
Objectives Towards Stable Adversarial Training Without Gradient Penalties

Christos Tsirigotis*
Mila, Université de Montréal,
Aristotle University of Thessaloniki

R. Devon Hjelm
Mila, Université de Montréal,
Microsoft Research

Aaron Courville†
Mila, Université de Montréal

Pericles A. Mitkas
Aristotle University of Thessaloniki

Abstract

Recent advances in adversarial learning observe that stabilization with gradient penalties trades off generated sample quality. We try to address this limitation by proposing novel objectives inspired by the logical XOR operation, which should not depend on gradient penalty regularization in order to be locally stable. In the sections that follow, we will present a theoretical study on this type of objective functions. We will prove global optimality conditions with similar assumptions as those made in the original GAN paper by Goodfellow et al., we will notice connections between the XOR-type objectives and the original GAN, and we will define non zero-sum parametric objectives based on that connection. Finally, we will attempt to study the local stability of the continuous-time training dynamical system around desirable equilibria.

1 Introduction

Generative Adversarial Networks (GAN) [1] offer a training methodology which has produced state-of-the-art generative models in terms of sample quality, enabling solutions that scale to vast datasets [2] as well as to large resolution images [3]. Though successful, stabilizing the training procedure still requires considerable effort from a practitioner’s point of view, and as a consequence, their training dynamics have attracted the research interest of the communities of optimization and game theory. Two complementary research directions are considered in this endeavour. The first tries to solve the stabilization problem by proposing, alternative to stochastic gradient descent, algorithms which are more suitable for saddle-point optimization [4, 5]. The second one tries instead to propose or modify the training objectives to be optimized, which in turn hopefully leads to better behaved training procedures [6, 7, 8]. This work proceeds in the second direction by proposing and studying novel adversarial objectives, that solve the same generative problem, utilizing the GAN framework, while trying to establish stable training and to avoid limitations introduced by existing methods.

A common limitation among existing methods is the usage of objectives which correspond to a zero-sum game between the generator and the critic network [9]. This is a known issue regarding the zero-sum game setting: An algorithm or regularized dynamics trying to reach a Nash equilibrium will eventually lead to “limit” cycles around desired equilibria. Modifications to the original objectives which depart from a zero-sum game, like the gradient penalty [10, 6, 7] or the non-saturating

*Work done for diploma thesis at Aristotle University of Thessaloniki. Correspondence to: christos.tsirigotis@umontreal.ca

†Canadian Institute for Advanced Research (CIFAR) Fellow

standard GAN objectives [1, 11], offer in certain cases stable training procedures in the cost of approximating solutions of a slightly different problem. The gradient penalty, in particular, tries to regularize the Lipschitz constant of the critic function and it is considered as standard technology for the robust training of GANs.

1.1 Discussion on constraining critic function’s Lipschitz constant

While originally devised [10] in order to satisfy the function class constraint of the Kantorovich metric between probability measures, which was introduced to the GAN framework by Arjovsky et al. (WGAN) [12], variations of the gradient penalty were soon proposed to stabilize adversarial training [6, 7]. These methods regularize the GAN or WGAN objectives so that certain theoretical properties about local convergence can be shown and they have been deployed in practice successfully. However, there have been experiments which suggest that utilizing a gradient penalty does not yield state-of-the-art performance, even though the training procedure is effectively stabilized [2]. Similar concerns have been attempted to be put forward by theoretical arguments [13]. Looking closely at proposition 1.1, we postulate that this behaviour is due to the final critic having zero gradient with respect to its input everywhere on the support of the real distribution. Mescheder et al. [7] made a seemingly necessary assumption for studying the equilibria of GAN and WGAN training; that the critic will have a constant value in a local neighbourhood of the real distribution’s support at equilibrium. Thus, the final critic is locally constant in all directions on every point in the real distribution’s support. So by combining this, with the hypothesis that the real distribution lies on a lower dimensional manifold [11], we hypothesize in particular that this condition predicts the limitation of the critic’s capacity to locally discriminate between real and fake samples, which are nearby, but not on, the support of the real distribution.

Proposition 1.1 (Stationary points of GAN and WGAN (Mescheder et al. [7]))

Let C_ψ be a parametric model of the critic function and G_θ be the parametric model of the generator function. Also, $\mathbb{Q}_\theta := G_\theta\#Z$ the induced measure in sample space by pushing random variable Z through model G_θ . \mathbb{P} is the target probability measure. Points (ψ^*, θ^*) of the joint parameter space consist equilibria of a system optimizing with GAN (1a) or WGAN (1b) parametric objectives.

$$\begin{aligned} \text{(GAN)} \quad & \mathbb{Q}_{\theta^*} = \mathbb{P} \quad \text{and} \quad C_{\psi^*}(x) = 0 \quad \text{and} \quad \nabla_x C_{\psi^*}(x) = 0 \quad \forall x \in \text{supp}\{\mathbb{P}\} \quad (1a) \\ \text{(WGAN)} \quad & \mathbb{Q}_{\theta^*} = \mathbb{P} \quad \text{and} \quad \nabla_x C_{\psi^*}(x) = 0 \quad \forall x \in \text{supp}\{\mathbb{P}\} \quad (1b) \end{aligned}$$

A parallel successful attempt to satisfy a Lipschitz constraint in the critic’s class is spectral normalization [8]. It modifies critic’s architecture in such a way so that its Lipschitz constant is approximately equal to 1. This has been deployed successfully in practice [2] for the original GAN objectives, leading to superior sample quality and easier optimization, however it does not seem to suffice for guaranteeing training’s stability. This could be because such normalization does not guarantee that the final condition for the critic’s gradients (proposition 1.1) will be met in the case of GAN’s or WGAN’s desired training equilibria. So, we believe that this method acts complimentary to those guaranteeing local stability around training’s equilibria.

2 Methodology

We motivate our search for alternative adversarial objective functions in methods which utilize more than one sample as inputs to the discriminator function. Some of them [14, 15, 16, 17, 18, 19, 13] derive the expressions for the objectives from the Maximum Mean Discrepancy [20] metric between probability measures, while others motivate their objective functions in an attempt to mitigate the mode dropping problem [17, 21, 22]. A recent method in this line of thought [23] suggests experimentally that defining objectives which act on the relative discrimination between real and fake samples may have actual benefits on the stability of a GAN’s training.

In this work we will use a differentiable analogue of the XOR logical operation, in order to define relative discrimination objective functions for adversarial generation. In particular, we think of a discriminator function $D: \mathbb{R}^d \times \mathbb{R}^d \rightarrow [0, 1]$ which accepts two sample inputs and tries to discriminate whether these two samples have been drawn from the same or from a different distribution. We are going to have D learn an XOR-type relation between samples from the two distribution.

this sense, D is trained to assign the "true" value when two samples are drawn from different distributions, and the "false" value when they are drawn from the same. Consequently, the adversarial game, we will initially consider, is formulated as such:

$$\min_{\mathbb{Q}} \max_D \mathbb{E}_{x \sim \mathbb{P}} \log(1 - D(x, y)) + \mathbb{E}_{x \sim \mathbb{P}} \log(D(x, y)) + \mathbb{E}_{x \sim \mathbb{Q}} \log(1 - D(x, y)) + \mathbb{E}_{x \sim \mathbb{Q}} \log(D(x, y)) \quad (2)$$

3 Theoretical Results

3.1 Global Optimality Analysis

For solving analytically the game described in eq. (2), we will make the following assumptions: First, analysis concerns the realizable case. Second, measures \mathbb{P} and \mathbb{Q} are absolutely continuous between themselves. Third, they both admit probability density functions, p and q respectively, under a common measure of reference, ν .

Lemma 3.1 (Optimal discriminator of (2))

$$D^*(x, y) = \frac{a(x, y)}{a(x, y) + b(x, y)} \quad \forall x, y \in \text{supp}\{\nu \times \nu\} \quad (3a)$$

$$a(x, y) = \frac{1}{2} (p(x)q(y) + q(x)p(y)) \quad (3b)$$

$$b(x, y) = \frac{1}{2} (p(x)p(y) + q(x)q(y)) \quad (3c)$$

Outside $\text{supp}\{\nu \times \nu\}$, D^* can take any real value.

Proof: Upon expanding the expectations in (2) and after algebraic manipulations, we get:

$$D^* = \arg \max_D 2 \iint \left\{ \log(D(x, y)) a(x, y) + \log(1 - D(x, y)) b(x, y) \right\} d\nu(x) d\nu(y) \quad (4)$$

Expressions (3b) and (3c) have an integral equal to 1 and consequently a and b can be considered as "mixture" densities of probability measures \mathbb{A} and \mathbb{B} on $\text{supp}\{\nu\} \times \text{supp}\{\nu\}$, absolutely continuous with respect to $\nu \times \nu$. Likewise to the analysis in Goodfellow et al. [1], the integral takes maximum value iff each integrated term is maximal for each $(x, y) \in \text{supp}\{\nu \times \nu\}$. So:

$$D^*(x, y) = \arg \max_d \left\{ \log(d) a(x, y) + \log(1 - d) b(x, y) \right\} \quad \forall (x, y) \in \text{supp}\{\nu \times \nu\} \quad (5)$$

The function of d to be maximized has unique extremal at $\frac{a}{a+b}$ and is strictly concave in $(0, 1)$. \square

Proposition 3.2 (Optimal generator of (2))

$$\mathbb{Q}^* = \mathbb{P} \quad (6)$$

Proof: Substituting (3) in (2), we get the following equivalent optimization problem for \mathbb{Q} :

$$\mathbb{Q}^* = \arg \min_{\mathbb{Q}} \left\{ 4 \text{JSD}(\mathbb{A}(\mathbb{Q}), \mathbb{B}(\mathbb{Q})) - 4 \log 2 \right\} \iff \mathbb{A}(\mathbb{Q}^*) = \mathbb{B}(\mathbb{Q}^*) \iff \quad (7a)$$

$$p(x)q^*(y) + q^*(x)p(y) = p(x)p(y) + q^*(x)q^*(y) \iff \quad (7b)$$

$$(p(x) - q^*(x))(p(y) - q^*(y)) = 0 \iff q^* = p \iff \mathbb{Q}^* = \mathbb{P} \quad (7c)$$

\square

3.2 XORGAN Parametric Objectives

We notice that we can factor the expression of optimal discriminator in lemma 3.1 with terms of the original GAN's optimal discriminator expression [1], $D_1^*(x)$, as such:

$$D_{\text{XOR}}^*(x, y) = \frac{p(x)q(y) + q(x)p(y)}{(p(x) + q(x))(p(y) + q(y))} = D_1^*(x)(1 - D_1^*(y)) + (1 - D_1^*(x))D_1^*(y) \quad (8)$$

This motivates us to define the XORGAN’s objectives using a critic function of a single sample input. Furthermore, this allows us to state the following equivalent non zero-sum game between the critic and the generator function:

Definition 3.3 (Parametric XORGAN objective functions)

$$\arg \max_{\psi} \mathcal{L}_{\text{GAN}}(\psi, \theta) := \mathbb{E}_{x \sim \mathbb{P}} \log(\sigma(C_{\psi}(x))) + \mathbb{E}_{x \sim \mathbb{Q}_{\theta}} \log(\sigma(-C_{\psi}(x))) \quad (9)$$

$$\arg \min_{\theta} \mathcal{L}_{\text{XOR}}(\psi, \theta) := \mathcal{S}(\mathbb{P}, \mathbb{P}, C_{\psi}) + \mathcal{S}(\mathbb{Q}_{\theta}, \mathbb{Q}_{\theta}, C_{\psi}) + 2\mathcal{D}(\mathbb{P}, \mathbb{Q}_{\theta}, C_{\psi}) \quad (10)$$

$$\mathcal{D}(\mathbb{P}, \mathbb{Q}, C) := \mathbb{E}_{\substack{x \sim \mathbb{P} \\ y \sim \mathbb{Q}}} \log(\sigma(C(x))\sigma(-C(y)) + \sigma(-C(x))\sigma(C(y))) \quad (11)$$

$$\mathcal{S}(\mathbb{P}, \mathbb{Q}, C) := \mathbb{E}_{\substack{x \sim \mathbb{P} \\ y \sim \mathbb{Q}}} \log(\sigma(C(x))\sigma(C(y)) + \sigma(-C(x))\sigma(-C(y))) \quad (12)$$

The expressions above can be derived by modelling $D_1(x; \psi)$ by $\sigma(C_{\psi}(x))$ as in [1] and substituting appropriately into a model for $D_{\text{XOR}}(x, y; \psi)$, and then in eq. (2).

3.3 Stability Analysis in Continuous-Time Training Dynamics

We now present the training of a XORGAN as a continuous-time dynamical system, as in [24, 7], by considering gradient descent/ascent as its optimization algorithm:

$$v(\psi, \theta) := \begin{pmatrix} \dot{\psi} \\ \dot{\theta} \end{pmatrix} = \begin{pmatrix} \nabla_{\psi} \mathcal{L}_{\text{GAN}}(\psi, \theta) \\ -\nabla_{\theta} \mathcal{L}_{\text{XOR}}(\psi, \theta) \end{pmatrix} \quad (13)$$

Then, the following statements can be proven regarding its stability analysis:

Proposition 3.4 (Stationary points of XORGAN)

Points (ψ^*, θ^*) of the joint parameter space consist equilibria of the system which occurs from the optimization of XORGAN parametric objectives, as described in definition 3.3.

$$(XORGAN) \quad \mathbb{Q}_{\theta^*} = \mathbb{P} \quad \text{and} \quad C_{\psi^*}(x) = 0 \quad \forall x \in \text{supp}\{\mathbb{P}\} \quad (14)$$

Also, the system’s Jacobian at these points is negative semi-definite and it has real eigenvalues only.

Proof of proposition 3.4 can be found in supplementary material.

4 Conclusions and Future Work

We hope that these primary findings can serve as motivation for further investigation towards this direction. First of all, however, we are still seeking to finalize an argument towards asymptotic stability to a forward invariant set of the training dynamics, using Lyapunov arguments. That would reassure us about XORGAN’s local stability properties compared to existing analyses of other objectives. Nevertheless, the finding of desirable and provable equilibria with negative semi-definite Jacobian indicates that we are searching in a good direction. This view can be reinforced by the simulation experiments on toy problems, presented in fig. 1 and in supplementary material. Second, in future work, we seek to compare experimentally in image generation tasks against benchmark adversarial objectives, like GAN [1] or WGAN [12], as well as with more recent approaches like the RGAN [23]. Finally, we hope that current work will motivate simple experiments which would investigate the limitations of gradient penalty regularization in adversarial learning, if any.

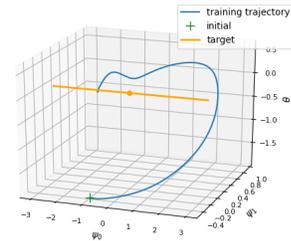


Figure 1: Simulation with $C_{\psi}(x) = \psi_0 x + \psi_1$ and $\mathbb{Q}_{\theta} = \delta_{\theta}$

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] Andrew Brock, Jeff Donahue, and Karen Simonyan. Large Scale GAN Training for High Fidelity Natural Image Synthesis. In *International Conference on Learning Representations*, 2019.
- [3] Tero Karras, Samuli Laine, and Timo Aila. A Style-Based Generator Architecture for Generative Adversarial Networks. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 4401–4410, 2019.
- [4] Constantinos Daskalakis, Andrew Ilyas, Vasilis Syrgkanis, and Haoyang Zeng. Training GANs with Optimism. In *International Conference on Learning Representations*, 2018.
- [5] Gauthier Gidel, Hugo Berard, Gaëtan Vignoud, Pascal Vincent, and Simon Lacoste-Julien. A Variational Inequality Perspective on Generative Adversarial Networks. In *International Conference on Learning Representations*, 2019.
- [6] Kevin Roth, Aurelien Lucchi, Sebastian Nowozin, and Thomas Hofmann. Stabilizing Training of Generative Adversarial Networks through Regularization. In I. Guyon, U. V. Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett, editors, *Advances in Neural Information Processing Systems 30*, pages 2018–2028. Curran Associates, Inc., 2017.
- [7] Lars Mescheder, Andreas Geiger, and Sebastian Nowozin. Which Training Methods for GANs do actually Converge? In *Proceedings of the 35th International Conference on Machine Learning*, volume 80 of *Proceedings of Machine Learning Research*, pages 3481–3490. PMLR, Stockholm, Sweden, July 2018.
- [8] Takeru Miyato, Toshiki Kataoka, Masanori Koyama, and Yuichi Yoshida. Spectral Normalization for Generative Adversarial Networks. In *International Conference on Learning Representations*, 2018.
- [9] Panayotis Mertikopoulos, Christos Papadimitriou, and Georgios Piliouras. Cycles in Adversarial Regularized Learning. In *Proceedings of the Twenty-Ninth Annual ACM-SIAM Symposium on Discrete Algorithms, SODA '18*, pages 2703–2717, Philadelphia, PA, USA, 2018. Society for Industrial and Applied Mathematics.
- [10] Ishaan Gulrajani, Faruk Ahmed, Martin Arjovsky, Vincent Dumoulin, and Aaron C Courville. Improved Training of Wasserstein GANs. In I. Guyon, U. V. Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett, editors, *Advances in Neural Information Processing Systems 30*, pages 5767–5777. Curran Associates, Inc., 2017.
- [11] Martin Arjovsky and Léon Bottou. Towards Principled Methods for Training Generative Adversarial Networks. In *International Conference on Learning Representations*, 2017.
- [12] Martin Arjovsky, Soumith Chintala, and Léon Bottou. Wasserstein Generative Adversarial Networks. In *Proceedings of the 34th International Conference on Machine Learning*, volume 70 of *Proceedings of Machine Learning Research*, pages 214–223. PMLR, International Convention Centre, Sydney, Australia, 2017.
- [13] Wei Wang, Yuan Sun, and Saman Halgamuge. Improving MMD-GAN Training with Repulsive Loss Function. In *International Conference on Learning Representations*, 2019.
- [14] Yujia Li, Kevin Swersky, and Rich Zemel. Generative Moment Matching Networks. In *Proceedings of the 32nd International Conference on Machine Learning*, volume 37 of *Proceedings of Machine Learning Research*, pages 1718–1727. PMLR, Lille, France, 2015.
- [15] Gintare Karolina Dziugaite, Daniel M. Roy, and Zoubin Ghahramani. Training generative neural networks via maximum mean discrepancy optimization. *arXiv preprint arXiv:1505.03906*, 2015.

- [16] Chun-Liang Li, Wei-Cheng Chang, Yu Cheng, Yiming Yang, and Barnabas Poczos. MMD GAN: Towards Deeper Understanding of Moment Matching Network. In I. Guyon, U. V. Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett, editors, *Advances in Neural Information Processing Systems 30*, pages 2203–2213. Curran Associates, Inc., 2017.
- [17] Thomas Unterthiner, Bernhard Nessler, Günter Klambauer, Martin Heusel, Hubert Ramsauer, and Sepp Hochreiter. Coulomb GANs: Provably Optimal Nash Equilibria via Potential Fields. In *International Conference on Learning Representations*, 2018.
- [18] Mikołaj Bińkowski, Dougal J. Sutherland, Michael Arbel, and Arthur Gretton. Demystifying MMD GANs. In *International Conference on Learning Representations*, 2018.
- [19] Michael Arbel, Dougal Sutherland, Mikołaj Bińkowski, and Arthur Gretton. On gradient regularizers for MMD GANs. In S. Bengio, H. Wallach, H. Larochelle, K. Grauman, N. Cesa-Bianchi, and R. Garnett, editors, *Advances in Neural Information Processing Systems 31*, pages 6701–6711. Curran Associates, Inc., 2018.
- [20] Bharath K. Sriperumbudur, Arthur Gretton, Kenji Fukumizu, Gert Lanckriet, and Bernhard Schölkopf. Injective Hilbert space embeddings of probability measures. In *Proceedings of the 21st Annual Conference on Learning Theory*, Helsinki, Finland, 2008.
- [21] Zinan Lin, Ashish Khetan, Giulia Fanti, and Sewoong Oh. PacGAN: The power of two samples in generative adversarial networks. In S. Bengio, H. Wallach, H. Larochelle, K. Grauman, N. Cesa-Bianchi, and R. Garnett, editors, *Advances in Neural Information Processing Systems 31*, pages 1505–1514. Curran Associates, Inc., 2018.
- [22] Thomas Lucas, Corentin Tallec, Yann Ollivier, and Jakob Verbeek. Mixed batches and symmetric discriminators for GAN training. In *Proceedings of the 35th International Conference on Machine Learning*, volume 80 of *Proceedings of Machine Learning Research*, pages 2844–2853. PMLR, Stockholm, Sweden, July 2018.
- [23] Alexia Jolicoeur-Martineau. The relativistic discriminator: A key element missing from standard GAN. In *International Conference on Learning Representations*, 2019.
- [24] Vaishnavh Nagarajan and J. Zico Kolter. Gradient descent GAN optimization is locally stable. In I. Guyon, U. V. Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett, editors, *Advances in Neural Information Processing Systems 30*, pages 5585–5595. Curran Associates, Inc., 2017.
- [25] Jingyi Xu, Zilu Zhang, Tal Friedman, Yitao Liang, and Guy Broeck. A Semantic Loss Function for Deep Learning with Symbolic Knowledge. In *Proceedings of the 35th International Conference on Machine Learning*, volume 80 of *Proceedings of Machine Learning Research*, pages 5502–5511. PMLR, Stockholm, Sweden, July 2018.
- [26] Diederik Kingma and Jimmy Ba. Adam: A Method for Stochastic Optimization. In *International Conference on Learning Representations*, 2015.

5 Supplementary Material

5.1 Remarks on XORGAN Optimal Discriminator

5.1.1 In relation to the original GAN objectives

Let us look into eq. (3) of the optimal discriminator. We could try to represent this with a model of two samples inputs. However, this would be more computationally inefficient than representing it with a model of one sample input.

Remark 5.1 *optimal GAN discriminator appears in optimal XORGAN discriminator:* Substitute expressions for \mathbb{A} (3b) and \mathbb{B} (3c) and after algebraic manipulations we get:

$$D_{\text{XOR}}^*(x, y) = \frac{p(x)q(y) + q(x)p(y)}{(p(x) + q(x))(p(y) + q(y))} \quad (15a)$$

$$= \frac{p(x)}{p(x) + q(x)} \frac{q(y)}{p(y) + q(y)} + \frac{q(x)}{p(x) + q(x)} \frac{p(y)}{p(y) + q(y)} \quad (15b)$$

$$= \frac{p(x)}{p(x) + q(x)} \left(1 - \frac{p(y)}{p(y) + q(y)}\right) + \left(1 - \frac{p(x)}{p(x) + q(x)}\right) \frac{p(y)}{p(y) + q(y)} \quad (15c)$$

$$= D_{\text{GAN}}^*(x)(1 - D_{\text{GAN}}^*(y)) + (1 - D_{\text{GAN}}^*(x))D_{\text{GAN}}^*(y) \quad (15d)$$

This way we can isolate the same term multiple times within the same expression, $D_{\text{GAN}}^* = \frac{p}{p+q}$, hence simplifying the representation into an operation between functions of a single sample input.

We could have isolated the term $\frac{q}{p+q}$ instead. This would mean that an original GAN targeted 1 for the fake data and 0 for the real. However, such a choice does not affect the final expression for XORGAN's optimal discriminator as it is symmetric with respect to p and q .

Furthermore, it was possible to derive our relativistic objectives in the same way for the NXOR logical operation. The optimal discriminator, we would find in that case, would be:

$$D_{\text{NXOR}}^*(x, y) = \frac{p(x)p(y) + q(x)q(y)}{(p(x) + q(x))(p(y) + q(y))} \quad (16)$$

where again we could have isolated the same terms of a single sample input without affecting the final expression due to symmetry. It is true, however, that, independently from such choices, the objective functional expression remains invariant for every adversarial objectives studied (GAN, WGAN, XORGAN, and others). In any case, the final objective expression is irrelevant to the choice of predetermined targets for the discriminator. In particular, this seems to be the case for XORGAN as well because the following relation holds:

$$D_{\text{XOR}}^*(x, y) + D_{\text{NXOR}}^*(x, y) = 1 \quad (17)$$

5.1.2 In relation to the logical events: same vs different distribution

We chose to represent bijectively the interval $[0, 1]$, which is the image of probability measures, by using the sigmoid function. By substituting the optimal discriminator for the GAN objectives with the composite function $\sigma \circ C^*$, where as $C: \mathbb{R}^d \rightarrow \mathbb{R}$ we define the critic model:

$$\begin{aligned} D_{\text{XOR}}^*(x, y) &= \sigma(C^*(x)) \left(1 - \sigma(C^*(y))\right) + \left(1 - \sigma(C^*(x))\right) \sigma(C^*(y)) \\ &= \sigma(C^*(x)) \sigma(-C^*(y)) + \sigma(-C^*(x)) \sigma(C^*(y)) \end{aligned} \quad (18)$$

Similarly, for the NXOR alternative, we can show that:

$$D_{\text{NXOR}}^*(x, y) = \sigma(C^*(x)) \sigma(C^*(y)) + \sigma(-C^*(x)) \sigma(-C^*(y)) \quad (19)$$

We will attempt to give semantics to the derived discriminator expressions. We warn the reader, however, that the following attempt is rushed and immature, more intuitive than well-defined mathematically. More mature and general connection of objective functions to probabilistic semantics can be found in the work of Xu et al. [25]. We find however the following perspective to be useful to our intuition, so we are going to develop it in the current text.

We like to think that the expression for $D_{\text{XOR}}^*(x, y)$ (18) corresponds to the logical proposition:

$$\mathbf{D}xy := (\mathbf{P}x \wedge \mathbf{Q}y) \vee (\mathbf{Q}x \wedge \mathbf{P}y) \quad (20)$$

Proposition **D** (20) is evaluated as true if and only if its inputs are samples of different distributions. Similarly, the expression for $D_{\text{NXOR}}^*(x, y)$ (19) corresponds to the logical proposition:

$$\mathbf{S}xy := (\mathbf{P}x \wedge \mathbf{P}y) \vee (\mathbf{Q}x \wedge \mathbf{Q}y) \quad (21)$$

Proposition **S** (21) is evaluated as true if and only if its inputs are samples from the same distribution. We further see that the disjunction of these two proposition is a tautology (22), which is aligned with the observation made in eq. (17). Due to symmetry, we would derive the same remarks even if we considered $\sigma(C^*(x))$ to represent $\frac{q}{p+q}$.

$$\models (\mathbf{D}xy \vee \mathbf{S}xy) \Leftrightarrow (\mathbf{P}x \wedge \mathbf{Q}y) \vee (\mathbf{Q}x \wedge \mathbf{P}y) \vee (\mathbf{P}x \wedge \mathbf{P}y) \vee (\mathbf{Q}x \wedge \mathbf{Q}y) \quad (22)$$

Having remarked these intuitive relations, we define two objective functionals with respect to two probability measure inputs and a critic function, which the discriminator is consisted of, C . We name them by the first letters of the english words "Different" and "Same", reminding their intuitive utility:

$$\mathcal{D}(\mathbb{P}, \mathbb{Q}, C) := \mathbb{E}_{\substack{x \sim \mathbb{P} \\ y \sim \mathbb{Q}}} \log \left(\sigma(C(x)) \sigma(-C(y)) + \sigma(-C(x)) \sigma(C(y)) \right) \quad (23)$$

$$\mathcal{S}(\mathbb{P}, \mathbb{Q}, C) := \mathbb{E}_{\substack{x \sim \mathbb{P} \\ y \sim \mathbb{Q}}} \log \left(\sigma(C(x)) \sigma(C(y)) + \sigma(-C(x)) \sigma(-C(y)) \right) \quad (24)$$

The initial problem (2) is written as such, when expressed by the functionals (23) and (24):

$$\arg \min_{\mathbb{Q}} \arg \max_C \mathcal{S}(\mathbb{P}, \mathbb{P}, C) + \mathcal{D}(\mathbb{P}, \mathbb{Q}, C) + \mathcal{D}(\mathbb{Q}, \mathbb{P}, C) + \mathcal{S}(\mathbb{Q}, \mathbb{Q}, C) \quad (25)$$

By observing that the functional \mathcal{D} (23) is symmetric with respect to its probability measure inputs, we can simplify the last expression:

$$\arg \min_{\mathbb{Q}} \arg \max_C \mathcal{S}(\mathbb{P}, \mathbb{P}, C) + \mathcal{S}(\mathbb{Q}, \mathbb{Q}, C) + 2\mathcal{D}(\mathbb{P}, \mathbb{Q}, C) \quad (26)$$

5.2 Proofs about Stability of Non Zero-sum XORGAN Training

As it is described in definition 3.3, we chose to have the original GAN objective for training the discriminator in the non zero-sum game formulation. We remind that this choice does not alter the solutions of the initial problem, due to the symmetry of the optimal XOR discriminator as we have seen in section 5.1.1. We are restating the non zero-sum objectives (definition 3.3) that we are going to study, as well as the definitions for the functionals \mathcal{D} (23) and \mathcal{S} (24).

Definition 5.2 (Non zero-sum XORGAN objective functions)

$$\arg \max_{\psi} \mathcal{L}_{\text{GAN}}(\psi, \theta) := \mathbb{E}_{x \sim \mathbb{P}} \log \left(\sigma(C_{\psi}(x)) \right) + \mathbb{E}_{x \sim \mathbb{Q}_{\theta}} \log \left(\sigma(-C_{\psi}(x)) \right) \quad (27a)$$

$$\arg \min_{\theta} \mathcal{L}_{\text{XOR}}(\psi, \theta) := \mathcal{S}(\mathbb{P}, \mathbb{P}, C_{\psi}) + \mathcal{S}(\mathbb{Q}_{\theta}, \mathbb{Q}_{\theta}, C_{\psi}) + 2\mathcal{D}(\mathbb{P}, \mathbb{Q}_{\theta}, C_{\psi}) \quad (27b)$$

Training with gradient descent/ascent, stated as a continuous-time dynamical system, becomes:

$$\begin{aligned} v(\psi, \theta) &:= \begin{pmatrix} \dot{\psi} \\ \dot{\theta} \end{pmatrix} = \begin{pmatrix} \nabla_{\psi} \mathcal{L}_{\text{GAN}}(\psi, \theta) \\ -\nabla_{\theta} \mathcal{L}_{\text{XOR}}(\psi, \theta) \end{pmatrix} \quad (28) \\ &= \begin{pmatrix} \mathbb{E}_{x \sim \mathbb{P}} \left[\sigma(-C_{\psi}(x)) \nabla_{\psi} C_{\psi}(x) \right] - \mathbb{E}_{x \sim \mathbb{Q}_{\theta}} \left[\sigma(C_{\psi}(x)) \nabla_{\psi} C_{\psi}(x) \right] \\ -\nabla_{\theta_1} \mathcal{S}(\mathbb{Q}_{\theta_1}, \mathbb{Q}_{\theta_1}, C_{\psi}) \big|_{\theta_1=\theta} - \nabla_{\theta_2} \mathcal{S}(\mathbb{Q}_{\theta_2}, \mathbb{Q}_{\theta_2}, C_{\psi}) \big|_{\theta_2=\theta} - 2\nabla_{\theta} \mathcal{D}(\mathbb{P}, \mathbb{Q}_{\theta}, C_{\psi}) \end{pmatrix} \end{aligned}$$

Proposition 5.3 (Non zero-sum XORGAN equilibria)

Points (ψ^*, θ^*) of the joint parameter space consist stationary points of the system which occurs from the optimization of non zero-sum XORGAN parametric objectives, as described in eq. (28).

$$\mathbb{Q}_{\theta^*} = \mathbb{P} \quad \text{and} \quad C_{\psi^*}(x) = 0 \quad \forall x \in \text{supp}\{\mathbb{P}\} \quad (29)$$

Proof: We can easily verify that points (29) are indeed stationary points, as the time-derivative of the states equals to zero:

$$\begin{aligned} \nabla_{\psi} \mathcal{L}_{\text{GAN}}(\psi, \theta) \Big|_{\substack{\psi=\psi^* \\ \theta=\theta^*}} &= \quad (30) \\ &= \mathbb{E}_{x \sim \mathbb{P}} \left[\sigma(-C_{\psi^*}(x)) \nabla_{\psi} C_{\psi}(x) \Big|_{\psi=\psi^*} \right] - \mathbb{E}_{x \sim \mathbb{Q}_{\theta^*}} \left[\sigma(C_{\psi^*}(x)) \nabla_{\psi} C_{\psi}(x) \Big|_{\psi=\psi^*} \right] \\ &= \mathbb{E}_{x \sim \mathbb{P}} \left[(\sigma(0) - \sigma(0)) \nabla_{\psi} C_{\psi}(x) \Big|_{\psi=\psi^*} \right] = 0 \end{aligned}$$

For finding the derivative with respect to θ , let us remember that the functionals \mathcal{S} (24) and \mathcal{D} (23) are symmetric with respect to their measure inputs, as also that $D_{\text{NXOR}} = 1 - D_{\text{XOR}}$, which is apparent from the following definitions:

$$D_{\text{NXOR}}(x, y, \psi) := \sigma(C_{\psi}(x)) \sigma(C_{\psi}(y)) + \sigma(-C_{\psi}(x)) \sigma(-C_{\psi}(y)) \quad (31a)$$

$$D_{\text{XOR}}(x, y, \psi) := \sigma(C_{\psi}(x)) \sigma(-C_{\psi}(y)) + \sigma(-C_{\psi}(x)) \sigma(C_{\psi}(y)) \quad (31b)$$

Calculating the derivatives of the constituent expressions:

$$\nabla_{\theta} \mathcal{S}(\mathbb{P}, \mathbb{Q}_{\theta}, C_{\psi}) = \nabla_{\theta} \left\{ \mathbb{E}_{\substack{x \sim \mathbb{P} \\ z \sim Z}} \log \left(D_{\text{NXOR}}(x, G_{\theta}(z), \psi) \right) \right\} \quad (32a)$$

$$= \mathbb{E}_{\substack{x \sim \mathbb{P} \\ z \sim Z}} \left[\frac{1}{D_{\text{NXOR}}(x, G_{\theta}(z), \psi)} \nabla_{\theta} G_{\theta}(z) \nabla_y D_{\text{NXOR}}(x, y, \psi) \Big|_{y=G_{\theta}(z)} \right]$$

$$\nabla_{\theta} \mathcal{D}(\mathbb{P}, \mathbb{Q}_{\theta}, C_{\psi}) = \mathbb{E}_{\substack{x \sim \mathbb{P} \\ z \sim Z}} \left[\frac{1}{D_{\text{XOR}}(x, G_{\theta}(z), \psi)} \nabla_{\theta} G_{\theta}(z) \nabla_y D_{\text{XOR}}(x, y, \psi) \Big|_{y=G_{\theta}(z)} \right] \quad (32b)$$

Also for D_{NXOR} and D_{XOR} :

$$\nabla_y D_{\text{NXOR}}(x, y, \psi) = \sigma(C_{\psi}(y)) \sigma(-C_{\psi}(y)) \left[\sigma(C_{\psi}(x)) - \sigma(-C_{\psi}(x)) \right] \nabla_y C_{\psi}(y) \quad (33a)$$

$$\nabla_y D_{\text{XOR}}(x, y, \psi) = -\nabla_y D_{\text{NXOR}}(x, y, \psi) \quad (33b)$$

As $\forall x^*, y^* \in \text{supp}\{\mathbb{P}\} \quad C_{\psi^*}(x^*) = 0$, according to eq. (29), we find the following values for eqs. (31) and (33):

$$D_{\text{NXOR}}(x^*, y^*, \psi^*) = \sigma(0) \sigma(0) + \sigma(0) \sigma(0) = \frac{1}{2} \quad (34a)$$

$$D_{\text{XOR}}(x^*, y^*, \psi^*) = \frac{1}{2} \quad (34b)$$

$$\nabla_y D_{\text{NXOR}}(x^*, y, \psi^*) \Big|_{y=y^*} = \nabla_y D_{\text{NXOR}}(x^*, y, \psi^*) \Big|_{y=y^*} = 0 \quad (34c)$$

Now to find the values for (32), we make the following observation, beginning from eq. (29):

$$\mathbb{Q}_{\theta^*} = \mathbb{P} \implies \text{supp}\{\mathbb{Q}_{\theta^*}\} = \text{supp}\{\mathbb{P}\} \quad (35a)$$

$$\implies G_{\theta^*}(z) \in \text{supp}\{\mathbb{P}\} \quad \forall z \in \text{supp}\{Z\} \quad (35b)$$

$$\implies C_{\psi^*}(G_{\theta^*}(z)) = 0 \quad \forall z \in \text{supp}\{Z\} \quad (35c)$$

By using eqs. (29), (34) and (35c) in (32), we observe that those are zeroed out, and consequently by substituting into (13), we verify that indeed (ψ^*, θ^*) are stationary points. \square

We will examine the behaviour of the system (13) locally around the stationary points (29). For this reason we will linearize the system by expanding its first order. So we are extracting the Jacobian of the system, $J(\psi, \theta)$:

$$\begin{pmatrix} \dot{\psi} \\ \dot{\theta} \end{pmatrix} \simeq J(\psi^*, \theta^*) \begin{pmatrix} \psi - \psi^* \\ \theta - \theta^* \end{pmatrix} = \begin{pmatrix} \nabla_{\psi\psi}^2 \mathcal{L}_{\text{GAN}}(\psi, \theta) \Big|_{\substack{\psi=\psi^* \\ \theta=\theta^*}} & \nabla_{\psi\theta}^2 \mathcal{L}_{\text{GAN}}(\psi, \theta) \Big|_{\substack{\psi=\psi^* \\ \theta=\theta^*}} \\ -\nabla_{\theta\psi}^2 \mathcal{L}_{\text{XOR}}(\psi, \theta) \Big|_{\substack{\psi=\psi^* \\ \theta=\theta^*}} & -\nabla_{\theta\theta}^2 \mathcal{L}_{\text{XOR}}(\psi, \theta) \Big|_{\substack{\psi=\psi^* \\ \theta=\theta^*}} \end{pmatrix} \begin{pmatrix} \psi - \psi^* \\ \theta - \theta^* \end{pmatrix} \quad (36)$$

Lemma 5.4 (Jacobian at SP (29) of non zero-sum XORGAN)

The Jacobian of the system (28), J , at the stationary points (29) is:

$$J(\psi^*, \theta^*) = \begin{pmatrix} J_{\psi\psi} & J_{\psi\theta} \\ J_{\theta\psi} & J_{\theta\theta} \end{pmatrix} \quad (37a)$$

$$J_{\psi\psi} = -\frac{1}{2} \mathbb{E}_{x \sim \mathbb{P}} \left[\nabla_{\psi} C_{\psi}(x) \Big|_{\psi=\psi^*} \nabla_{\psi} C_{\psi}(x) \Big|_{\psi=\psi^*}^T \right] \quad (37b)$$

$$J_{\psi\theta} = -\frac{1}{4} \nabla_{\theta} \mathbb{E}_{x \sim \mathbb{Q}_{\theta}} \left[(C_{\psi^*}(x) + 2) \nabla_{\psi} C_{\psi}(x) \Big|_{\psi=\psi^*} \right] \Big|_{\theta=\theta^*} \quad (37c)$$

$$J_{\theta\psi} = 0 \quad (37d)$$

$$J_{\theta\theta} = -\frac{1}{2} \nabla_{\theta} \mathbb{E}_{x \sim \mathbb{Q}_{\theta}} [C_{\psi^*}(x)] \Big|_{\theta=\theta^*} \nabla_{\theta} \mathbb{E}_{x \sim \mathbb{Q}_{\theta}} [C_{\psi^*}(x)] \Big|_{\theta=\theta^*}^T \quad (37e)$$

Proof: We will first find the values for $J_{\psi\psi}$ and $J_{\psi\theta}$ terms of J , which correspond to the second derivatives of the original GAN critic objective function (27a). We are going to reuse observations that we used in the proof of proposition 5.3.

$$\nabla_{\psi\psi}^2 \mathcal{L}_{\text{GAN}}(\psi, \theta) \Big|_{\substack{\psi=\psi^* \\ \theta=\theta^*}} = \quad (38)$$

$$\begin{aligned} &= \mathbb{E}_{x \sim \mathbb{P}} \left[-\sigma(C_{\psi^*}(x)) \sigma(-C_{\psi^*}(x)) \nabla_{\psi} C_{\psi}(x) \Big|_{\psi=\psi^*} \nabla_{\psi} C_{\psi}(x) \Big|_{\psi=\psi^*}^T + \sigma(-C_{\psi^*}(x)) \nabla_{\psi\psi}^2 C_{\psi}(x) \Big|_{\psi=\psi^*} \right] \\ &+ \mathbb{E}_{x \sim \mathbb{Q}_{\theta^*}} \left[-\sigma(C_{\psi^*}(x)) \sigma(-C_{\psi^*}(x)) \nabla_{\psi} C_{\psi}(x) \Big|_{\psi=\psi^*} \nabla_{\psi} C_{\psi}(x) \Big|_{\psi=\psi^*}^T - \sigma(-C_{\psi^*}(x)) \nabla_{\psi\psi}^2 C_{\psi}(x) \Big|_{\psi=\psi^*} \right] \\ &= -\frac{1}{2} \mathbb{E}_{x \sim \mathbb{P}} \left[\nabla_{\psi} C_{\psi}(x) \Big|_{\psi=\psi^*} \nabla_{\psi} C_{\psi}(x) \Big|_{\psi=\psi^*}^T \right] \end{aligned}$$

$$\nabla_{\psi\theta}^2 \mathcal{L}_{\text{GAN}}(\psi, \theta) \Big|_{\substack{\psi=\psi^* \\ \theta=\theta^*}} = \quad (39)$$

$$\begin{aligned} &= \mathbb{E}_{z \sim Z} \left[\left(-\sigma(C_{\psi^*}(x)) \sigma(-C_{\psi^*}(x)) \nabla_{\psi} C_{\psi}(x) \Big|_{\psi=\psi^*} \nabla_x C_{\psi^*}(x)^T \right. \right. \\ &\quad \left. \left. - \sigma(C_{\psi^*}(x)) \nabla_{\psi x}^2 C_{\psi}(x) \Big|_{\psi=\psi^*} \right) \Big|_{x=G_{\theta^*}(z)} \nabla_{\theta} G_{\theta}(z) \Big|_{\theta=\theta^*} \right] \\ &= -\frac{1}{4} \nabla_{\theta} \mathbb{E}_{x \sim \mathbb{Q}_{\theta}} \left[(C_{\psi^*}(x) + 2) \nabla_{\psi} C_{\psi}(x) \Big|_{\psi=\psi^*} \right] \Big|_{\theta=\theta^*} \end{aligned}$$

Now we will proceed to calculate the values for $J_{\theta\psi}$ and $J_{\theta\theta}$. To this end, we will make extensive use of the property $D_{\text{NXOR}} + D_{\text{XOR}} = 1$, which implies that all derivatives of the summation are

equal to 0: $\nabla\{D_{\text{NXOR}} + D_{\text{XOR}}\} = \nabla\{1\} = 0$. We also observe that $\forall x^*, y^* \in \text{supp}\{\mathbb{P}\}$:

$$\begin{aligned} \nabla_{xy}^2 D_{\text{NXOR}}(x, y, \psi^*) \Big|_{\substack{x=x^* \\ y=y^*}} &= \\ &= 2 \sigma(C_{\psi^*}(x^*)) \sigma(-C_{\psi^*}(x^*)) \sigma(C_{\psi^*}(y^*)) \sigma(-C_{\psi^*}(y^*)) \nabla_x C_{\psi^*}(x) \Big|_{x=x^*} \nabla_y C_{\psi^*}(y) \Big|_{y=y^*}^T \\ &= \frac{1}{8} \nabla_x C_{\psi^*}(x) \Big|_{x=x^*} \nabla_y C_{\psi^*}(y) \Big|_{y=y^*}^T \end{aligned} \quad (40)$$

So:

$$\nabla_{\psi\theta}^2 \mathcal{L}_{\text{XOR}}(\psi, \theta) \Big|_{\substack{\psi=\psi^* \\ \theta=\theta^*}} = 2\nabla_{\psi\theta}^2 \mathcal{S}(\mathbb{Q}_{\theta^*}, \mathbb{Q}_{\theta^*}, C_{\psi^*}) \Big|_{\substack{\psi=\psi^* \\ \theta=\theta^*}} + 2\nabla_{\psi\theta}^2 \mathcal{D}(\mathbb{P}, \mathbb{Q}_{\theta^*}, C_{\psi^*}) \Big|_{\substack{\psi=\psi^* \\ \theta=\theta^*}} \quad (41a)$$

$$\begin{aligned} &= 2 \left\{ 2 \mathbb{E}_{\substack{x \sim \mathbb{P} \\ z \sim Z}} \left[\left(\nabla_{\psi y}^2 D_{\text{NXOR}}(x, y, \psi) \Big|_{\substack{\psi=\psi^* \\ y=G_{\theta^*}(z)}} + \nabla_{\psi y}^2 D_{\text{XOR}}(x, y, \psi) \Big|_{\substack{\psi=\psi^* \\ y=G_{\theta^*}(z)}} \right) \nabla_{\theta}^T G_{\theta}(z) \Big|_{\theta=\theta^*} \right] \right\} \\ &= 2 \left\{ 2 \mathbb{E}_{\substack{x \sim \mathbb{P} \\ z \sim Z}} \left[\nabla_{\psi\theta}^2 1 \quad \nabla_{\theta}^T G_{\theta}(z) \Big|_{\theta=\theta^*} \right] \right\} = 0 \end{aligned}$$

$$\nabla_{\theta\psi}^2 \mathcal{L}_{\text{XOR}}(\psi, \theta) \Big|_{\substack{\psi=\psi^* \\ \theta=\theta^*}} = \nabla_{\psi\theta}^2 \mathcal{L}_{\text{XOR}}(\psi, \theta) \Big|_{\substack{\psi=\psi^* \\ \theta=\theta^*}} = 0 \quad (41b)$$

$$\nabla_{\theta\theta}^2 \mathcal{L}_{\text{XOR}}(\psi, \theta) \Big|_{\substack{\psi=\psi^* \\ \theta=\theta^*}} = \quad (41c)$$

$$\begin{aligned} &= 2\nabla_{\theta_1\theta_2}^2 \mathcal{S}(\mathbb{Q}_{\theta_1}, \mathbb{Q}_{\theta_2}, C_{\psi^*}) \Big|_{\substack{\theta_1=\theta^* \\ \theta_2=\theta^*}} + 2\nabla_{\theta\theta}^2 \mathcal{S}(\mathbb{Q}_{\theta^*}, \mathbb{Q}_{\theta^*}, C_{\psi^*}) \Big|_{\theta=\theta^*} + 2\nabla_{\theta\theta}^2 \mathcal{D}(\mathbb{P}, \mathbb{Q}_{\theta^*}, C_{\psi^*}) \Big|_{\theta=\theta^*} \\ &= 4 \mathbb{E}_{\substack{z_1 \sim Z \\ z_2 \sim Z}} \left[\nabla_{\theta} G_{\theta}(z_1) \Big|_{\theta=\theta^*} \nabla_{xy}^2 D_{\text{NXOR}}(x, y, \psi^*) \Big|_{\substack{x=G_{\theta^*}(z_1) \\ y=G_{\theta^*}(z_2)}} \nabla_{\theta} G_{\theta}(z_2) \Big|_{\theta=\theta^*}^T \right] \\ &= \frac{4}{8} \mathbb{E}_{\substack{z_1 \sim Z \\ z_2 \sim Z}} \left[\nabla_{\theta} G_{\theta}(z_1) \Big|_{\theta=\theta^*} \nabla_x C_{\psi^*}(x) \Big|_{x=G_{\theta^*}(z_1)} \nabla_x C_{\psi^*}(x) \Big|_{x=G_{\theta^*}(z_2)}^T \nabla_{\theta} G_{\theta}(z_2) \Big|_{\theta=\theta^*}^T \right] \\ &= \frac{1}{2} \mathbb{E}_{z_1 \sim Z} \left[\nabla_{\theta} C_{\psi^*}(G_{\theta}(z_1)) \Big|_{\theta=\theta^*} \right] \mathbb{E}_{z_2 \sim Z} \left[\nabla_{\theta} C_{\psi^*}(G_{\theta}(z_2)) \Big|_{\theta=\theta^*}^T \right] \\ &= \frac{1}{2} \nabla_{\theta} \mathbb{E}_{x \sim \mathbb{Q}_{\theta}} [C_{\psi^*}(x)] \Big|_{\theta=\theta^*} \nabla_{\theta} \mathbb{E}_{x \sim \mathbb{Q}_{\theta}} [C_{\psi^*}(x)] \Big|_{\theta=\theta^*}^T \end{aligned}$$

□

In order to proceed with our analysis at this point, we will state some assumptions. Specifically, we are going to define the following *reparameterization* manifolds:

$$\mathcal{M}_D := \{\psi \mid C_{\psi}(x) = 0 \quad \forall x \in \text{supp}\{\mathbb{P}\}\} \quad (42)$$

$$\mathcal{M}_G := \left\{ \theta \mid \arg \min_{\theta} \left(\mathbb{E}_{x \sim \mathbb{Q}_{\theta}} C_{\psi^*}(x) \right)^2 \right\} \quad (43)$$

We observe that \mathcal{M}_D (42) simply describes ψ^* part out of the SP of non zero-sum XORGAN, according to proposition 5.3. In addition, we assume that for each (ψ^*, θ^*) there exist ϵ -balls, $B_{\epsilon}(\psi^*)$ and $B_{\epsilon}(\theta^*)$, around ψ^* and θ^* at their respective subspaces, such that $\mathcal{M}_D \cap B_{\epsilon}(\psi^*)$ and $\mathcal{M}_G \cap B_{\epsilon}(\theta^*)$ define C^1 -manifolds. Finally, we can express \mathcal{M}_D equivalently as:

$$\mathcal{M}_D := \left\{ \psi \mid \arg \min_{\psi} \mathbb{E}_{x \sim \mathbb{P}} |C_{\psi}(x)|^2 \right\} \quad (44)$$

To understand the equivalence, we notice that eq. (44) describes a condition which is minimized if and only if it is equal to 0, as it is non-negative. This condition is equal to 0, when all (non-negative) square terms in the integral are also equal to 0. The integral is evaluated on $\text{supp}\{\mathbb{P}\}$, thus, for all points in the support, the square terms, and consequently the critic C_ψ , must be equal to 0.

Lemma 5.5 (Condition for negative definite $J_{\psi\psi}$)

If vector $u \neq 0$ does not lie in the tangent space of \mathcal{M}_D (44) at ψ^* , $\mathcal{T}_{\psi^*}\mathcal{M}_D$, then $u^T J_{\psi\psi} u < 0$.

Proof: From lemma 5.4 we have

$$u^T J_{\psi\psi} u = -\frac{1}{2} \mathbb{E}_{x \sim \mathbb{P}} \left[\left(u^T \nabla_\psi C_\psi(x) \Big|_{\psi=\psi^*} \right)^2 \right] \quad (45)$$

which implies $u^T J_{\psi\psi} u \leq 0$. We get equality if and only if:

$$u^T \nabla_\psi C_\psi(x) \Big|_{\psi=\psi^*} = 0 \quad \forall x \in \text{supp}\{\mathbb{P}\}. \quad (46)$$

Let:

$$h(\psi) := \mathbb{E}_{x \sim \mathbb{P}} \left[|C_\psi(x)|^2 \right] \quad (47)$$

$$\implies \nabla_\psi h(\psi) = 2 \mathbb{E}_{x \sim \mathbb{P}} \left[C_\psi(x) \nabla_\psi C_\psi(x) \right] \quad (48)$$

$$\implies \nabla_{\psi\psi}^2 h(\psi) = 2 \mathbb{E}_{x \sim \mathbb{P}} \left[\nabla_\psi C_\psi(x) \nabla_\psi C_\psi(x)^T \right] + 2 \mathbb{E}_{x \sim \mathbb{P}} \left[C_\psi(x) \nabla_{\psi\psi}^2 C_\psi(x) \right] \quad (49)$$

By using the expression of ψ^* from proposition 5.3, we observe that:

$$h(\psi^*) = 0 \quad (50)$$

$$\nabla_\psi h(\psi) \Big|_{\psi=\psi^*} = 0 \quad (51)$$

$$\nabla_{\psi\psi}^2 h(\psi) \Big|_{\psi=\psi^*} = 2 \mathbb{E}_{x \sim \mathbb{P}} \left[\nabla_\psi C_\psi(x) \Big|_{\psi=\psi^*} \nabla_\psi C_\psi(x) \Big|_{\psi=\psi^*}^T \right] \geq 0 \quad (52)$$

so it achieves the minimum of the expression h and thus $\psi^* \in \mathcal{M}_D$. The vector $u \in \mathcal{T}_{\psi^*}\mathcal{M}_D$ if and only if the second directional derivative at $\psi = \psi^*$ is equal to 0, which is iff eq. (46) holds. \square

Lemma 5.6 (Condition for negative definite $J_{\theta\theta}$)

If vector $w \neq 0$ does not lie in the tangent space of \mathcal{M}_G (43) at θ^* , $\mathcal{T}_{\theta^*}\mathcal{M}_G$, then $w^T J_{\theta\theta} w < 0$.

Proof: From lemma 5.4 we have

$$w^T J_{\theta\theta} w = -\frac{1}{2} \left(w^T \nabla_\theta \mathbb{E}_{x \sim \mathbb{Q}_\theta} [C_{\psi^*}(x)] \Big|_{\theta=\theta^*} \right)^2 \quad (53)$$

which implies $w^T J_{\theta\theta} w \leq 0$. We get equality if and only if:

$$w^T \nabla_\theta \mathbb{E}_{x \sim \mathbb{Q}_\theta} [C_{\psi^*}(x)] \Big|_{\theta=\theta^*} = 0 \quad (54)$$

Let:

$$g(\theta) := \left(\mathbb{E}_{x \sim \mathbb{Q}_\theta} C_{\psi^*}(x) \right)^2 \quad (55)$$

$$\implies \nabla_\theta g(\theta) = 2 \left(\mathbb{E}_{x \sim \mathbb{Q}_\theta} C_{\psi^*}(x) \right) \nabla_\theta \mathbb{E}_{x \sim \mathbb{Q}_\theta} [C_{\psi^*}(x)] \quad (56)$$

$$\implies \nabla_{\theta\theta}^2 g(\theta) = 2 \nabla_\theta \mathbb{E}_{x \sim \mathbb{Q}_\theta} [C_{\psi^*}(x)] \nabla_\theta \mathbb{E}_{x \sim \mathbb{Q}_\theta} [C_{\psi^*}(x)]^T + 2 \left(\mathbb{E}_{x \sim \mathbb{Q}_\theta} C_{\psi^*}(x) \right) \nabla_{\theta\theta}^2 \mathbb{E}_{x \sim \mathbb{Q}_\theta} [C_{\psi^*}(x)] \quad (57)$$

By using proposition 5.3, we observe that:

$$g(\theta^*) = 0 \quad (58)$$

$$\nabla_{\theta} g(\theta) \Big|_{\theta=\theta^*} = 0 \quad (59)$$

$$\nabla_{\theta\theta}^2 g(\theta) \Big|_{\theta=\theta^*} = 2\nabla_{\theta} \mathbb{E}_{x \sim \mathbb{Q}_{\theta}} [C_{\psi^*}(x)] \Big|_{\theta=\theta^*} \nabla_{\theta} \mathbb{E}_{x \sim \mathbb{Q}_{\theta}} [C_{\psi^*}(x)] \Big|_{\theta=\theta^*}^T \geq 0 \quad (60)$$

so it achieves the minimum of the expression g and thus $\theta^* \in \mathcal{M}_G$. The vector $w \in \mathcal{T}_{\theta^*} \mathcal{M}_G$ if and only if the second directional derivative at $\theta = \theta^*$ is equal to 0, which is iff eq. (54) holds. \square

We will attempt to draw some conclusion about the eigenvalues of Jacobian $J(\psi^*, \theta^*)$. For this purpose we prove the following lemma.

Lemma 5.7 (Eigenvalues of block upper triangular matrix)

Let matrix

$$J := \begin{pmatrix} -A & B \\ 0 & -D \end{pmatrix} \quad (61)$$

where $A \in \mathbb{R}^{n \times n}$ $D \in \mathbb{R}^{m \times m}$. Then, $\lambda\{J\} = \lambda\{-A\} \cup \lambda\{-D\}$.

Proof: Initially we will show that $\lambda\{J\} = \lambda\{-A\} \cup \lambda\{-D\}$. Let $u^T = (x^T, y^T)$ be an eigenvector of matrix J , with $x \in \mathbb{R}^n$ and $y \in \mathbb{R}^m$, and the corresponding eigenvalue λ :

$$Ju = \lambda u \iff \begin{pmatrix} -Ax + By \\ -Dy \end{pmatrix} = \begin{pmatrix} \lambda x \\ \lambda y \end{pmatrix} \quad (62)$$

If $y \neq 0$, then from $-Dy = \lambda y$ we understand that y consists an eigenvector of $-D$ with eigenvalue λ . Thus, $\lambda \in \lambda\{-D\}$. Otherwise $y = 0$, and then $-Ax = \lambda x$ will hold, from which we understand that x consists an eigenvector of $-A$ with eigenvalue λ . In total, by combining the two possibilities, it ought to hold that $\lambda\{J\} \subseteq \lambda\{-A\} \cup \lambda\{-D\}$.

In inverse, let λ be an eigenvalue of $-A$ with corresponding eigenvector $x \neq 0$. Then $-Ax = \lambda x$ and we observe that the vector $(x^T, 0)^T$ is an eigenvector of J , because

$$J \begin{pmatrix} x \\ 0 \end{pmatrix} = \begin{pmatrix} -Ax \\ 0 \end{pmatrix} = \lambda \begin{pmatrix} x \\ 0 \end{pmatrix} \quad (63)$$

and with λ as its eigenvalue. Thus $\lambda\{-A\} \subseteq \lambda\{J\}$. Finally, we consider λ to be an eigenvalue of $-D$, which is not also an eigenvalue of $-A$, so $\lambda \in \lambda\{-D\} - \lambda\{-A\}$, and its corresponding eigenvector $y \neq 0$. Then $-Dy = \lambda y$, but also the matrix $-A - \lambda I$ is not singular, and hence invertible. Then we observe that the vector $(x^T, y^T)^T$, with $x := (A + \lambda I)^{-1} By$, consists an eigenvector of J with eigenvalue λ .

$$J \begin{pmatrix} x \\ y \end{pmatrix} = \begin{pmatrix} -A(A + \lambda I)^{-1} By + By \\ -Dy \end{pmatrix} = \begin{pmatrix} [-A + (A + \lambda I)] (A + \lambda I)^{-1} By \\ -Dy \end{pmatrix} = \lambda \begin{pmatrix} x \\ y \end{pmatrix} \quad (64)$$

Consequently, from all possibilities of the inverse, we have $\lambda\{-A\} \cup (\lambda\{-D\} - \lambda\{-A\}) \subseteq \lambda\{J\}$.

By combining the direct and the inverse, we have: $\lambda\{J\} = \lambda\{-A\} \cup \lambda\{-D\}$. \square

Corollary 5.7.1 (Eigenvalues of Jacobian at SP (29) of non zero-sum XORGAN)

Matrix J , found at lemma 5.4, is negative semi-definite in the parameter space $(\psi^T, \theta^T)^T$. However, in the subspace $V(\psi^*, \theta^*) := \text{span}\{(u^T, w^T)^T \mid u \notin \mathcal{T}_{\psi^*} \mathcal{M}_D, w \notin \mathcal{T}_{\theta^*} \mathcal{M}_G\}$, J is negative definite. In addition, all of its eigenvalues are real numbers.

Proof: From lemma 5.4, we see that the symmetric matrices $J_{\psi\psi}$ and $J_{\theta\theta}$ are negative semi-definite. In consequence, lemma 5.7 implies that J is negative semi-definite in the joint parameter space. In subspace V though, $J_{\psi\psi}$ and $J_{\theta\theta}$ are negative definite, according to lemmata 5.5 and 5.6. Thus, J has some eigenvalues with negative real part. Furthermore, all eigenvalues are real numbers because $J_{\psi\psi}$ and $J_{\theta\theta}$ are symmetric matrices \square

Stability analysis will be continued in future work with purpose to draw guarantees of asymptotic convergence for the dynamical system at the set of parameters $\mathcal{M}_D \times \mathcal{M}_G$, as those were defined at eqs. (46) and (54). Such analysis will be done by adapting properly the methodology utilized by [24, 7]. Specifically, we will first reparameterize the linearized system with respect to the tangent and the co-tangent subspaces, and second we will use Lyapunov argument, as well as perhaps Lasalle’s invariance principle, in order to prove some local stability properties at the set $\mathcal{M}_D \times \mathcal{M}_G$.

5.3 Simulation on Toy Problems

We are performing simulation on toy problem settings as proof of concept. We are generating two training sets on which we are going to train our models. The first one is composed by samples from a mixture of 8 2D gaussian distributions, whose means form a regular octagon, while the second one by samples from a mixture of 25 2D gaussians, whose means form a grid. Our generator and critic models are feed-forward networks, and spectral normalization [8] is used at the critic model. Training is performed by optimizing one gradient step at a time, alternatively, for each of XORGAN’s objectives. Adam [26] is the optimization algorithm used in these simulations, with $\beta_1 = 0.5$, $\beta_2 = 0.9$, and a fixed learning rate of $\eta = 1e-4$. Finally, we keep an exponential moving average of our generator model’s weights with $\beta = 0.999$. We use the averaged weights for the generator model in our visualizations.

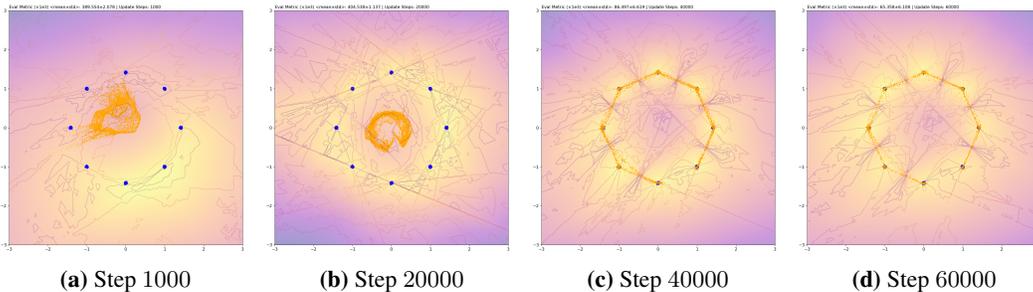


Figure 2: Training evolution of XORGAN on 8 gaussians synthetic dataset

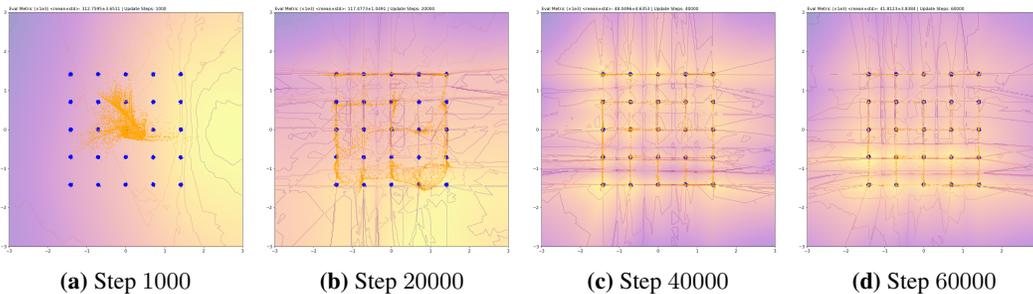


Figure 3: Training evolution of XORGAN on 25 gaussians synthetic dataset