# WGAN

Martin Arjovsky, Soumith Chintala, Léon Bottou, Wasserstein GAN, arXiv prepring, 2017

Ishaan Gulrajani, Faruk Ahmed, Martin Arjovsky, Vincent Dumoulin, Aaron Courville, "Improved Training of Wasserstein GANs", arXiv prepring, 2017

# JS divergence is not suitable

- In most cases, $P_G$ and $P_{data}$ are not overlapped.
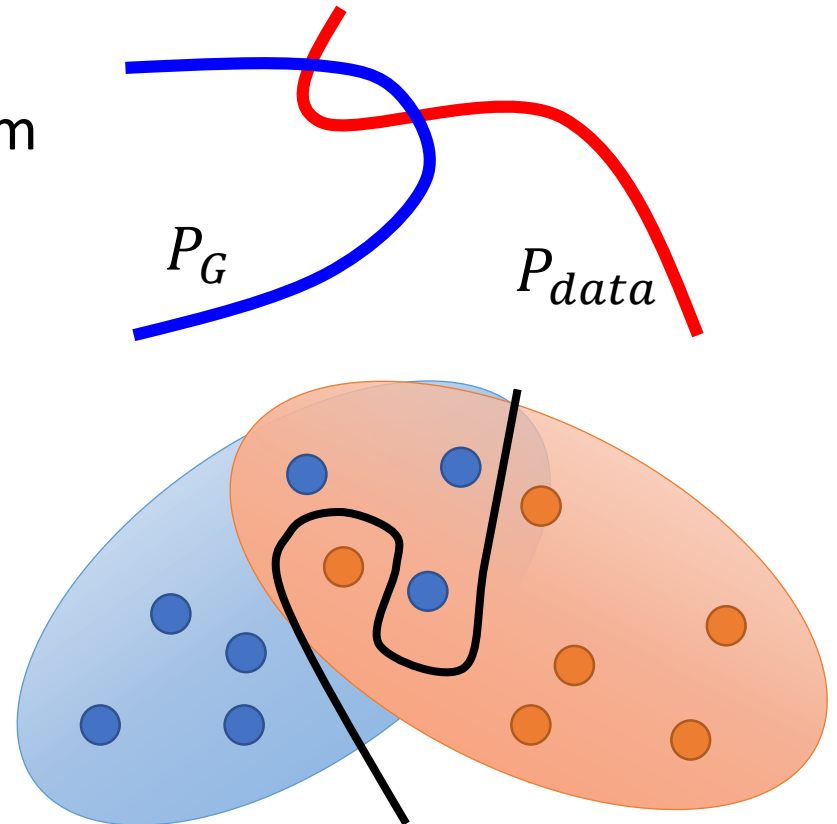
- 1. The nature of data

    Both $P_{data}$ and $P_G$ are low-dim manifold in high-dim space.
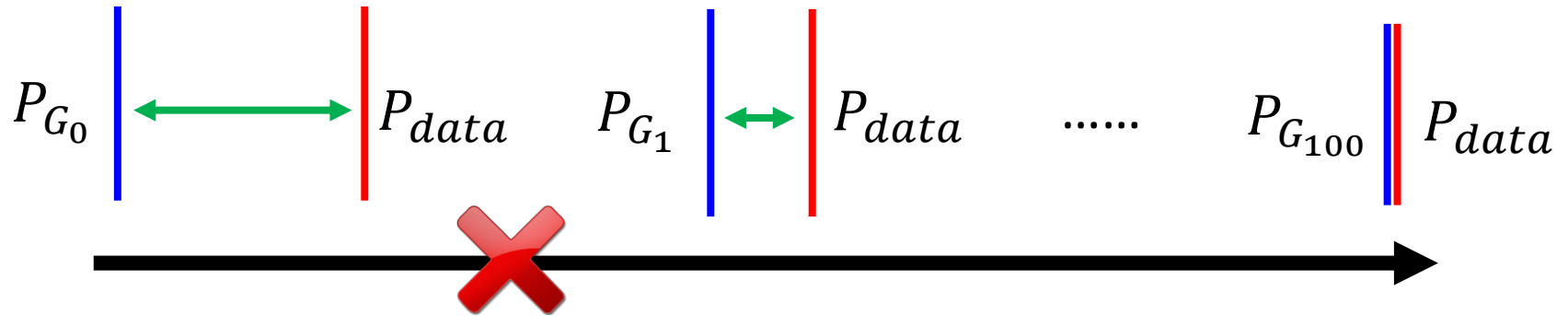
    The overlap can be ignored.

- 2. Sampling

    Even though $P_{data}$ and $P_G$ have overlap.

    If you do not have enough sampling ……

$P_G$

$P_{data}$

# *What is the problem of JS divergence?*

$$P_{G_0} \quad \xleftrightarrow{\hspace{2cm}} \quad P_{data} \qquad P_{G_1} \quad \xleftrightarrow{} \quad P_{data} \qquad \ldots\ldots \qquad P_{G_{100}} \quad P_{data}$$

Equally bad

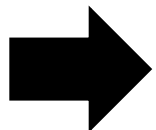$$JS(P_{G_0}, P_{data}) = log2 \qquad JS(P_{G_1}, P_{data}) = log2 \qquad \ldots\ldots \qquad JS(P_{G_{100}}, P_{data}) = 0$$
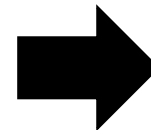
在GAN中的D，如果是train binary classifier，則會發現當兩坨資料可以完全分開的時候他們的loss都是一樣的

JS divergence is log2 if two distributions do not overlap.

Intuition: If two distributions do not overlap, binary classifier achieves 100% accuracy
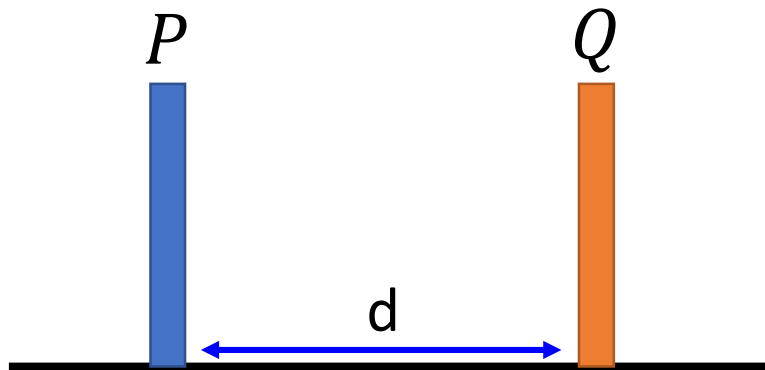
➡ Same objective value is obtained. ➡ Same divergence

可以改用LSGAN，將D的output的sigmoid拿掉，改為linear，這樣就變成regression problem
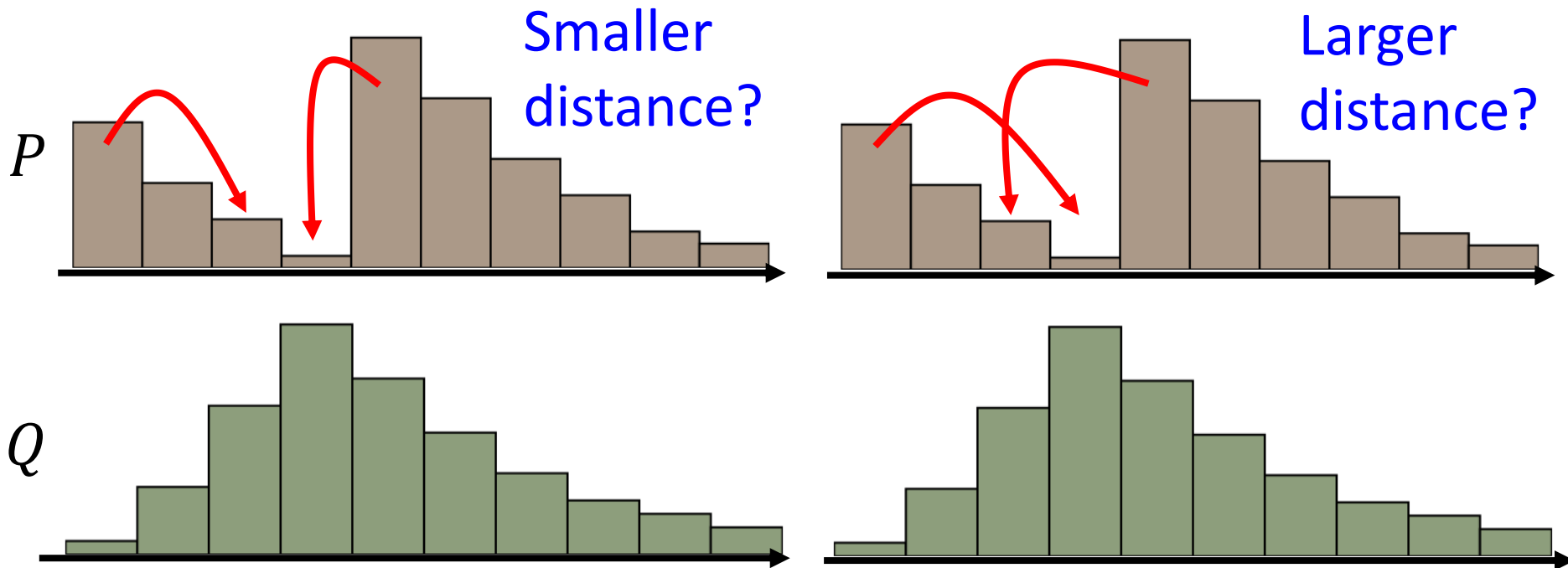
# Earth Mover's Distance

- Considering one distribution P as a pile of earth, and another distribution Q as the target

- The average distance the earth mover has to move the earth.

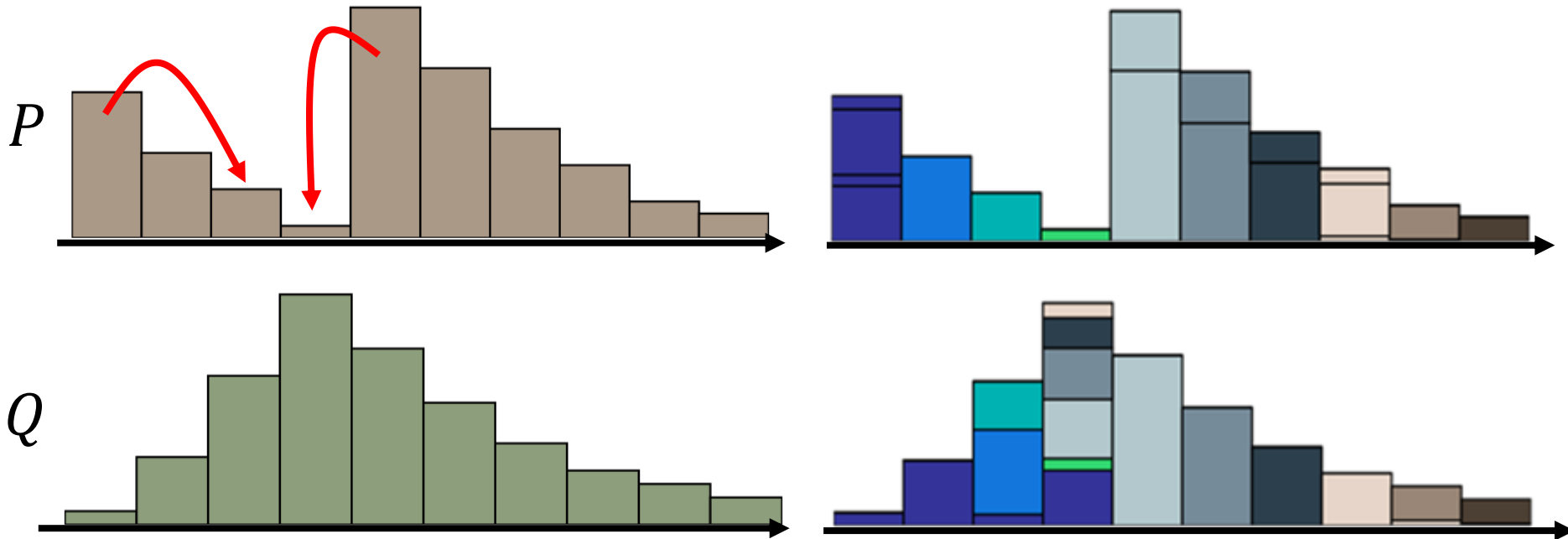$$W(P,Q) = d$$

# Earth Mover's Distance



P    Smaller distance?    Larger distance?

Q

There many possible "moving plans".    窮舉所有鏟土的方法找出最小distance

Using the "moving plan" with the smallest average distance to define the earth mover's distance.

Source of image: https://vincentherrmann.github.io/blog/wasserstein/
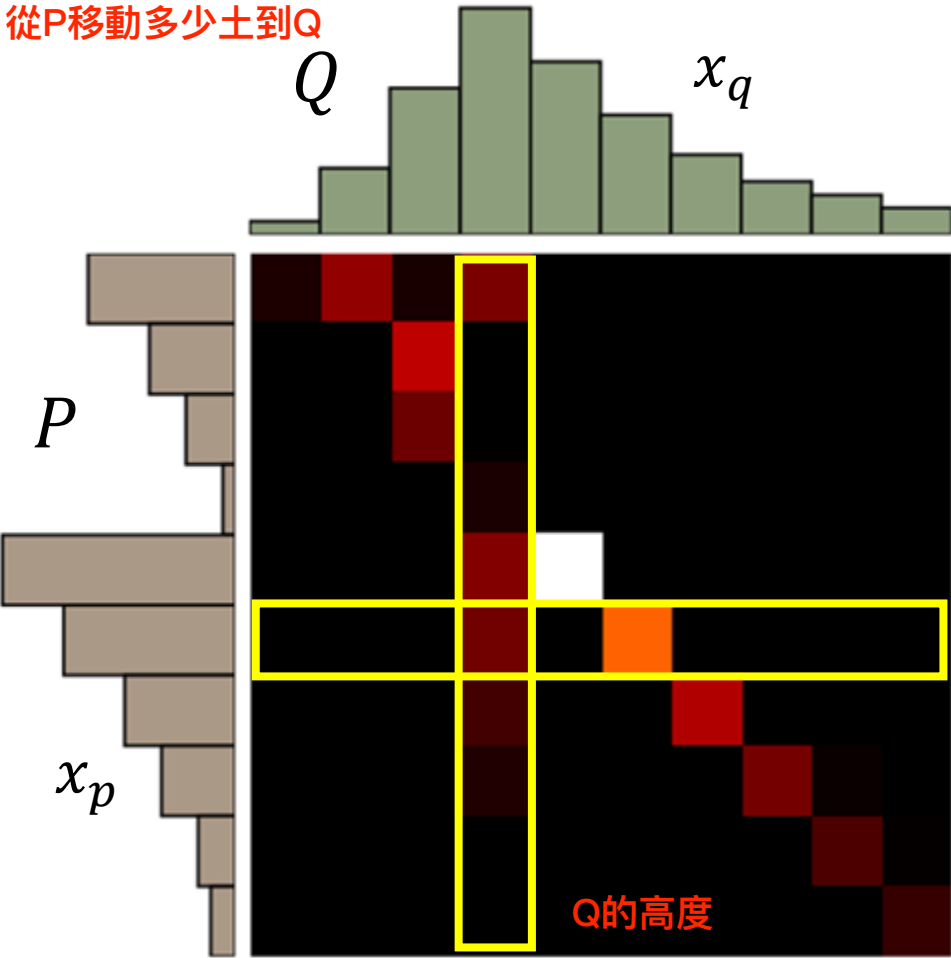
# Earth Mover's Distance

Best "moving plans" of this example



There many possible "moving plans".

Using the "moving plan" with the smallest average distance to define the earth mover's distance.

從P移動多少土到Q

$Q$ $x_q$

$P$

$x_p$

P的高度

Q的高度

moving plan $\gamma$
All possible plan $\Pi$

A "moving plan" is a matrix

The value of the element is the amount of earth from one position to another.
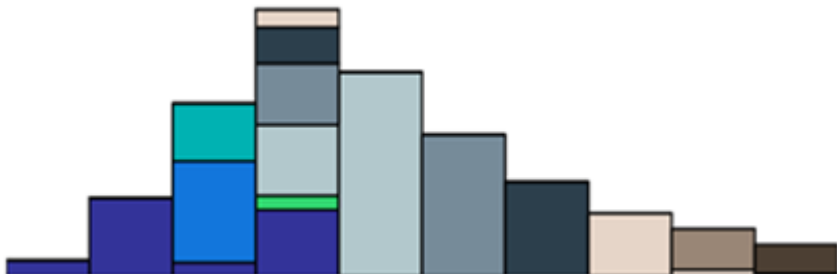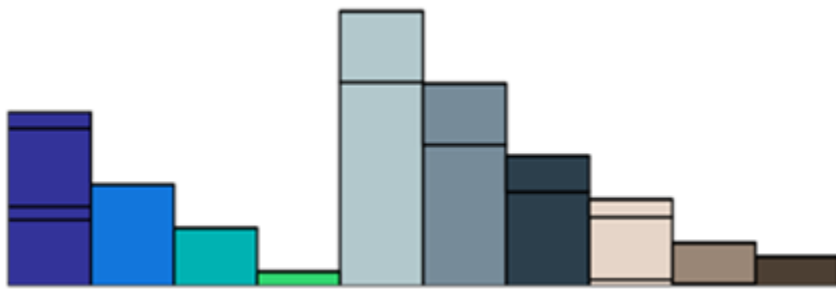
Average distance of a plan $\gamma$:

$$B(\gamma) = \sum_{x_p, x_q} \gamma(x_p, x_q)\|x_p - x_q\|$$

Earth Mover's Distance:

解optimization problem

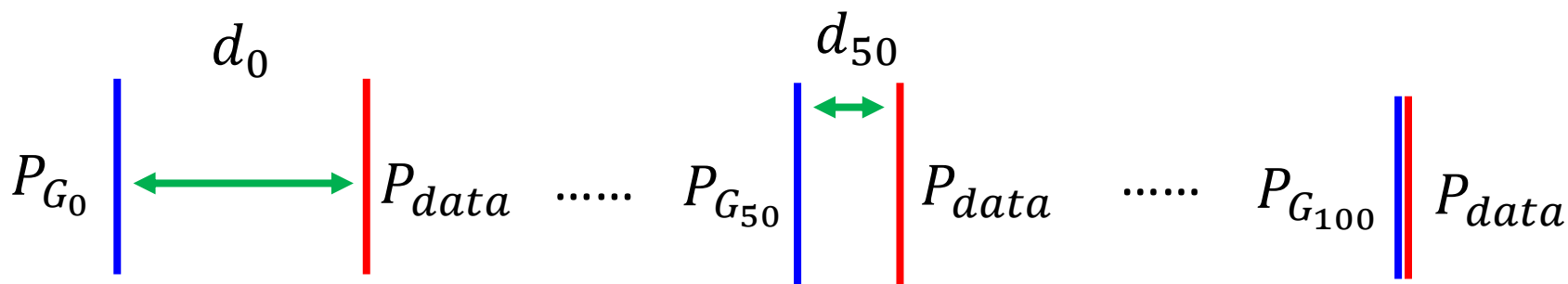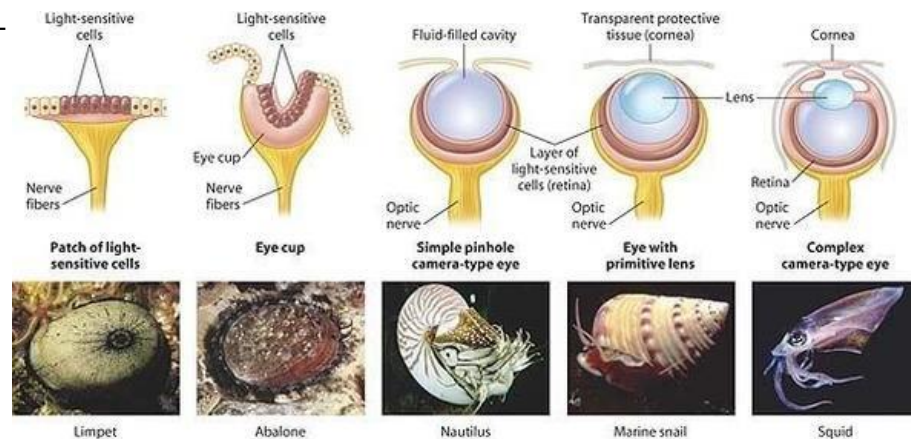$$W(P, Q) = \min_{\gamma \in \Pi} B(\gamma)$$

The best plan

# *Why Earth Mover's Distance?*

$$D_f(P_{data}||P_G)$$

$$\downarrow$$

$$W(P_{data}, P_G)$$



$d_0$       $d_{50}$

$P_{G_0}$   $P_{data}$   ......   $P_{G_{50}}$   $P_{data}$   ......   $P_{G_{100}}$   $P_{data}$

每一次都要有一點點的進步才能有好的演化結果

$$JS(P_{G_0}, P_{data}) = log2 \qquad JS(P_{G_{50}}, P_{data}) = log2 \qquad JS(P_{G_{100}}, P_{data}) = 0$$

$$W(P_{G_0}, P_{data}) = d_0 \qquad W(P_{G_{50}}, P_{data}) = d_{50} \qquad W(P_{G_{100}}, P_{data}) = 0$$

# Back to the GAN framework

$$D_f(P_{data}||P_G) \Rightarrow W(P_{data}, P_G)$$

原始的WGAN方法適用weight clipping

$$= \max_D \{E_{x \sim P_{data}}[D(x)] - E_{x \sim P_G}[f^*(D(x))]\}$$

$$W(P_{data}, P_G)$$
$$= \max_{D \in 1-Lipschitz} \{E_{x \sim P_{data}}[D(x)] - E_{x \sim P_G}[D(x)]\}$$
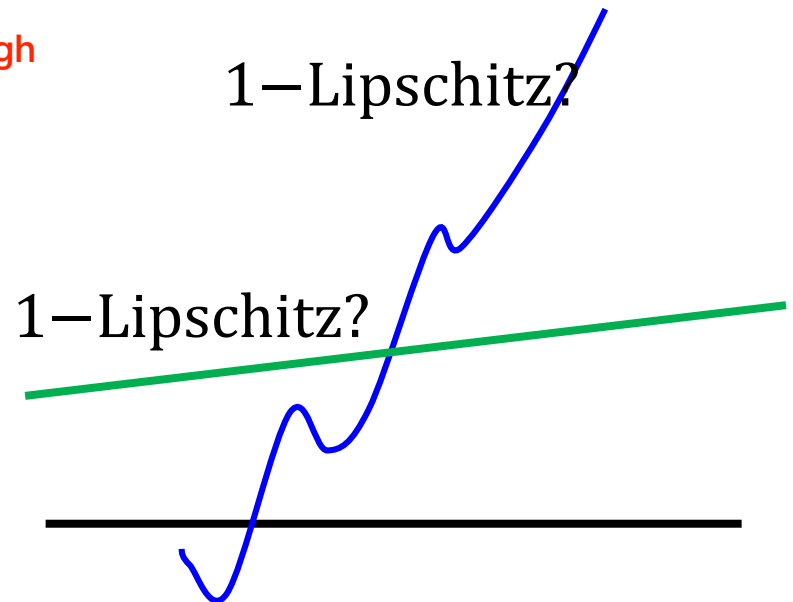
D has to be smooth enough

1−Lipschitz?

## *Lipschitz Function*

$$\|f(x_1) - f(x_2)\| \leq K\|x_1 - x_2\|$$

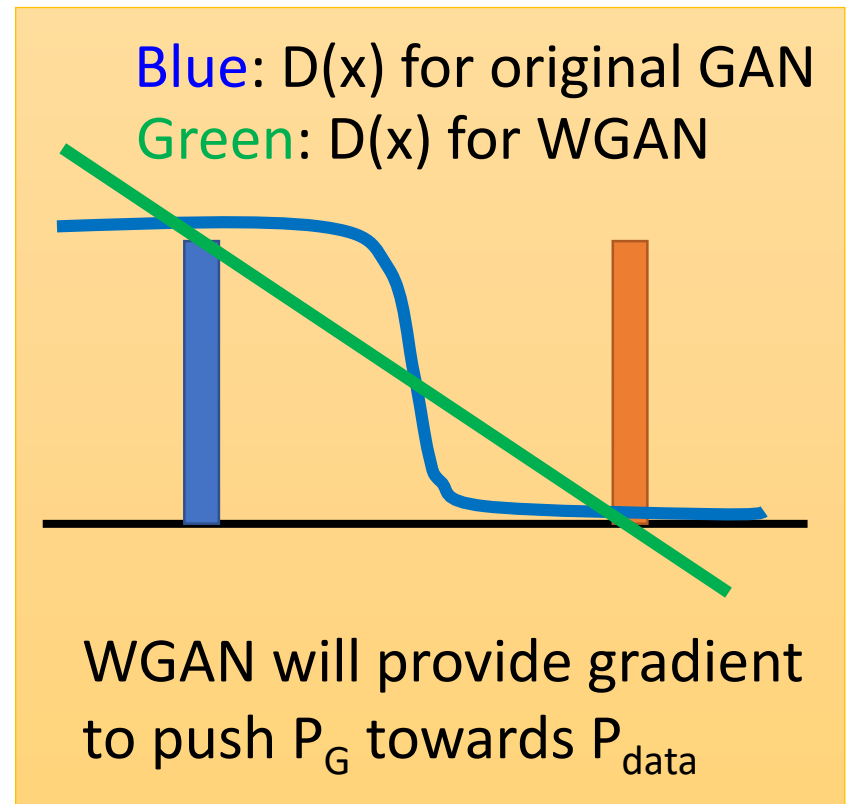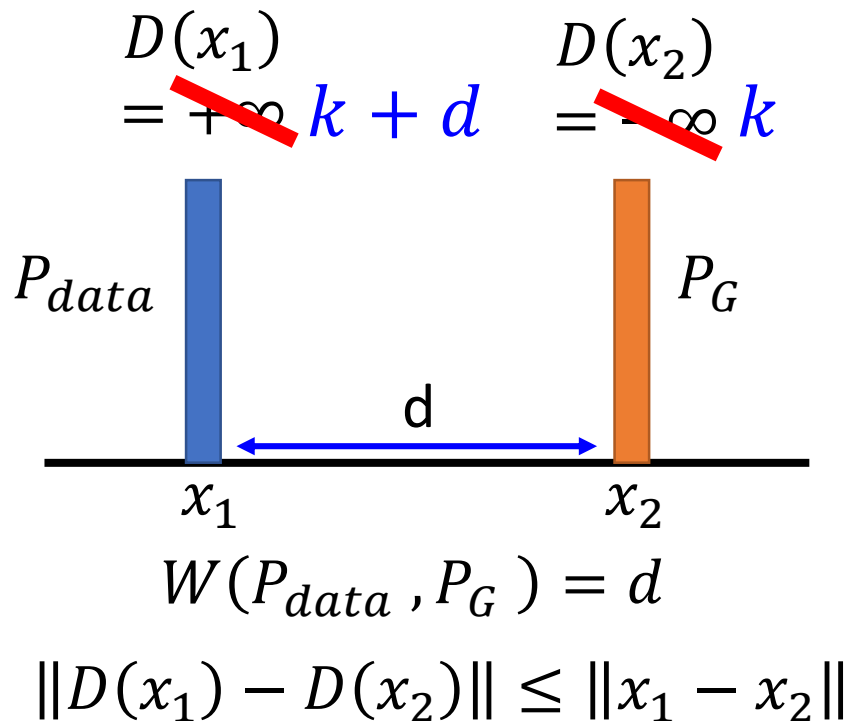Output change    Input change

1−Lipschitz?

K=1 for "$1 - Lipschitz$"

Do not change fast

# Back to the GAN framework

$$W(P_{data}, P_G)$$

$$= \max_{D \in 1-Lipschitz} \{E_{x \sim P_{data}}[D(x)] - E_{x \sim P_G}[D(x)]\}$$

$k + d$         $k$

$D(x_1)$
$= +\infty \quad k + d$

$D(x_2)$
$= -\infty \quad k$

$P_{data}$         $P_G$

d

$x_1$       $x_2$

$$W(P_{data}, P_G) = d$$

$$\|D(x_1) - D(x_2)\| \leq \|x_1 - x_2\|$$

Blue: D(x) for original GAN
Green: D(x) for WGAN

WGAN will provide gradient to push P$_G$ towards P$_{data}$

# Back to the GAN framework

K $W(P_{data}, P_G)$

$$= \max_{\substack{D \in 1-Lipschitz}} \{E_{x \sim P_{data}}[D(x)] - E_{x \sim P_G}[D(x)]\}$$

K

How to use gradient descent to optimize?

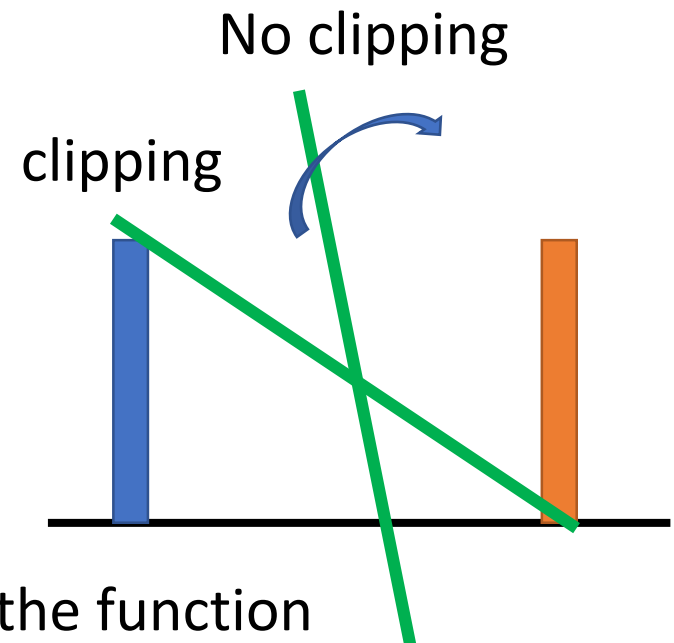**_Weight clipping:_**

Force the weights w between c and -c

After parameter update,
if w > c, then w=c; if w<-**c**, then w=-c

We only ensure that
$$\|D(x_1) - D(x_2)\| \le K\|x_1 - x_2\|$$

For some K

Do not truly find function D maximizing the function

No clipping

clipping

# *Algorithm of* WGAN

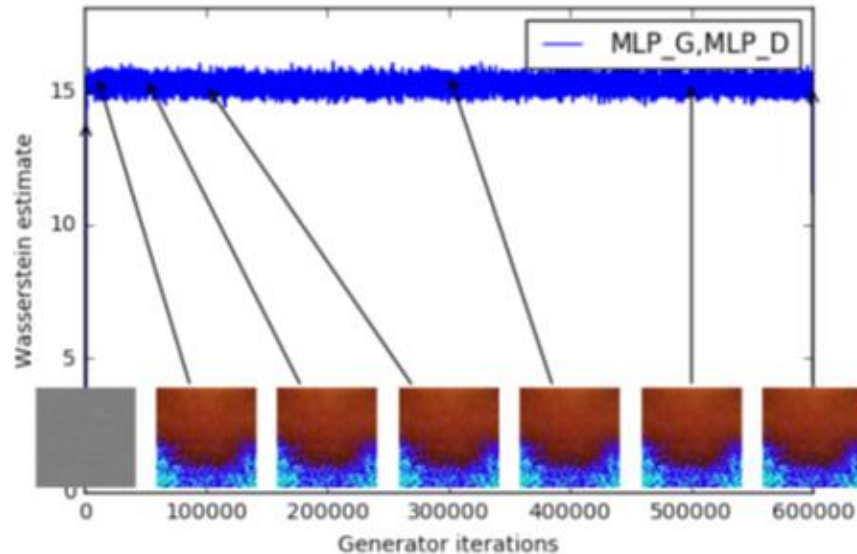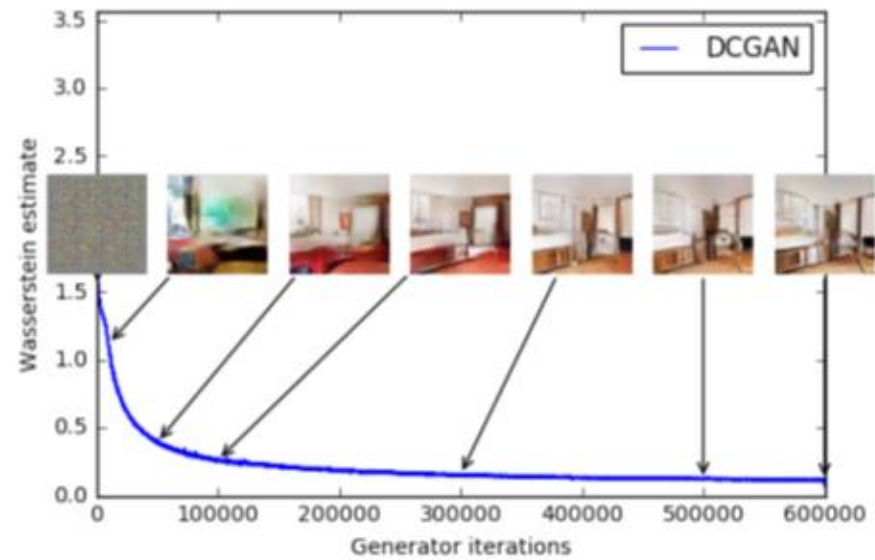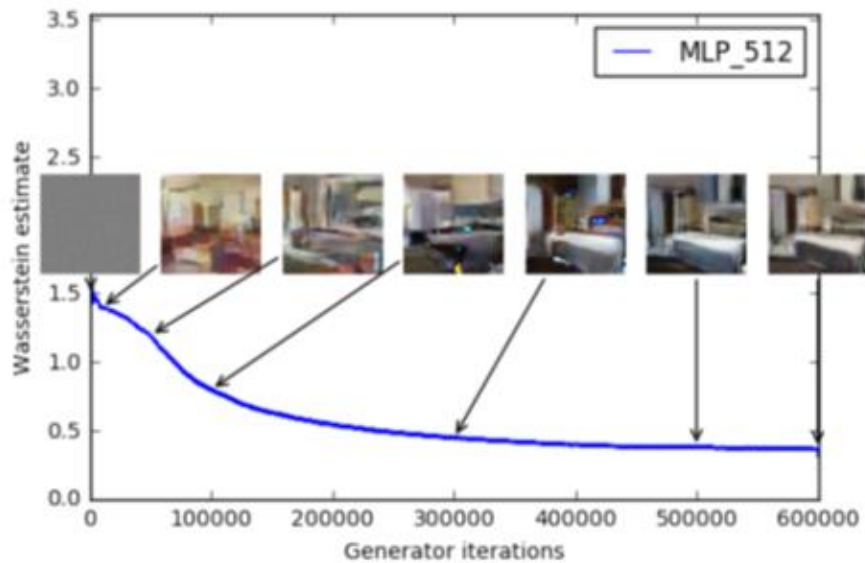- In each training iteration:   No sigmoid for the output of D

**Learning D**

**Repeat k times**

- Sample m examples $\{x^1, x^2, \ldots, x^m\}$ from data distribution $P_{data}(x)$
- Sample m noise samples $\{z^1, z^2, \ldots, z^m\}$ from the prior $P_{prior}(z)$
- Obtaining generated data $\{\tilde{x}^1, \tilde{x}^2, \ldots, \tilde{x}^m\}$, $\tilde{x}^i = G(z^i)$
- Update discriminator parameters $\theta_d$ to maximize
  - $\tilde{V} = \frac{1}{m}\sum_{i=1}^m D(x^i) - \frac{1}{m}\sum_{i=1}^m D(\tilde{x}^i)$
  - $\theta_d \leftarrow \theta_d + \eta \nabla \tilde{V}(\theta_d)$   Weight clipping

**Learning G**

**Only Once**

- Sample another m noise samples $\{z^1, z^2, \ldots, z^m\}$ from the prior $P_{prior}(z)$
- Update generator parameters $\theta_g$ to minimize
  - $\tilde{V} = \frac{1}{m}\sum_{i=1}^m logD(x^i) - \frac{1}{m}\sum_{i=1}^m D\left(G(z^i)\right)$
  - $\theta_g \leftarrow \theta_g - \eta \nabla \tilde{V}(\theta_g)$

**_Vertical_**

$$W(P_{data}, P_G)$$

$$= \max_{D \in 1-Lipschitz} \{ E_{x \sim P_{data}}[D(x)]$$

$$- E_{x \sim P_G}[D(x)] \}$$

https://arxiv.org/abs/1701.07875

# *Improved WGAN*

$$W(P_{data}, P_G)$$
$$= \max_{D \in 1-Lipschitz} \{E_{x \sim P_{data}}[D(x)] - E_{x \sim P_G}[D(x)]\}$$

A differentiable function is 1-Lipschitz if and only if it has gradients with norm less than or equal to 1 everywhere.

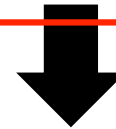$$D \in 1 - Lipschitz \quad \Longleftrightarrow \quad \|\nabla_x D(x)\| \leq 1 \text{ for all x}$$

等價的

$$W(P_{data}, P_G) \approx \max_D \{E_{x \sim P_{data}}[D(x)] - E_{x \sim P_G}[D(x)]$$

penalty，類似regularization

$$\lambda \int_x max(0, \|\nabla_x D(x)\| - 1)dx\}$$
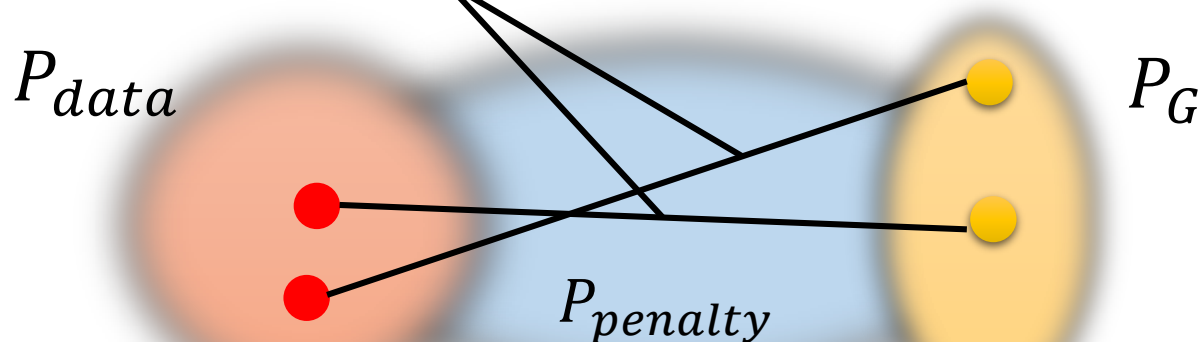
Prefer $\|\nabla_x D(x)\| \leq 1$ for all x

$$-\lambda E_{x \sim P_{penalty}}[max(0, \|\nabla_x D(x)\| - 1)]\}$$

但是無法真的對所有x(image)算gradient

Prefer $\|\nabla_x D(x)\| \leq 1$ for x sampling from $x \sim P_{penalty}$

# Improved WGAN

$$W(P_{data}, P_G) \approx \max_D \{E_{x \sim P_{data}}[D(x)] - E_{x \sim P_G}[D(x)]$$
$$-\lambda E_{x \sim P_{penalty}}[max(0, \|\nabla_x D(x)\| - 1)]\}$$

$P_{data}$

$P_G$

$P_{penalty}$

從兩個distribution中sample各兩個點連起來做interpolation，中間的值就是penalty

"Given that enforcing the Lipschitz constraint everywhere is intractable, enforcing it *only along these straight lines* seems sufficient and experimentally results in good performance."

Only give gradient constraint to the region between $P_{data}$ and $P_G$ because they influence how $P_G$ moves to $P_{data}$
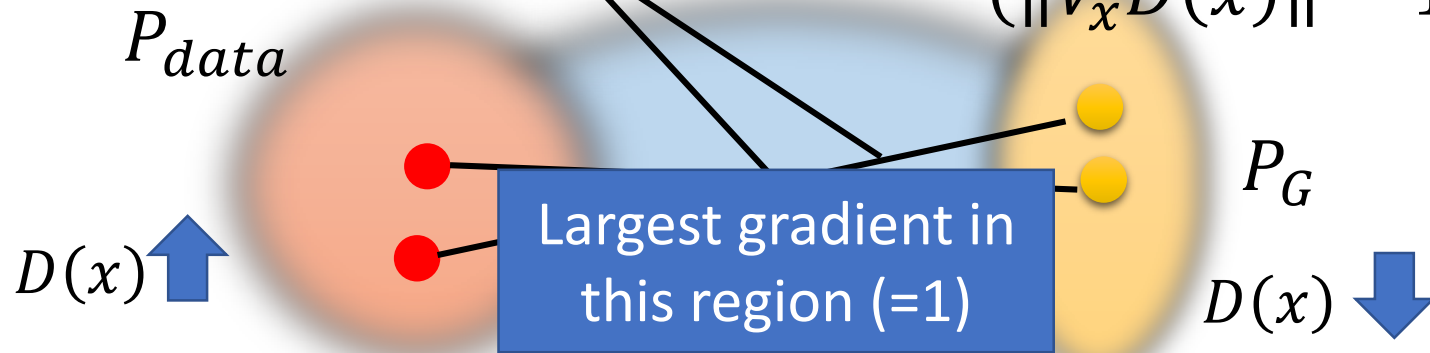
其實是滿合理的，因為要移動PG就是要參考中間這些區域做gradient descend

# *Improved WGAN*

$$W(P_{data}, P_G) \approx \max_D \{E_{x \sim P_{data}}[D(x)] - E_{x \sim P_G}[D(x)]$$

$$-\lambda E_{x \sim P_{penalty}}[\overline{max(0, \|\nabla_x D(x)\| - 1)}]\}$$

$$(\|\nabla_x D(x)\| - 1)^2$$

$P_{data}$

$P_G$

$D(x)$ ⬆

Largest gradient in this region (=1)

$D(x)$ ⬇

"One may wonder why we penalize the norm of the gradient for differing from 1, instead of just penalizing large gradients. The reason is that the optimal critic ... actually has gradients with norm 1 almost everywhere under Pr and Pg"
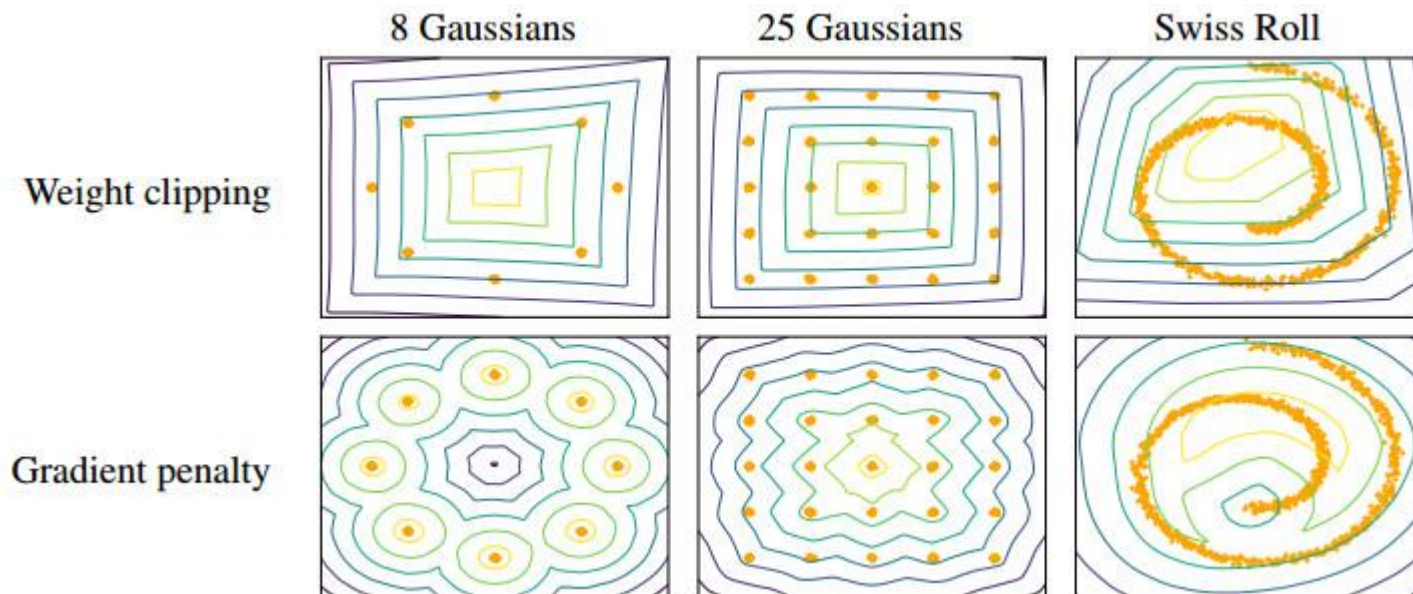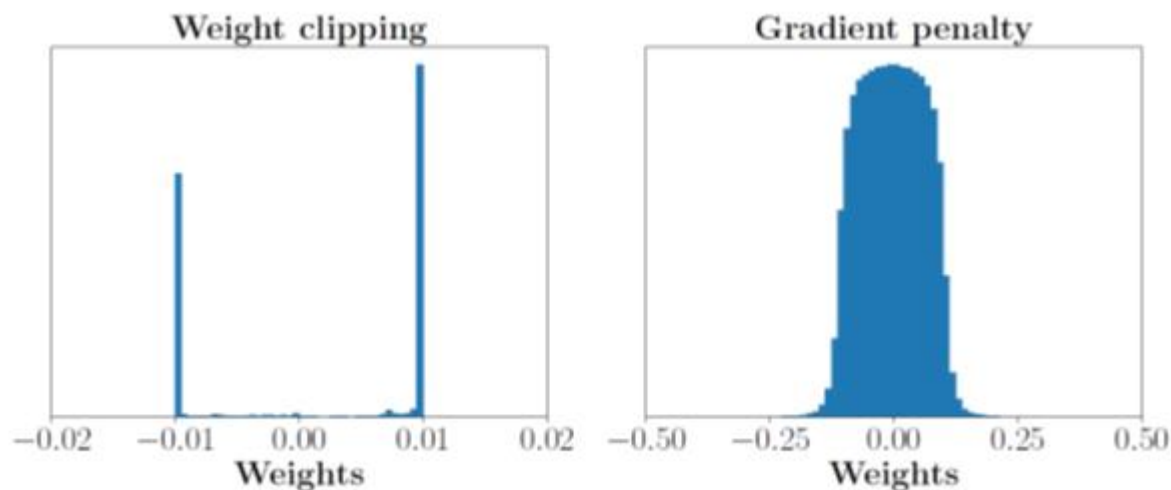
(check the proof in the appendix)

"Simply penalizing overly large gradients also works in theory, but experimentally we found that this approach converged faster and to better optima."

# Improved WGAN

| DCGAN | LSGAN | Original WGAN | Improved WGAN |
|---|---|---|---|

**G: CNN, D: CNN**



**G: CNN (no normalization), D: CNN (no normalization)**
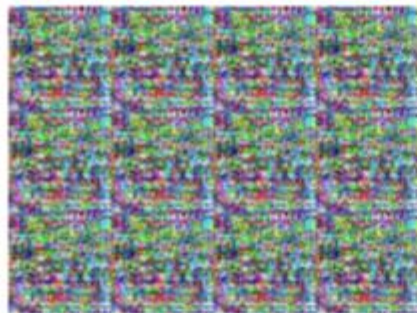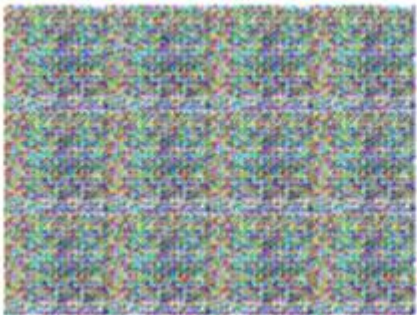


**G: CNN (tanh), D: CNN(tanh)**
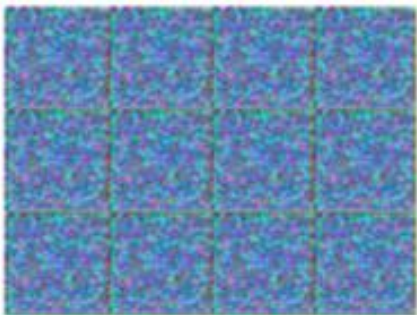
| DCGAN | LSGAN | Original WGAN | Improved WGAN |

G: MLP, D: CNN

G: CNN (bad structure), D: CNN

G: 101 layer, D: 101 layer

# Spectrum Norm

很強！！有空看一下

# Energy-based GAN

Ref: Junbo Zhao, Michael Mathieu, Yann LeCun, Energy-based Generative Adversarial Network, ICRL 2017
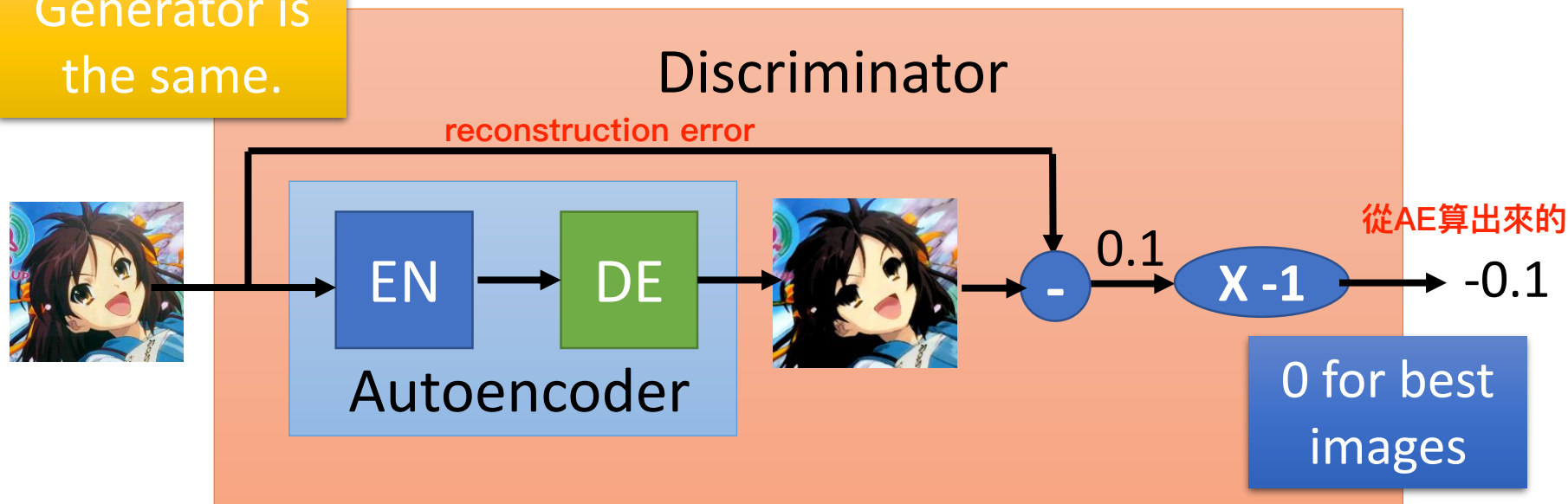
# Energy-based GAN (EBGAN)

- Using an autoencoder as discriminator D

An image is good. **=** It can be reconstructed by autoencoder.

Generator is the same.

**Discriminator**

reconstruction error

EN → DE

Autoencoder

**-** 0.1 → **X -1** → -0.1
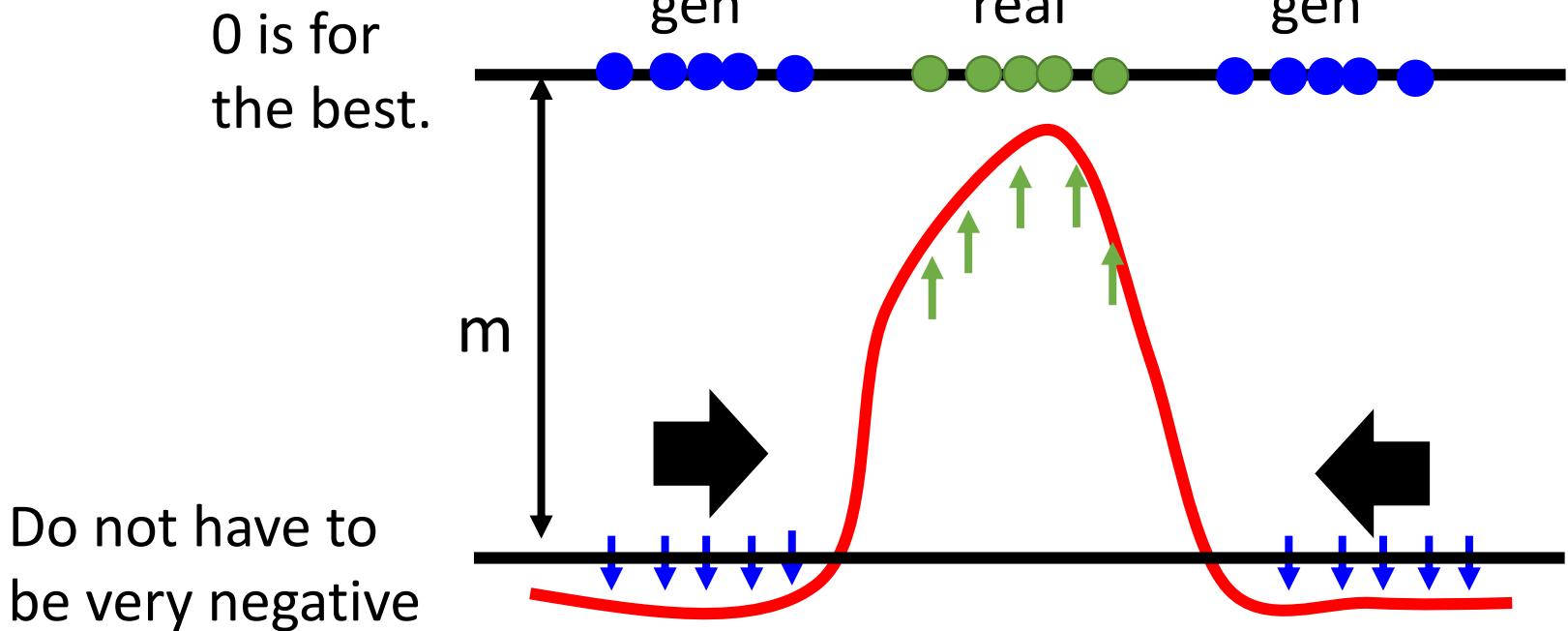
從AE算出來的

0 for best images

優點是Discriminator可以pre-trained，只要讓他看許多positive example就可以根據reconstruct error pretrain

# EBGAN

Auto-encoder based discriminator
only give limited region large value.

<span style="color:red">定義一個threshold，使得Generator產生的reconstruction error不會train壞掉</span>

<span style="color:red">因為擔心G發現real data部分無法上升太多，那至少上generative data 下降多一點，就會破壞到結果</span>

0 is for
the best.

gen          real          gen

m

Do not have to
be very negative

Hard to reconstruct, easy to destroy