



# A NEURAL TANGENT KERNEL PERSPECTIVE OF GANS

GdR ISIS – Théorie du deep learning

June 28th, 2021

Jean-Yves Franceschi,<sup>1</sup> Emmanuel de Bézenac,<sup>1</sup>  
Ibrahim Ayed,<sup>1,2</sup> Mickaël Chen,<sup>3</sup> Sylvain Lamprier,<sup>1</sup>  
Patrick Gallinari<sup>1,4</sup>

<sup>1</sup>Sorbonne Université, CNRS, LIP6, F-75005 Paris, France

<sup>2</sup>ThereSIS Lab, Thales, Palaiseau, France    <sup>3</sup>Valeo.ai, Paris, France

<sup>4</sup>Criteo AI Lab, Paris, France



- | Novel general framework of analysis of GANs taking into account discriminator architecture, unlike prior work.
- | Based on the theory of infinite-width neural networks.
- | Closes theoretical gaps to better understand GAN training.
- | Allows experimental studies of various GAN settings.
- | Highlights the interesting properties of neural tangent kernels for generative modeling.

## Preprint

<https://arxiv.org/abs/2106.05566>

## Principle

Generative model with a generator  $g$  competing with a discriminator  $f$  trained to distinguish between fake ( $\alpha_g$ ) and target ( $\beta$ ) samples.

## Principle

Generative model with a generator  $g$  competing with a discriminator  $f$  trained to distinguish between fake ( $\alpha_g$ ) and target ( $\beta$ ) samples.

Examples:

- | Vanilla GAN (Goodfellow et al., 2014):

$$\min_g \max_f \mathbb{E}_{x \sim \alpha_g} \left[ \log \left( 1 - \sigma(f(x)) \right) \right] + \mathbb{E}_{y \sim \beta} \left[ \log \sigma(f(y)) \right].$$

## Principle

Generative model with a generator  $g$  competing with a discriminator  $f$  trained to distinguish between fake ( $\alpha_g$ ) and target ( $\beta$ ) samples.

Examples:

- | Vanilla GAN (Goodfellow et al., 2014):

$$\min_g \max_f \mathbb{E}_{x \sim \alpha_g} \left[ \log \left( 1 - \sigma(f(x)) \right) \right] + \mathbb{E}_{y \sim \beta} \left[ \log \sigma(f(y)) \right].$$

- | IPM, e.g. WGAN (Arjovsky et al., 2017):

$$\min_g \max_f \mathbb{E}_{x \sim \alpha_g} [f(x)] - \mathbb{E}_{y \sim \beta} [f(y)].$$

## Setting

Generator  $g \in G: \mathbb{R}^d \rightarrow \mathbb{R}^n$  optimizing:

$$\inf_{g \in G} \left\{ C_{f_{\alpha_g}^*}(\alpha_g), \mathbb{E}_{x \sim \alpha_g} \left[ (c \circ f_{\alpha_g}^*)(x) \right] \right\},$$

with the discriminator  $f_{\alpha_g}^* \in F: \mathbb{R}^n \rightarrow \mathbb{R}$  optimizing:

$$\sup_{f \in F} \left\{ L_{\alpha_g}(f), \mathbb{E}_{x \sim \alpha_g} [(a \circ f)(x)] - \mathbb{E}_{y \sim \beta} [(b \circ f)(y)] \right\}.$$

Includes many GAN formulations, e.g. WGAN, with  $a = b = c = \text{id}$ :

$$\min_g \max_f \mathbb{E}_{x \sim \alpha_g} [f(x)] - \mathbb{E}_{y \sim \beta} [f(y)].$$

I  $G$  and  $F$  are families of neural networks.

- |  $G$  and  $F$  are families of neural networks.
- | Generator  $g = g_\theta$  optimization:

$$\theta_{j+1} = \theta_j - \eta \mathbb{E}_{z \sim p_z} \left[ r_{\theta} g_{\theta_j}(z)^T r_x \left( c - f_{\alpha_g}^*(x) \right) \Big|_{x=g_{\theta_j}(z)} \right].$$

$r_{\alpha_g}^*$  is crucial for generator training.



- |  $G$  and  $F$  are families of neural networks.
- | Generator  $g = g_\theta$  optimization:

$$\theta_{j+1} = \theta_j - \eta \mathbb{E}_{z \sim p_z} \left[ r_{\theta g_{\theta_j}(z)}^T r_x \left( c - f_{\alpha_g}^* \right) (x) \Big|_{x=g_{\theta_j}(z)} \right].$$

$r_{f_{\alpha_g}^*}$  is crucial for generator training.

- | The discriminator is trained at each step using a finite number of samples from  $\alpha$  and  $\beta$ .

The training distribution  $\hat{\gamma} = \frac{\alpha + \beta}{2}$  is a finite mixture of Diracs.

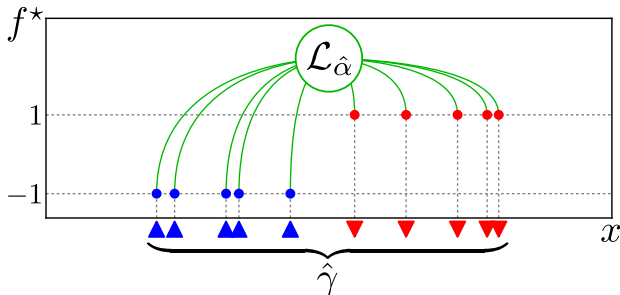
| Analysis by Arjovsky et al. (2017) of LSGAN (Mao et al., 2017):

$$L_{\hat{\alpha}}(f) = \mathbb{E}_{x \sim \hat{\alpha}} \left[ (f(x) + 1)^2 \right] - \mathbb{E}_{y \sim \hat{\beta}} \left[ (f(y) - 1)^2 \right].$$

- I Analysis by Arjovsky et al. (2017) of LSGAN (Mao et al., 2017):

$$L_{\hat{\alpha}}(f) = \mathbb{E}_{x \sim \hat{\alpha}} \left[ (f(x) + 1)^2 \right] - \mathbb{E}_{y \sim \beta} \left[ (f(y) - 1)^2 \right].$$

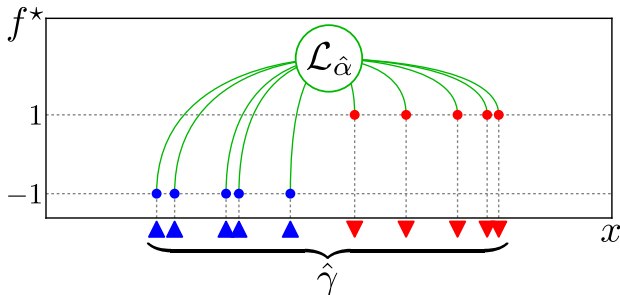
- I Considering  $f \in L^2(\hat{\gamma})$ , the optimal discriminator  $f^*$  has values independent from the generated distribution.



- Analysis by Arjovsky et al. (2017) of LSGAN (Mao et al., 2017):

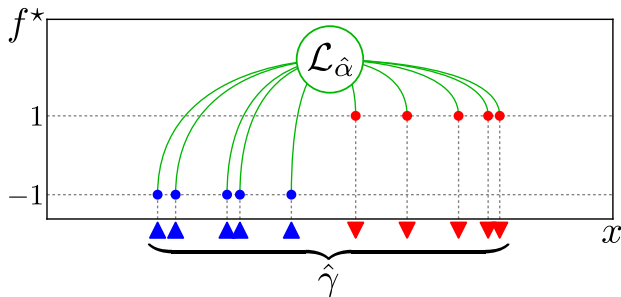
$$L_{\hat{\alpha}}(f) = \mathbb{E}_{x \sim \hat{\alpha}} \left[ (f(x) + 1)^2 \right] - \mathbb{E}_{y \sim \beta} \left[ (f(y) - 1)^2 \right].$$

- Considering  $f \in L^2(\hat{\gamma})$ , the optimal discriminator  $f^*$  has values independent from the generated distribution.
- Conclusion: the generator loss is constant and the generator receives null gradients.



## Problem

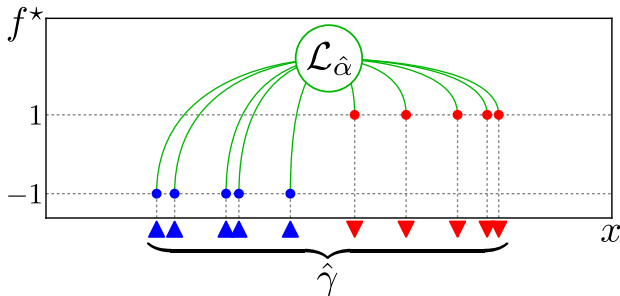
In practice, LSGAN and vanilla GAN are trainable.



## Problem

In practice, LSGAN and vanilla GAN are trainable.

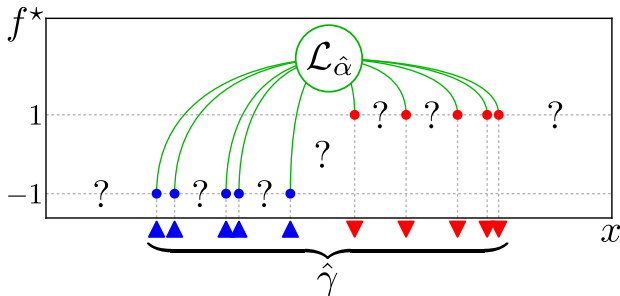
- When updating the generator, the discriminator is assumed fixed. What about  $r f^*$ ?



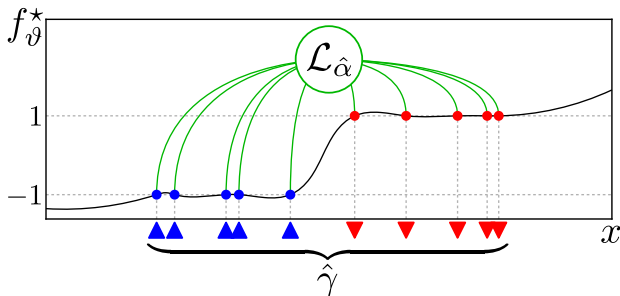
## Problem

In practice, LSGAN and vanilla GAN are trainable.

- | When updating the generator, the discriminator is assumed fixed. What about  $r f^*$ ?
- | Considering  $f \in L^2(\hat{\gamma})$ ,  $r f^*$  is undefined as  $f^*$  is only defined over  $\text{supp } \hat{\gamma}$ .



- | Considering  $f = f_\vartheta$  as a neural network of a given architecture trained by gradient ascent to maximize  $\mathcal{L}_{\hat{\alpha}}$ .
- | With this modelization, the discriminator is defined everywhere.
- | In our framework, we model this parameterization in the infinite-width regime using neural tangent kernels (NTKs).





- | Generally applicable proxy representing the learning capacity of the architecture of  $f$ .
- | Makes neural networks act like kernel methods and simplifies their theoretical study.

- | Generally applicable proxy representing the learning capacity of the architecture of  $f$ .
- | Makes neural networks act like kernel methods and simplifies their theoretical study.

## NTK of $f_{\vartheta}$ (Jacot et al., 2018)

$$k_{f_{\vartheta}}(x, y) = \langle \partial_{\vartheta} f_{\vartheta}(x), \partial_{\vartheta} f_{\vartheta}(y) \rangle.$$

- | Generally applicable proxy representing the learning capacity of the architecture of  $f$ .
- | Makes neural networks act like kernel methods and simplifies their theoretical study.

## NTK of $f_{\vartheta}$ (Jacot et al., 2018)

$$k_{f_{\vartheta}}(x, y) = \langle \partial_{\vartheta} f_{\vartheta}(x), \partial_{\vartheta} f_{\vartheta}(y) \rangle.$$

## Constant NTK (Jacot et al., 2018)

In the infinite-width limit of  $f$ , during training,  $k_{f_{\vartheta}}(x, y) = k(x, y)$ .

## Kernel Integral Operator and RKHS

$$T_{k, \hat{\gamma}}: L^2(\hat{\gamma}) \rightarrow H_k^{\hat{\gamma}}, \quad h \mapsto \int_x k(\cdot, x)h(x) d\hat{\gamma}(x),$$

where  $H_k^{\hat{\gamma}} \subset L^2(\hat{\gamma})$  is the RKHS of  $k$  generated by  $\hat{\gamma}$ .

## Discriminator Training

$$\partial_t f_t = T_{k, \hat{\gamma}} \left( r^{\hat{\gamma}} L_{\hat{\alpha}}(f_t) \right).$$

## Kernel Integral Operator and RKHS

$$T_{k, \hat{\gamma}}: L^2(\hat{\gamma}) \rightarrow H_k^{\hat{\gamma}}, \quad h \mapsto \int_x k(\cdot, x)h(x) d\hat{\gamma}(x),$$

where  $H_k^{\hat{\gamma}} \subset L^2(\hat{\gamma})$  is the RKHS of  $k$  generated by  $\hat{\gamma}$ .

## Discriminator Training

$$\partial_t f_t = T_{k, \hat{\gamma}} \left( r^{\hat{\gamma}} L_{\hat{\alpha}}(f_t) \right).$$

In short:

- | Simple expression of the evolution of  $f$  during training.
- | Confers strong regularity to  $f$  (corresponding to its parameterization as a neural network).

The following results hold under mild assumptions on the architecture of  $f$  and  $L_{\hat{\alpha}}$ .

## Theorem

$f_t$  is uniquely defined from the previous ODE, and:

$$\forall t \in \mathbb{R}_+, f_t = f_0 + T_{k, \hat{\gamma}} \left( \int_0^t r^{\hat{\gamma}} L_{\hat{\alpha}}(f_s) ds \right) \in f_0 + H_k^{\hat{\gamma}}.$$

The following results hold under mild assumptions on the architecture of  $f$  and  $L_{\hat{\alpha}}$ .

## Theorem

$f_t$  is uniquely defined from the previous ODE, and:

$$\forall t \in \mathbb{R}_+, f_t = f_0 + T_{k, \hat{\gamma}} \left( \int_0^t r^{\hat{\gamma}} L_{\hat{\alpha}}(f_s) ds \right) \in f_0 + H_k^{\hat{\gamma}}.$$

## Theorem

The NTK  $k$  of  $f$  and  $f_t$  are infinitely differentiable (almost) everywhere.

$$L_{\hat{\alpha}}(f) = E_{x \sim \hat{\alpha}}[f(x)] - E_{y \sim \hat{\beta}}[f(y)]$$



$$L_{\hat{\alpha}}(f) = \mathbb{E}_{x \sim \hat{\alpha}}[f(x)] - \mathbb{E}_{y \sim \hat{\beta}}[f(y)]$$

## IPM Discriminator and MMD

In the infinite-width regime:

$$\begin{aligned} f_t &= f_0 + t f_{\hat{\alpha}}, \\ f_{\hat{\alpha}} &= \mathbb{E}_{x \sim \hat{\alpha}}[k(x, \cdot)] - \mathbb{E}_{y \sim \hat{\beta}}[k(y, \cdot)], \\ L_{\hat{\alpha}}(f_t) &= L_{\hat{\alpha}}(f_0) + t \text{MMD}_k^2(\hat{\alpha}, \hat{\beta}). \end{aligned}$$

$$L_{\hat{\alpha}}(f) = \mathbb{E}_{x \sim \hat{\alpha}}[f(x)] - \mathbb{E}_{y \sim \hat{\beta}}[f(y)]$$

## IPM Discriminator and MMD

In the infinite-width regime:

$$\begin{aligned} f_t &= f_0 + t f_{\hat{\alpha}}, \\ f_{\hat{\alpha}} &= \mathbb{E}_{x \sim \hat{\alpha}}[k(x, \cdot)] - \mathbb{E}_{y \sim \hat{\beta}}[k(y, \cdot)], \\ L_{\hat{\alpha}}(f_t) &= L_{\hat{\alpha}}(f_0) + t \text{MMD}_k^2(\hat{\alpha}, \hat{\beta}). \end{aligned}$$

## Consequence

The generator minimizes the MMD, defined with the NTK of  $f$ , between fake and target distributions (thanks to Arbel et al. (2019)).

$$L_{\hat{\alpha}}(f) = \mathbb{E}_{x \sim \hat{\alpha}} \left[ (f(x) + 1)^2 \right] - \mathbb{E}_{y \sim \hat{\beta}} \left[ (f(y) - 1)^2 \right]$$

## LSGAN Discriminator

In the infinite-width regime:

$$f_t = \exp(-4tT_{k,\hat{\gamma}})(f_0 - \rho) + \rho, \quad \rho(x) = \begin{cases} 1 & \text{if } x \in \text{supp } \hat{\alpha} \\ -1 & \text{if } x \in \text{supp } \hat{\beta} \end{cases}.$$

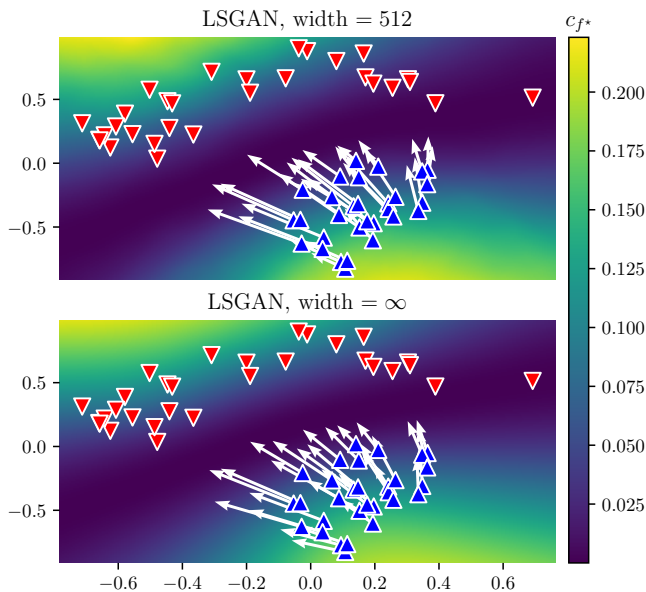
## Consequence

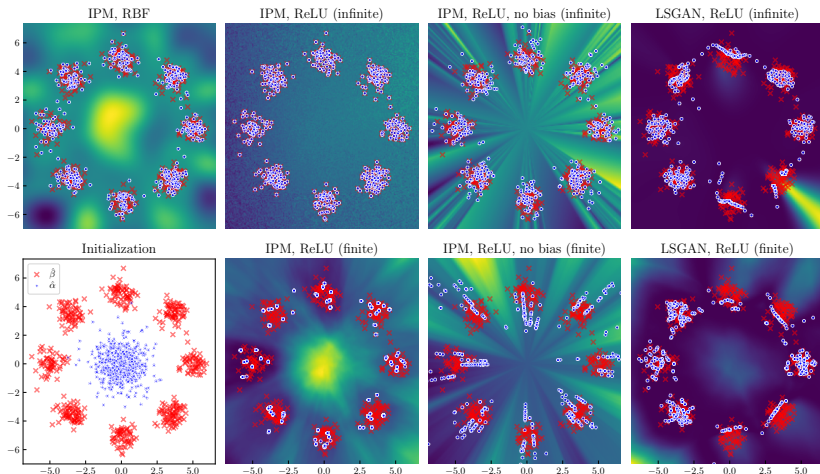
$f_1$  is optimal when  $k$  is strictly positive definite ( $f_1|_{\text{supp } \hat{\gamma}} = \rho$ ) but  $T_{k,\hat{\gamma}}$  smoothes  $\rho$  on the whole domain.

## Code

<https://github.com/emited/gantk2>

- | Various finite- and infinite-width implementations of GANs.
- | Based on the Neural Tangents library (Novak et al., 2020).
- | Could be used for further experimental studies.





- Arbel, Michael et al. (2019). “Maximum Mean Discrepancy Gradient Flow”. In: *Advances in Neural Information Processing Systems*. Ed. by Hanna Wallach et al. Vol. 32. Curran Associates, Inc., pp. 6481–6491.
- Arjovsky, Martin, Soumith Chintala, and Léon Bottou (Aug. 2017). “Wasserstein Generative Adversarial Networks”. In: *Proceedings of the 34th International Conference on Machine Learning*. Ed. by Doina Precup and Yee Whye Teh. Vol. 70. Proceedings of Machine Learning Research. PMLR, pp. 214–223.
- Goodfellow, Ian et al. (2014). “Generative Adversarial Nets”. In: *Advances in Neural Information Processing Systems*. Ed. by Zoubin Ghahramani et al. Vol. 27. Curran Associates, Inc., pp. 2672–2680.
- Jacot, Arthur, Franck Gabriel, and Clement Hongler (2018). “Neural Tangent Kernel: Convergence and Generalization in Neural Networks”. In: *Advances in Neural Information Processing Systems*. Ed. by Samy Bengio et al. Vol. 31. Curran Associates, Inc., pp. 8580–8589.
- Mao, Xudong et al. (Oct. 2017). “Least Squares Generative Adversarial Networks”. In: *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*, pp. 2813–2821.
- Novak, Roman et al. (2020). “Neural Tangents: Fast and Easy Infinite Neural Networks in Python”. In: *International Conference on Learning Representations*.