



>>> Paper Reading: LSGAN

>>> Least Squares Generative Adversarial Networks, ICCV 2017.  
On the Effectiveness of LSGANs, TPAMI 2019.

Name: 李喆琛 信息科学技术学院  
初济群 数学科学学院

Date: 2020.05.14 → 第十三周·第七场



## >>> Outline

1. Regular GAN and Least Squares GAN
2. Why is LSGAN better?
3. Theoretical Analysis
4. Deficiencies of LSGAN



# >>> Regular Generative Adversarial Networks I

## Objective function of regular GAN

$$\min_G \max_D V_{\text{GAN}}(D, G) = \mathbb{E}_{x \sim p_{\text{data}}(x)}[\log D(x)] + \mathbb{E}_{z \sim p_x(z)}[\log(1 - D(G(z)))]$$

### Why using sigmoid cross entropy loss function?

Real distribution  $p$ , Fake distribution  $q$ .

- \* *Information entropy*:  $H(p) = -\sum_i p(i) \cdot \log p(i)$ ;
  - \*  $H(p) \uparrow$ , Uncertainty  $\uparrow$ .
- \* *Cross entropy*:  $H(p, q) = -\sum_i p(i) \cdot \log q(i)$ ;
  - \*  $H(p, q) \uparrow$ , Difference  $\uparrow$ .

The KL divergence  $D(p|q) = H(p, q) - H(p)$  is a way to measure the distance between two distributions  $p$  and  $q$ .

- \* it comes to its minimum point when  $p = q$ .

Thus, GAN minimizes KL divergence.



## >>> Regular Generative Adversarial Networks II

### A problem: vanishing gradients

When  $D$  is closing to its optimal point  $D^*$ , the gradient of  $G$  is also closing to zeros, that is :

$$\nabla_x E_{z \sim p(z)} [\log(1 - D(G(z)))] \approx 0$$

As  $P_r$  and  $P_g$  are two low dimension manifolds, discriminator  $D$  is easy to train, and thus closing to optimal point  $D^*$  quickly.

### How to get over it ?

- \* Method: improved GAN
- \* problem: Oscillations and Mode Collapse

However, LSGAN don't have both the problems of vanilla GAN and improved GAN theorecitally.



# >>> Least Squares Generative Adversarial Networks

Sigmoid cross entropy  $\rightarrow$  Least squares loss function

## Objective function of LSGAN

$$\min_D V_{\text{LSGAN}}(D) = \frac{1}{2} \mathbb{E}_{x \sim p_{\text{data}}(x)} [(D(x) - b)^2] + \frac{1}{2} \mathbb{E}_{z \sim p_z(z)} [(D(G(z)) - a)^2]$$
$$\min_G V_{\text{LSGAN}}(G) = \frac{1}{2} \mathbb{E}_{z \sim p_z(z)} [(D(G(z)) - c)^2]$$

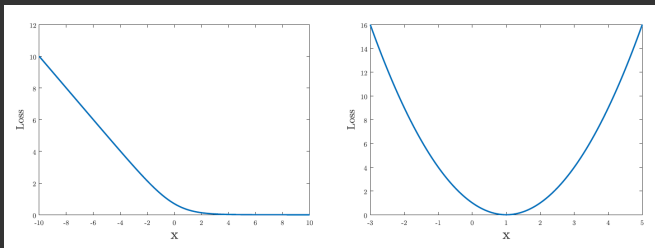
## The Parameters

- \*  $a$ : labels for fake data; Discriminator
- \*  $b$ : labels for real data; Discriminator
- \*  $c$ : values that  $G$  wants  $D$  to believe. Generator



# >>> Why is LSGAN better? I

More difficult to saturate → Better stability



**Figure:** Sigmoid  $y = -\log(1 - \frac{1}{1+e^x})$  vs. Least Square  $y = (x - 1)^2$

- \* Least squares loss function is flat only at one point;
- \* Sigmoid cross entropy will saturate when  $x$  is large.



## >>> Why is LSGAN better? II

Tougher penalties  $\rightarrow$  Higher Quality

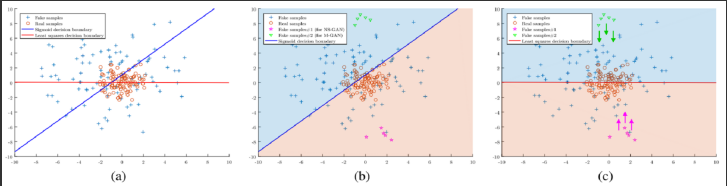
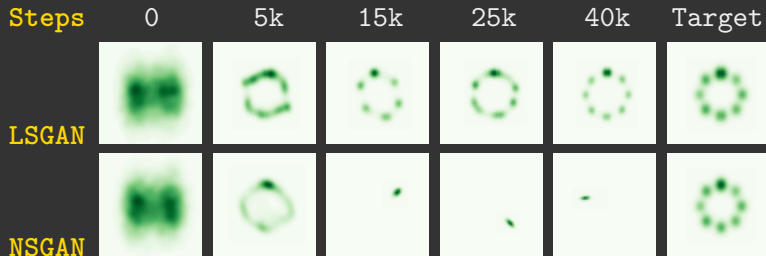


Figure: Real Samples  $\rightarrow$  Orange, Fake Samples  $\rightarrow$  Blue

- \* *Vanilla GAN*:  $\nabla$  in GREEN  $\rightarrow$  Little error
  - \* Leads to the problem of **Vanishing Gradients**.
- \* *Improved GAN*:  $\star$  in PINK  $\rightarrow$  Little error
  - \* Leads to the problem of **Mode Collapse**.
- \* *Least Squares*: Penalize samples far from the boundary.
  - \* Forces  $G$  to generate samples toward decision boundary.

# >>> Comparison of the results



## FID Results on Four Datasets

Method	LSUN	Cat	ImageNet	CIFAR10
NS-GAN	28.04	15.81	74.15	35.25
WGAN-GP	22.77	29.03	<b>62.05</b>	40.83
LSGAN <sub>(011)</sub>	27.21	15.46	72.54	36.46
LSGAN <sub>(-110)</sub>	<b>21.55</b>	<b>14.28</b>	68.95	<b>35.19</b>





## >>> Relation to Pearson $\chi^2$ Divergence I

Consider the following extensions:

$$\min_D V_{\text{LSGAN}}(D) = \frac{1}{2} \mathbb{E}_{x \sim p_{\text{data}}(x)} [(D(x) - b)^2] + \frac{1}{2} \mathbb{E}_{z \sim p_z(z)} [(D(G(z)) - a)^2]$$

$$\min_G V_{\text{LSGAN}}(G) = \frac{1}{2} \mathbb{E}_{x \sim p_{\text{data}}(x)} [(D(x) - c)^2] + \frac{1}{2} \mathbb{E}_{z \sim p_z(z)} [(D(G(z)) - c)^2]$$

\* Note that  $\mathbb{E}_{x \sim p_{\text{data}}(x)} [(D(x) - c)^2]$  does not contain  $G$ .

### Optimal Discriminator

For a fixed  $G$ , the optimal discriminator  $D$  is:

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



# >>> Relation to Pearson $\chi^2$ Divergence II

## Proof of the optimal discriminator

In fact, we are trying to minimize  $V(D)$ :

$$\begin{aligned}
V(D) &= \frac{1}{2} \mathbb{E}_{x \sim p_{\text{data}}} [(D(x) - b)^2] + \frac{1}{2} \mathbb{E}_{z \sim p_z} [(D(G(z)) - a)^2] \\
&= \frac{1}{2} \mathbb{E}_{x \sim p_{\text{data}}} [(D(x) - b)^2] + \frac{1}{2} \mathbb{E}_{x \sim p_g} [(D(x) - a)^2] \\
&= \int_{\mathcal{X}} \frac{1}{2} (p_{\text{data}}(x)(D(x) - b)^2 + p_g(x)(D(x) - a)^2) dx \leftarrow \text{Denoted by } Y
\end{aligned}$$

Let its derivative be zero:

$$\frac{dY}{dD(x)} = p_{\text{data}}(x)(D(x) - b) + p_g(x)(D(x) - a) = 0$$

Then we have:

$$D(x) = \frac{bp_{\text{data}}(x) + ap_g(x)}{p_{\text{data}}(x) + p_g(x)} \leftarrow \text{Denoted by } D^*(x)$$

In other word,  $D^*(x)$  minimizes  $V(D)$ . □



## >>> Relation to Pearson $\chi^2$ Divergence III

**Theorem.** Optimizing LSGANs yields minimizing Pearson  $\chi^2$  divergence between  $p_d + p_g$  and  $p_g$ , if  $b - c = 1$ , and  $b - a = 2$ .

**Proof.** Substitute  $D^*(x)$  into the equation:

$$\begin{aligned} 2C(G) &= \mathbb{E}_{x \sim p_d} \left[ (D^*(x) - c)^2 \right] + \mathbb{E}_{z \sim p_z} \left[ (D^*(G(z)) - c)^2 \right] \\ &= \mathbb{E}_{x \sim p_d} \left[ (D^*(x) - c)^2 \right] + \mathbb{E}_{x \sim p_g} \left[ (D^*(x) - c)^2 \right] \\ &= \mathbb{E}_{x \sim p_d} \left[ \left( \frac{bp_d(x) + ap_g(x)}{p_d(x) + p_g(x)} - c \right)^2 \right] + \mathbb{E}_{x \sim p_g} \left[ (\dots)^2 \right] \\ &= \int_{\mathcal{X}} p_d(x) \left( \frac{(b-c)p_d(x) + (a-c)p_g(x)}{p_d(x) + p_g(x)} \right)^2 dx + \int_{\mathcal{X}} p_g(x) (\dots)^2 dx \\ &= \int_{\mathcal{X}} \frac{((b-c)p_d(x) + (a-c)p_g(x))^2}{p_d(x) + p_g(x)} dx \end{aligned}$$



## >>> Relation to Pearson $\chi^2$ Divergence IV

Let  $b - c = 1$ ,  $b - a = 2$ , we have:

$$\begin{aligned} 2C(G) &= \int_{\mathcal{X}} \frac{((b - c)(p_d(x) + p_g(x)) - (b - a)p_g(x))^2}{p_d(x) + p_g(x)} dx \\ &= \int_{\mathcal{X}} \frac{(2p_g(x) - (p_d(x) + p_g(x)))^2}{p_d(x) + p_g(x)} dx \\ &= \chi_{\text{Pearson}}^2(p_d + p_g \| 2p_g) \end{aligned}$$

That proves the theorem. □



## >>> Parameters Selection

Let  $b - c = 1$ ,  $b - a = 2 \Rightarrow$  Minimizing Pearson  $\chi^2$  Divergence

For example,  $a = -1$ ,  $b = 1$ ,  $c = 0$ :

$$\min_D V_{\text{LSGAN}}(D) = \frac{1}{2} \mathbb{E}_{x \sim p_{\text{data}}(x)} [(D(x) - 1)^2] + \frac{1}{2} \mathbb{E}_{z \sim p_z(z)} [(D(G(z)) + 1)^2]$$

$$\min_G V_{\text{LSGAN}}(G) = \frac{1}{2} \mathbb{E}_{z \sim p_z(z)} [(D(G(z)))^2]$$

Let  $b = c \Rightarrow$  Generating samples as real as possible

For example,  $a = 0$ ,  $b = -1$ ,  $c = -1$ :

$$\min_D V_{\text{LSGAN}}(D) = \frac{1}{2} \mathbb{E}_{x \sim p_{\text{data}}(x)} [(D(x) - 1)^2] + \frac{1}{2} \mathbb{E}_{z \sim p_z(z)} [(D(G(z)))^2]$$

$$\min_G V_{\text{LSGAN}}(G) = \frac{1}{2} \mathbb{E}_{z \sim p_z(z)} [(D(G(z)) - 1)^2]$$



>>> Better than regular GAN, but not good enough I

If  $D$  is good enough, the problem still exists

Without loss of generality, let  $c = 0$ ,  
then the optimal discriminator is:

$$D^* = \frac{bP_d + aP_g}{P_d + P_g}$$

Plug it into the extended loss function:

$$V_{\text{LSGAN}}(G) = \frac{1}{2} \mathbb{E}_{x \sim p_d(x)} [D^*(x)^2] + \frac{1}{2} \mathbb{E}_{x \sim p_g(x)} [D^*(x)^2]$$

If **supp**  $p_d$  and **supp**  $p_g$  are low dimensional manifolds in high dimensional space, we have

$$\mathbb{P}[\mu(\text{supp } p_d \cap \text{supp } p_g) = 0] = 1 \tag{1}$$



## >>> Better than regular GAN, but not good enough II

For a given  $x$ , there are only **four** cases:

1.  $p_d(x) = 0, p_g(x) = 0$ : We can ignore this case;
2.  $p_d(x) \neq 0, p_g(x) = 0$ :  $D^*(x) = a$ ,  $V_{\text{LSGAN}}$  is a constant;
3.  $p_d(x) = 0, p_g(x) \neq 0$ :  $D^*(x) = b$ ,  $V_{\text{LSGAN}}$  is a constant;
4.  $p_d(x) \neq 0, p_g(x) \neq 0$ : Will not happen due to Equ (1).

### Gradients vanish again!

For LSGANs, if:

- \* *Equ (1) holds*; (Support sets are disjoint)
- \* *Discriminator  $D$  is good enough*; (very close to  $D^*$ )

then the loss will be zero  $\rightarrow$  Gradient vanish.

### Another point of view: WGAN (Wasserstein metric)

Not Lipschitz continuous  $\rightarrow$  Vanishing gradients.

To conclude, LSGAN cannot solve the problem completely.



## >>> To summarize

### 1. *What is LSGAN:*

1.1 Sigmoid cross entropy  $\rightarrow$  Least squares;

### 2. *Benifits of LSGAN:*

2.1 Better stability;

2.2 Higher quailty;

2.3 Partially solves the problem of vanishing gradients;

### 3. *Theoretical Properties:*

3.1 Convergence: LSGAN minimizes Pearson  $\chi^2$  Divergence;

### 4. *Deficiencies of LSGAN:*

4.1 Cannot solve the problem of completely.



Thanks For Listening!