

# The analysis of GANs from mathematics perspective

model part

2017-7-14

# Table

- 1 GANs model
  - Perspective of Analysis
  - Perspective of Statistics
  - Perspective of Info. Theory

# Generative Adversarial Networks<sup>1</sup>

Generator:  $G$ , Parameterization  $\theta_g$ , The true distribution of data  $P_r(PDF : p_r(x))$ , Generated distribution  $P_g(PDF : p_g(x))$

Discriminator  $D$ , Parameterization  $\theta_d$ , discriminate the probability of a data comes from  $P_g$

The objective function is a classical saddle point problem

$$\min_G \max_D V(D, G) = \mathbb{E}_{x \sim P_r} [\log D(x)] + \mathbb{E}_{x \sim P_g} [\log(1 - D(x))] \quad (1)$$

using variational method to  $D$ , it is not hard to get  $D^*(x) = \frac{p_r(x)}{p_g(x) + p_r(x)}$ , substitute it into the objective function, we have:

$$\min 2JSD(P_r || P_g) - 2 \log 2 \quad (2)$$

**FACT:** The optimization of the primal problem is equivalent to the optimization of the JS divergence between  $P_r$  and  $P_g$

<sup>1</sup> Ian J. Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, NIPS, 2014

## The Generalization of JS divergence<sup>2</sup>

The methods to measure the distance between two distributions are various, some work focus on the generalization of the JSD.

f-divergence family(Ali-Silvey distances):

$$D_f(P_r||P_g) = \int_x p_g(x) f\left(\frac{p_r(x)}{p_g(x)}\right) dx \quad (3)$$

in which, generative function  $f: \mathbb{R}_+ \rightarrow \mathbb{R}$  is convex, lower-semicontinuous, with  $f(1) = 0$

Based on the powerful theorem that if  $f$  is convex and lower semicontinuous, then  $f^{**} = f$ , it can be rewritten as

$$f(u) = \sup_{t \in \text{dom}_{f^*}} \{tu - f^*(t)\} \quad (4)$$

Thus, the lower bound for the divergence of  $P_r$  and  $P_g$  comes as:

$$D_f(P_r||P_g) \geq \sup_{T \in \mathcal{T}} (\mathbb{E}_{x \sim P_r}[T(x)] - \mathbb{E}_{x \sim P_g}[f^*(T(x))]) \quad (5)$$

---

<sup>2</sup>f-GAN: Training Generative Neural Samplers using Variational Divergence Minimization.

Still by variational method

$$T^*(x) = f' \left( \frac{p_r(x)}{p_g(x)} \right) \quad (6)$$

The primal objective function can be rewritten as:

$$\min_{\theta_g} \max_w F(\theta_g, w) = \min_{\theta_g} \max_w \mathbb{E}_{x \sim P_r} [T_w(x)] - \mathbb{E}_{x \sim P_g} [f^*(T_w(x))] \quad (7)$$

Name	$D_f(P\ Q)$	Generator $f(u)$	$T^*(x)$
Kullback-Leibler	$\int p(x) \log \frac{p(x)}{q(x)} dx$	$u \log u$	$1 + \log \frac{p(x)}{q(x)}$
Reverse KL	$\int q(x) \log \frac{q(x)}{p(x)} dx$	$-\log u$	$-\frac{q(x)}{p(x)}$
Pearson $\chi^2$	$\int \frac{(q(x)-p(x))^2}{p(x)} dx$	$(u-1)^2$	$2(\frac{p(x)}{q(x)} - 1)$
Squared Hellinger	$\int (\sqrt{p(x)} - \sqrt{q(x)})^2 dx$	$(\sqrt{u}-1)^2$	$(\sqrt{\frac{p(x)}{q(x)}} - 1) \cdot \sqrt{\frac{q(x)}{p(x)}}$
Jensen-Shannon	$\frac{1}{2} \int p(x) \log \frac{2p(x)}{p(x)+q(x)} + q(x) \log \frac{2q(x)}{p(x)+q(x)} dx$	$-(u+1) \log \frac{1+u}{2} + u \log u$	$\log \frac{2p(x)}{p(x)+q(x)}$
GAN	$\int p(x) \log \frac{2p(x)}{p(x)+q(x)} + q(x) \log \frac{2q(x)}{p(x)+q(x)} dx - \log(4)$	$u \log u - (u+1) \log(u+1)$	$\log \frac{p(x)}{p(x)+q(x)}$

it is not hard to prove that if  $f$  is taken as the last function in the chart, the objective function is exactly same as that in vanilla GAN. By the way, for  $Pearson\chi^2$  divergence, there is also a corresponding paper<sup>3</sup>

<sup>3</sup>Least square GAN

## The disadvantages of using f-divergence family

- the optimization process is unstable<sup>4</sup>
- it may lead to mode collapse<sup>5</sup>
- Generator didn't learn anything useful, just memorize some samples<sup>6</sup>

### Possible reasons

- it isn't appropriate to assume that both Generator and Discriminator has infinite modeling capacity

### FACT:

It isn't appropriate to use the JS divergence to measure the distance of two lower-dimensional distributions if both of them lie on a high-dimensional space<sup>7</sup>

---

<sup>4</sup>Improved techniques for training GANS

<sup>5</sup>Loss-Sensitive Generative Adversarial Networks on Lipschitz Densities

<sup>6</sup>Generalization and Equilibrium in Generative Adversarial Nets

<sup>7</sup>Towards Principle Methods for Training GAN



## Why JSD is not good?

### Lemma 1

True distribution  $P_r$  and generated distribution  $P_g$  lie on two lower-dimensional submanifolds of a high-dimensional space.

### Lemma 2

if  $A$  and  $B$  are two lower-dimensional, compact submanifolds,  $A$  doesn't intersect with  $B$ , then, there is always a perfect discriminator  $D$  ( smooth, continuous, and takes two different constants on  $A$  and  $B$  )that can separate  $A$  and  $B$ .

(Elegant Proof: Urysohn's lemma)

### Lemma 3: A weaker condition

As long as two lower-dimensional submanifolds doesn't perfect align (the probability for such circumstance is 1), there is always a perfect discriminator can separate these two submanifolds ( almost everywhere)

## Theorem 1

As long as these two lower-dimensional submanifolds don't perfectly align, their JS divergence always takes a constant, which is  $\log 2$ , neglecting the real distance of them.

## Theorem 2:

As long as the distribution  $P_g$  is generated by a continuous differentiable function, what's more  $P_g$ ,  $P_r$  lies on two lower-dimensional submanifolds, which don't perfectly align, then as the Discriminator converges to the perfect one, the gradient for the generator converges to 0

## FACT:

$JSD(P_r || P_{g_\theta})$  is a constant almost everywhere thus not continuous w.r.t  $P_{g_\theta}$ , the back-propagated gradient for  $G$  is close to 0

**Solution:** Find a distance measurement  $\rho$  that satisfies,

- Distribution  $P_{g_\theta}$  is continuous w.r.t  $g_\theta$
- Measurement  $\rho(P_r || P_{g_\theta})$  is continuous w.r.t  $P_{g_\theta}$



## Wasserstein Distance<sup>8</sup>

$$W(P_r, P_g) = \inf_{\gamma \in \Pi(P_r, P_g)} \mathbb{E}_{(x,y) \sim \gamma} [|x - y|] \quad (8)$$

### Properties

- ❶ if  $g_\theta$  is continuous w.r.t  $\theta$ , then  $W(P_r || P_{g_\theta})$  is continuous w.r.t  $\theta$
- ❷ if  $g_\theta$  is local Lipschitz, then  $W(P_r || P_{g_\theta})$  is continuous almost everywhere, differentiable almost everywhere.

### How weak

- TV Distance:  $\delta(P_r || P_g) = \sup_{A \in \Sigma} |p_r(A) - p_g(A)|$
- KL Divergence:  $KL(P_r || P_g) = \int \log\left(\frac{p_r(x)}{p_g(x)}\right) p_r(x) d\mu(x)$
- JS Divergence:  $JSD(P_r || P_g) = KL(P_r || \frac{1}{2}(P_r + P_g)) + KL(P_g || \frac{1}{2}(P_r + P_g))$

<sup>8</sup>Wasserstein GAN

- $\delta(P_r||P_g) \rightarrow 0 \iff JSD(P_r||P_g) \rightarrow 0$  (the two norms induced by TV divergence and JS divergence are euqivalent)
- $KL(P_g||P_r) \rightarrow 0 \implies JSD(P_r||P_g) \rightarrow 0 \implies W(P_r||P_g) \rightarrow 0$
- the topology induced by KL divergence is the strongest one, then comes JSD and TV, the weakest one is induced by Wasserstein

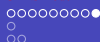
## Why it is weak

- $W(P_r||P_{g_\theta}) \rightarrow 0$  means  $P_{g_\theta}$  convergence in distribution to  $P_r$
- JSD and TV induce the strongest topology in  $C_b(\chi)^*$ , while the Wasserstein distance is the *weak*<sup>\*</sup> topology corresponding to that ◦

## How to use?

$$W(P_r||P_{g_\theta}) = \sup_{||f||_L \leq 1} \mathbb{E}_{x \sim P_r} [f(x)] - \mathbb{E}_{x \sim P_{g_\theta}} [f(x)] \quad (9)$$

Minimize  $W(P_r||P_{g_\theta})$  w.r.t  $\theta$ .



The objective function change into

$$\min_{\theta} \max_w \mathbb{E}_{x \sim P_r} [f_w(x)] - \mathbb{E}_{z \sim P(z)} [f_w(g_{\theta}(z))] \quad (10)$$

How to realize  $\|f\|_L \leq 1$ ?

a trick used is called weight clipping, just force all the parameters of  $f$  fall into  $w \in [-t, t]$ .

## Disadvantages

- This trick greatly reduce the set for Lipschitz functions.
- the experiment results show that all parameters of  $f$  are around  $t$  and  $-t$

Some subsequent work are came up with <sup>9</sup> to solve these two problems

---

<sup>9</sup>Improved Training of Wasserstein GANs

## Perspective of Statistics

### Lemma

$$P_r = P_g \iff \forall \Phi \in C, ||E_{x \sim P_r} \Phi(x) - E_{x \sim P_g} \Phi(x)||^2 = 0$$

GAN: view  $\Phi$  as the Discriminator  $D$ , the first step is to maximize the discrepancy by adjusting  $D$ :

$$\max_w ||E_{x \sim P_r} \Phi_w(x) - E_{x \sim P_g} \Phi_w(x)||^2 \quad (11)$$

then minimize the maximum w.r.t Generator:

$$\min_{\theta_g} \max_w ||E_{x \sim P_r} \Phi_w(x) - E_{x \sim P_g} \Phi_w(x)||^2 \quad (12)$$

MMD<sup>10</sup>: using the samples mean after nonlinear transformation as a substitute for Expectation:

$$||\frac{1}{N} \sum_{i=1}^N \Phi(x_i) - \frac{1}{M} \sum_{j=1}^M \Phi(x_j)||^2 \quad (13)$$

Make using of kernel trick, no Discriminator, optimize the objective function w.r.t Generator

<sup>10</sup>Training generative neural networks via Maximum Mean Discrepancy optimization

## Perspective of Info. Theory<sup>11 12</sup>

True distribution  $P_r$ , generated distribution  $P_{g_\theta}$ ,  $Z \sim \text{Ber}(\pi)$ , the random variable  $X$  satisfies:

$$P(X|Z=0) = P_g, P(X|Z=1) = P_r \quad (14)$$

Generator  $G_\theta$  hopes to inference the value of  $\pi$ , try to independent with  $Z$ , it can be achieved by minimizing the mutual information:

$$I(X, Z) = KL(p(x, z) || p(x)p(z)) \quad (15)$$

minimize mutual information  $\iff X, Z$  are independent,  $\iff P_r = P_g$

$$\begin{aligned} I(X, Z) &= H(Z) + \mathbb{E}_X \mathbb{E}_{Z|X} \log q(z|x) + \mathbb{E}_X KL[p(z|x) || q(y|x)] \\ &= \max_q H(Z) + \mathbb{E}_X \mathbb{E}_{Z|X} \log q(z|x) \end{aligned} \quad (16)$$

$$\begin{aligned} I(X, Z) &\geq H(Z) + \max_{\Psi} \mathbb{E}_{X,Z} \log q(z|x; \Psi) \\ &= H(Z) + \max_{\Psi} \pi \mathbb{E}_{P_r} \log q(1|x; \Psi) + (1 - \pi) \mathbb{E}_{P_g} \log q(0|x; \Psi) \end{aligned} \quad (17)$$

$$\min I(X, Z) \implies \min_{g_\theta} \max_{\Psi} \pi \mathbb{E}_{P_r} \log q(1|x; \Psi) + (1 - \pi) \mathbb{E}_{P_g} (1 - \log q(1|x; \Psi))$$

## Conclusion

The Non-parametric Estimation of the distance between  $P_r$  and  $P_g$ <sup>13</sup>:

### Definition

Integral Probability Metrics:

$$\gamma_F(P_r, P_g) := \sup_{f \in F} \left| \int_M f dP_r - \int_M f dP_g \right| \quad (18)$$

in which  $F$  is a set which contains all real valued, bounded measurable functions on  $M$

All the foregoing models differs in the choice of  $F$

- Wasserstein distance:  $F = \{f: \|f\|_L \leq 1\}$
- TV distance or Kolmogorov distance:  $F = \{f: \|f\|_\infty \leq 1\}$
- MMD:  $F = \{f: \|f\|_H \leq 1\}$

<sup>13</sup>Non-parametric Estimation of Integral Probability Metrics