeiters trying to produce fake currency and the police trying to detect it. A graphical representation of the process is depicted in fig.1. This kind of interaction between the two entities can naturally be modeled with a game theoretical approach, where each player has its own strategies and payoffs, but in this paper we will rather talk about costs, as will be discussed in II. This framework anyway has a major drawback: computing a Nash Equilibrium (NE) requires assuming that the networks have enough capacity, i.e. it has to be done in the non-parametric limit for the networks [1]. This implies that the equilibrium point found in practice can be different from the theoretical one, because of intrinsic limits imposed by implementing a network with a finite set of parameters.

In this paper we review some of the literature and explain our need to go back to the origins of GAN, implementing our own version of the code and simulating different scenarios, where discriminator is passed different fake-to-true ratios of images.

The remainder of the paper is organized as follows: a brief overview of the literature is presented in II; a description of our work is then presented in III; the obtained results are presented in IV; finally we discuss our conclusions in V.

## II. RELATED WORK

Since their appearance in literature, GANs have been successfully applied to problems of image generation, editing and semi-supervised learning [5] [6]. Results obtained were so promising that the new framework captured the interest of many researchers, leading to a proliferation of various flavors of GAN, each claiming to have better performances on a specific domain. It's difficult anyway to understand how to compare different GAN models, because of the lack of a consistent metric and the different architectures networks can be designed with, which for each project are related to the corresponding computational budget. A tentative to define some guidelines to avoid these problems, together with a fair and comprehensive comparison of state-of-the-art GANs, is discussed in [3]: what emerges here is that the computational budget plays a major role, allowing bad algorithms to outperform good ones if given enough time; plus, despite the many claims of superiority, there's no empirical evidence that those algorithms are better across all datasets, in fact in most of the cases the original model outperforms the others. We thus report here a formalization of the original problem [1], modeling it with a game theory approach in a more precise way, as done in [4].

GANs exploit as players (named G and D) the learning capability of neural networks. The two anyway have different purposes, so in general they are designed with different architectures (e.g. as in fig.2 and fig.3). As suggested in [1], the generator's distribution  $p_g$  over data  $x$  is learnt defining a prior on input noise variables  $p_z(z)$  and representing the mapping to the data space as  $G(z; \theta_g)$ , where  $\theta_g$  stands for the weights of the network. Similarly for the discriminator, a mapping  $D(x; \theta_d)$  can be defined from the data space to a scalar, also referred to as  $D(x)$ , representing the probability that the input belongs to the true distribution  $p_{data}$ , i.e. to the training set. In the formalization of the game then, pure strategies are defined by the sets of possible  $\theta_g$  and  $\theta_d$ , while the utilities functions are the opposite of the loss functions that the networks have to minimize.

For classification tasks, cross-entropy loss function is universally accepted as the best choice, giving good results in terms of learning speed: this is what is usually selected for D

in GANs. The simplest design of GAN uses as loss function for G the opposite of D's, defining thus a zero-sum game (MM-GAN). Ideally, the number of strategies for each player would be infinite but, as pointed out in [4], when using floating point numbers it becomes finite, albeit very large: the game is then finite, implying the existence of a NE, at least in mixed strategies. It has to be kept into account also the fact that the optimization of a neural network can lead to a Local NE (LNE) because the problem is non-convex. We can then solve the equilibrium problem by computing the minimax of the payoffs. Loss functions are defined as:

$$\begin{aligned} L^{(D)} &= -\mathbb{E}_{x \sim p_{data}(x)}[\log D(x)] + \\ &\quad -\mathbb{E}_{z \sim p_z(z)}[\log(1 - D(G(z)))]; \\ L^{(G)} &= -L^{(D)}, \end{aligned}$$

where, the expectations are computed on mini-batches of data. For the very definition of cross-entropy loss function,  $L^{(G)}$  can be simplified because all the samples are generated according to the same distribution:

$$L^{(G)} = \mathbb{E}_{z \sim p_z(z)}[\log(1 - D(G(z)))].$$

The value of the game is defined as  $V(G, D) = -L^{(D)}$ . Fixing a generator  $\bar{G}$ , the optimal discriminator can be computed maximizing D's utility, which corresponds to  $V(\bar{G}, D)$ :

$$D^*(x) = \arg \max_D \{V(\bar{G}, D)\}.$$

Assuming that the model has infinite capacity in representing pdfs, we can study the convergence to a NE point in the space of pdf:

$$\begin{aligned} V(\bar{G}, D) &= \int_x p_{data}(x) \log(D(x)) dx + \\ &\quad + \int_z p_z(z) \log(1 - D(\bar{G}(z))) dz \\ &= \int_x \left[ p_{data}(x) \log(D(x)) + \right. \\ &\quad \left. + p_g(x) \log(1 - D(x)) \right] dx \\ &= \int_x L(x, D(x), \dot{D}(x)) dx. \end{aligned}$$

$D^*(x)$  must satisfy the Euler-Lagrange equation:

$$\frac{\partial L}{\partial D} = \frac{dL}{dx} \frac{\partial L}{\partial \dot{D}}$$

and since  $\partial L / \partial \dot{D} = 0$  we get:

$$\frac{\partial L}{\partial D} = \frac{p_{data}(x)}{D(x)} - \frac{p_g(x)}{1 - D(x)} = 0$$

that implies:

$$D^*(x) = \frac{p_{data}(x)}{p_{data}(x) + p_g(x)}.$$

Minimizing then  $V(G, D^*)$  over  $G$ , the result is:

$$p_g = p_{data},$$

from which  $D^*(x) = 1/2$  and  $V(G^*, D^*) = -\log(4)$  as proved in [1].

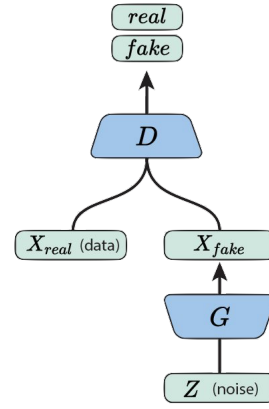


Figure 1: a block-diagram representing generative adversarial networks. A generator network ( $G$ ) creates samples in data space from noise  $z$ . A discriminator network ( $D$ ) then compares data, trying to distinguish true from fake samples.

Cross-entropy loss functions are particularly effective for discrimination tasks, but they're not for generation ones because, in the beginning of the training phase, when  $G$  is poor and  $D$  is able to correctly recognize generated samples,  $\log(1 - D(G(z)))$  saturates to zero, thus resulting in a poor gradient for backpropagation. The learning in that case is too slow, but this problem can be avoided changing the loss function for  $G$ : instead of minimizing  $\log(1 - D(G(z)))$ , we can maximize  $\log(D(G(z)))$ . This is formally defined as a Non-Saturating GAN (NS-GAN), where the losses to be minimized are:

$$\begin{aligned} L^{(D)} &= -\mathbb{E}_{x \sim p_{data}(x)}[\log D(x)] + \\ &\quad -\mathbb{E}_{z \sim p_z(z)}[\log(1 - D(G(z)))]; \\ L^{(G)} &= -\log(D(G(z))), \end{aligned}$$

This new game isn't zero-sum any more, but the same equilibrium point of the dynamics can be found as before.

### III. EXPERIMENTAL SETTING

For our project, we implemented a version of NS-GAN which can be found at [github.com/RicLincio/gt-project]. The training of the two neural networks, depicted in fig.2 and fig.3, occurs after mini-batches of data are passed to them: for the discriminator network  $D$ , in a mini-batch there are usually the same number of true and of fake images. We instead carried out the training with different configurations, where  $D$  is passed different true-to-fake data ratios: this is done in practice defining a parameter  $\alpha \in [0, 1]$  that represents the portion of true data with respect to the size of the mini-batch. Considering a zero-sum game, this results in a parametrisation of the loss functions as follows:

$$\begin{aligned} L^{(D)} &= -\alpha \cdot \mathbb{E}_{x \sim p_{data}(x)}[\log D(x)] + \\ &\quad -(1 - \alpha) \cdot \mathbb{E}_{z \sim p_z(z)}[\log(1 - D(G(z)))]; \\ L^{(G)} &= -L^{(D)}. \end{aligned}$$

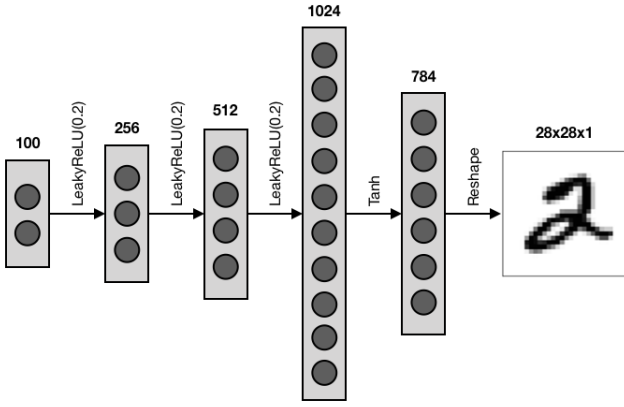


Figure 2: a high-level representation of G, the generator network, with 4 hidden layers. Each circle represents a perceptron unit in the network. A random vector of dimensionality 2, i.e. input noise, is transformed into a 28x28 grayscale image. The architecture presented is also used in our implementation of the model.

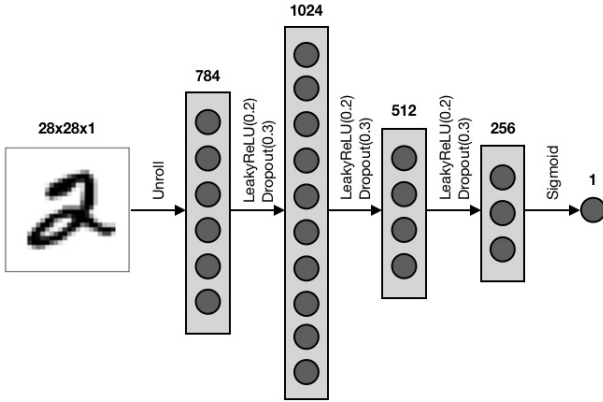


Figure 3: a high-level representation of D, the discriminator network, with 4 hidden layers. Each circle represents a perceptron unit in the network. 28x28 grayscale images are transformed into a scalar in range  $[0, 1]$ , representing the probability of the image being true. The architecture presented is also used in our implementation of the model.

Performing the same analysis as before, we found that the optimal discriminator should have the form:

$$D_{\alpha}^*(x) = \frac{\alpha \cdot p_{data}(x)}{\alpha \cdot p_{data}(x) + (1 - \alpha) \cdot p_g(x)}$$

and that the loss of the generator is minimized, again, when

$$p_g = p_{data},$$

from which the optimal D can be rewritten as

$$D_{\alpha}^*(x) = \alpha.$$

Supporting these results, it can also be noticed that when setting  $\alpha = 0.5$ , i.e. when there are the same amount of true and fake images in a mini-batch, results are consistent with what found

in II. The value of the game here is:

$$V(G^*, D_{\alpha}^*) = -\alpha \cdot \log(\alpha) - (1 - \alpha) \cdot \log(1 - \alpha).$$

As before the equilibrium point can be extended to the case of NS-GAN, though with different losses:

$$L^{(D)} = V(G^*, D_{\alpha}^*) = -\alpha \cdot \log(\alpha) - (1 - \alpha) \cdot \log(1 - \alpha),$$

$$L^{(G)} = -\log(\alpha).$$

To evaluate the effects of this parametrisation, we preliminarily trained the NS-GAN models on two different 2D artificial datasets: first on a "sigmoid-shaped" line and then on a circle, as can be seen in 4. The values of  $\alpha$  that we tested are  $\alpha = \{0.1, 0.3, 0.5, 0.7, 0.9\}$ . Once evaluated on the two artificial distribution, we tested our code also on the MNIST dataset on the same values of  $\alpha$ . All results are reported in IV.

#### IV. RESULTS

We first report here the results obtained on the two 2D datasets. In fig.4 points generated after 3000 epochs are displayed, while in fig.5 the losses for both D and G are plotted, as functions of the epochs, i.e. cycles over the entire dataset. As we can see in fig.5, the G loss converges to the theoretical optimum loss value, which is  $-\log(\alpha)$ , and the same happens for D, to the respective values.

It can be immediately noticed that for both data distributions, the highest loss for the discriminator is obtained with  $\alpha = 0.5$  which foretells the training of a good generator (notice that for this  $\alpha$ , as expected,  $L^{(D)}$  converges to  $\ln(2) = 0.693$ ). Anyway also with  $\alpha = \{0.3, 0.7\}$  good results are achieved in this sense. From the generator's losses it can be noticed that, the lower  $\alpha$ , the faster the loss convergence to the final value. An insight on how well are points generated is given in FIG.X, where for different values of  $\alpha$  we can see how points are generated after 3000 training epochs: these results confirm what expected, that is better performances from more balanced number of true images per mini-batch.

When applying the model to the MNIST dataset, similar results are obtained FIG.X, with losses for the discriminator approaching this time 0.5 for  $\alpha = \{0.3, 0.5, 0.7\}$ . Also in this data space, as pointed out in the figure, when mini-batches are more balanced the D loss is higher, indicating better generators. The evolution of how numbers are generated is depicted in FIG.X, where for each  $\alpha$ , every 20 and up to 100 epochs, a grid with 36 outputs are displayed. Also here for lower values of  $\alpha$ , the generator outputs acceptable images much faster, as observed for the 2D case.

#### V. CONCLUSIONS

In this report we presented a review of GAN framework, as presented in the original paper [1], with a game theoretical approach, going through the resolution method to find a Nash equilibrium. We then proposed a parametrization of the loss function by varying the percentage of true images used for the training in each mini-batch, adapting the analysis previously presented to this new case. Then we tried to apply our own implementation of the model to three different datasets, two of which are generated ad-hoc and used as toy examples, while the

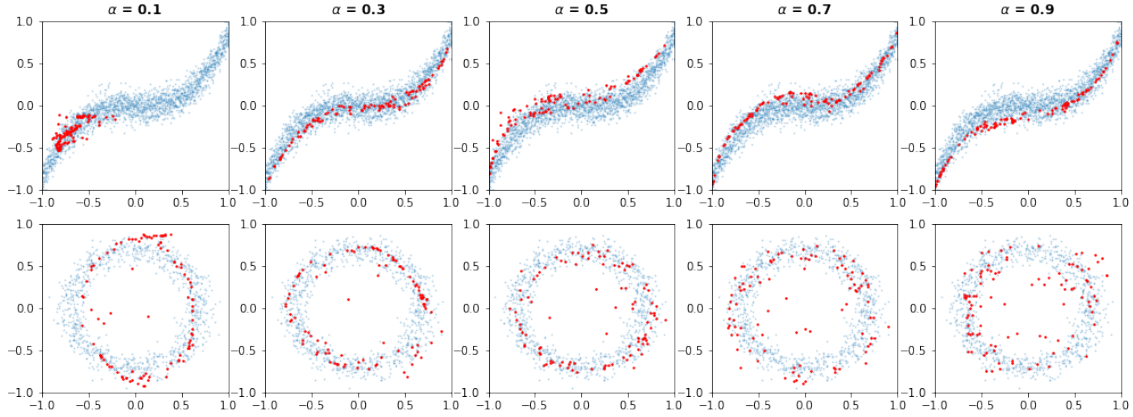


Figure 4: blue points are samples from the true distribution  $p_{data}$ , red dots are 100 samples generated from  $p_g$ , after the model was trained on 3000 epochs with mini-batch size 128 and a noise dimensionality of 10.

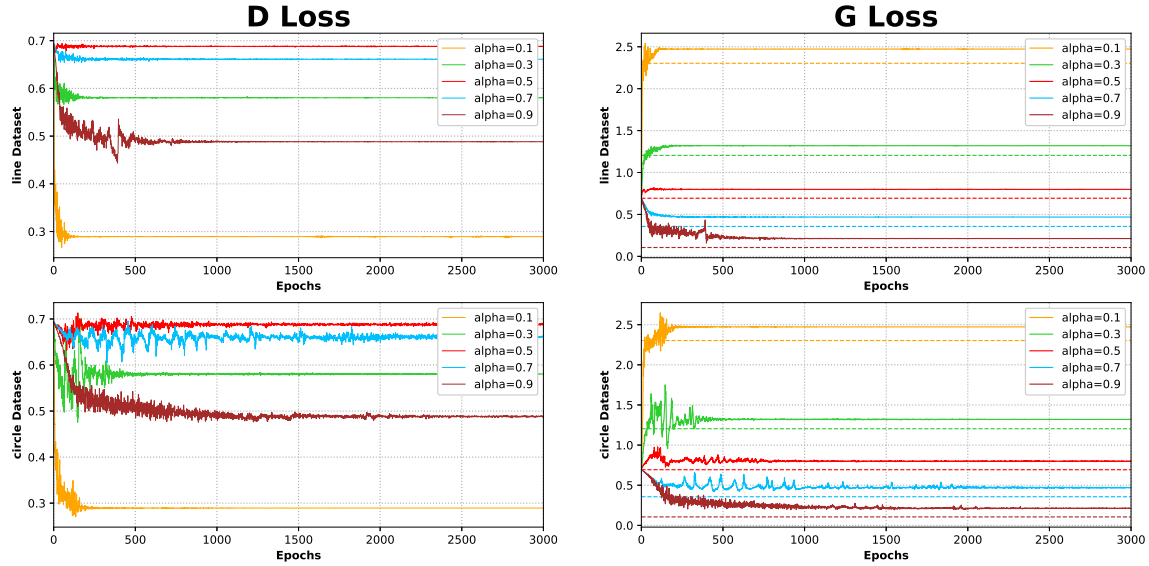


Figure 5: losses for 2D datasets of both G and D for all values of  $\alpha$ . Dashed lines in 'G Loss' plots represent the theoretical optimal values  $-\log(\alpha)$

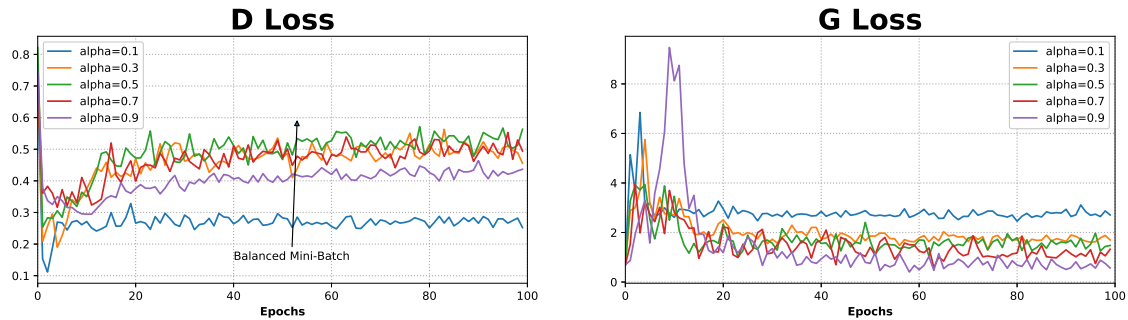


Figure 6: losses for 2D datasets of both G and D for all values of  $\alpha$ . Dashed lines in 'G Loss' plots represent the theoretical optimal values  $-\log(\alpha)$

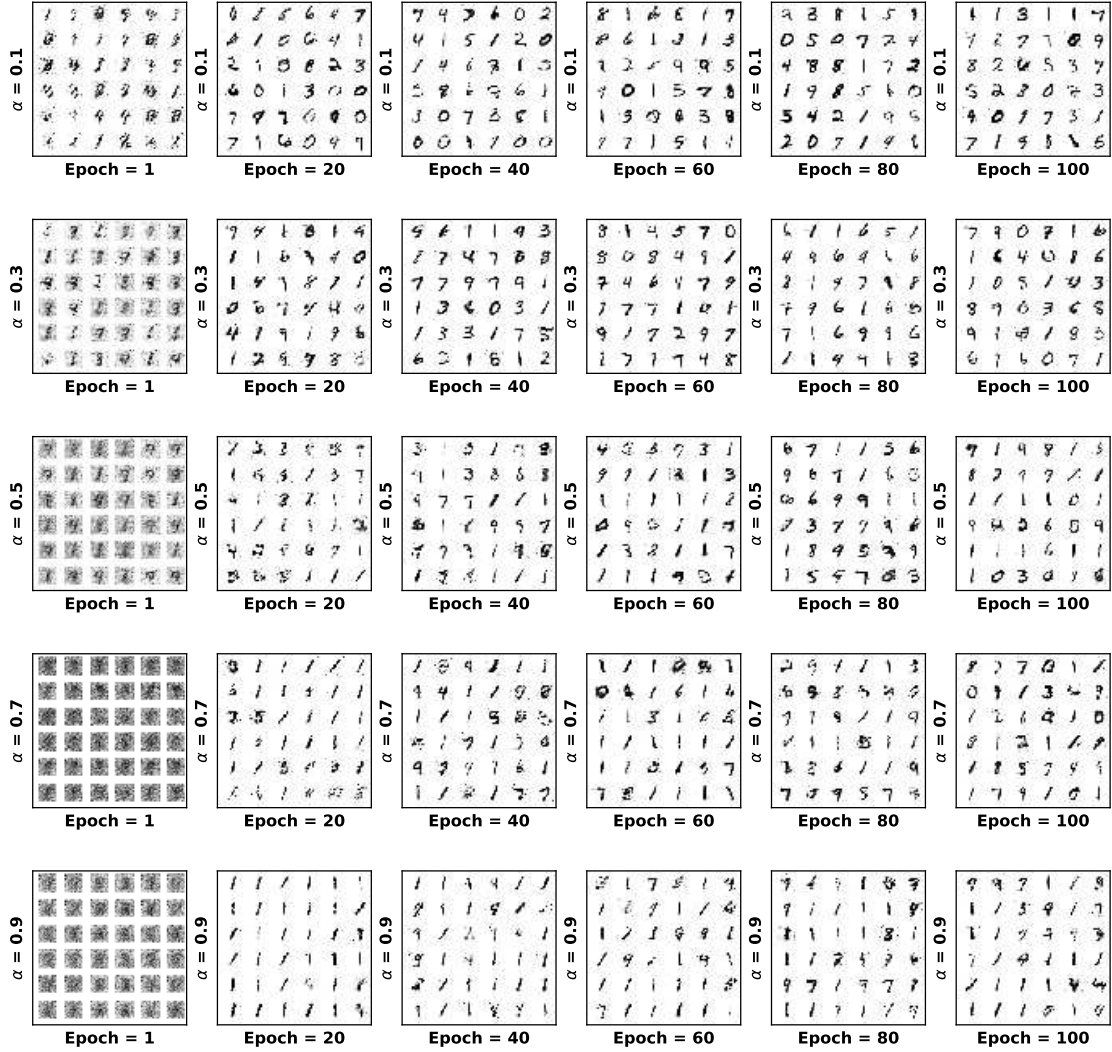


Figure 7: blue points are samples from the true distribution  $p_{data}$ , red dots are 100 samples generated from  $p_g$ , after the model was trained on 3000 epochs with mini-batch size 128 and a noise dimensionality of 10.

last one is the MNIST dataset, well known in literature for this kind of tasks. Results obtained show that when mini-batches are more balanced, the generator model performs better, but also trainings with less true samples can be performed with good results, which opens the door to the use of GANs even for slightly smaller datasets.

Another possible way to organise mini-batches, that could be investigated in the future, is to introduce randomness in the percentage of true samples per mini-batch:  $\alpha$  for example could stand for the probability of passing a true sample. The behaviour in this case should be similar, on average, to the one with  $\alpha$  indicating a deterministic ratio.

#### REFERENCES

- [1] Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. Generative adversarial nets. In Z. Ghahramani, M. Welling, C. Cortes, N. D. Lawrence, and K. Q. Weinberger, editors, *Advances in Neural Information Processing Systems 27*, pages 2672–2680. Curran Associates, Inc., 2014.
- [2] Ian J. Goodfellow. NIPS 2016 tutorial: Generative adversarial networks. *CoRR*, abs/1701.00160, 2017.
- [3] Mario Lucic, Karol Kurach, Marcin Michalski, Sylvain Gelly, and Olivier Bousquet. Are gans created equal? a large-scale study. *arXiv*, 2017.
- [4] F. A. Oliehoek, R. Savani, J. Gallego-Posada, E. van der Pol, E. D. de Jong, and R. Gross. GANGs: Generative Adversarial Network Games. *ArXiv e-prints*, December 2017.
- [5] Alec Radford, Luke Metz, and Soumith Chintala. Unsupervised representation learning with deep convolutional generative adversarial networks. *CoRR*, abs/1511.06434, 2015.
- [6] Han Zhang, Tao Xu, Hongsheng Li, Shaoting Zhang, Xiaolei Huang, Xiaogang Wang, and Dimitris N. Metaxas. Stackgan: Text to photo-realistic image synthesis with stacked generative adversarial networks. *CoRR*, abs/1612.03242, 2016.