

# Improved GAN

王宗威

# 回顾

$$G^* = \arg \min_G \max_D V(G, D)$$

$$V = E_{x \sim P_{data}} [\log D(x)] + E_{x \sim P_G} [\log(1 - D(x))]$$

**Algorithm** Initialize  $\theta_d$  for D and  $\theta_g$  for G

- In each training iteration:

Can only find lower bound of  $\max_D V(G, D)$

Learning  
D

Repeat  
k times

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

Learning  
G

Only  
Once

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

# Original GAN的缺点

问题：现在的GAN是否能指定生成想要的图片呢？

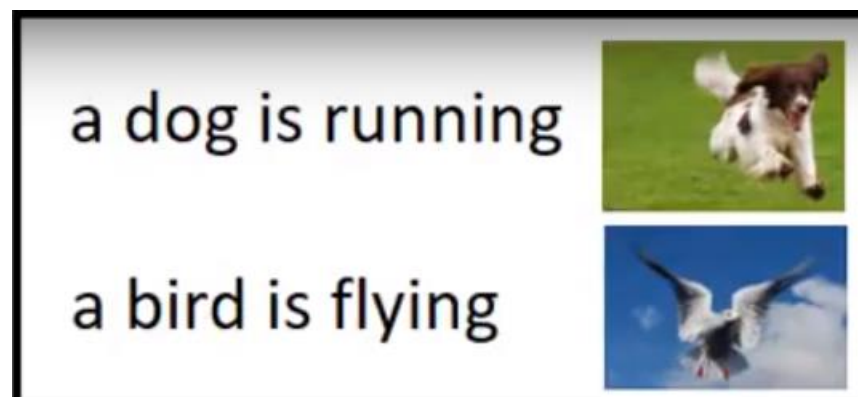
比如，

(1) case one:

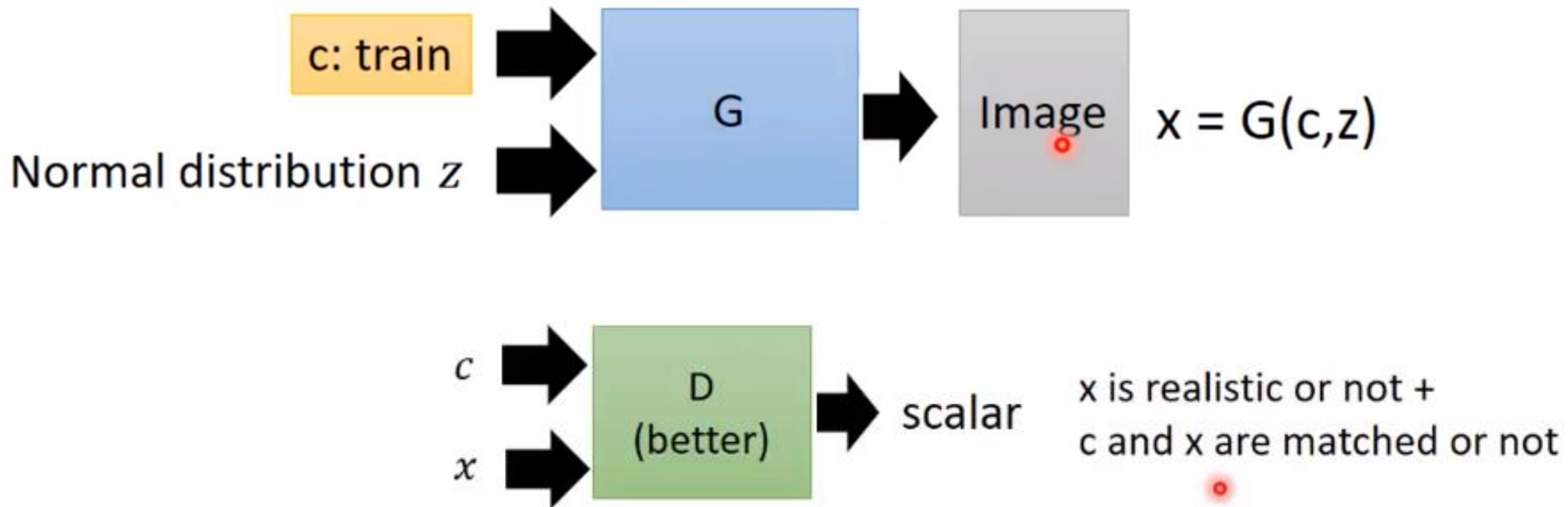
我们现在的任务是只需要9这个数字

(2) case two:

我们给一句话， a dog is running  
希望他生成一只小狗



# Conditional GAN



生成器

输入：随机噪声 $z$ ，标签  
向量 $c$

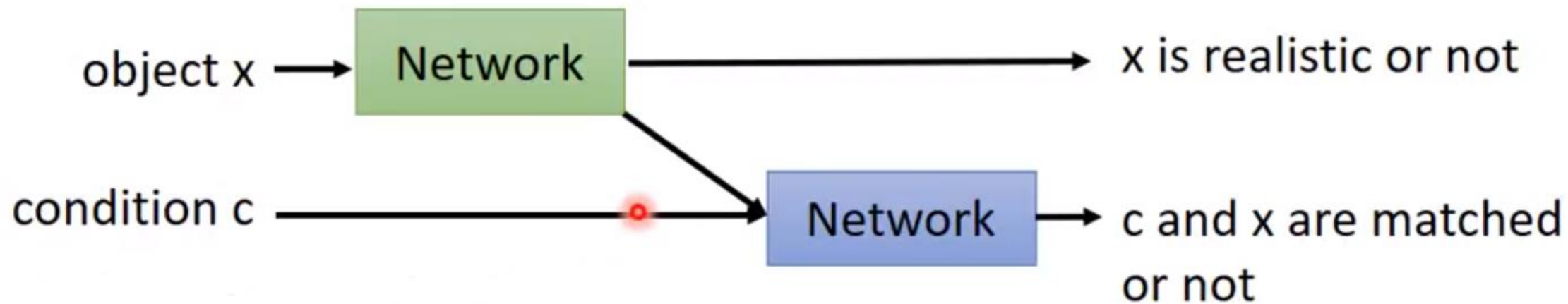
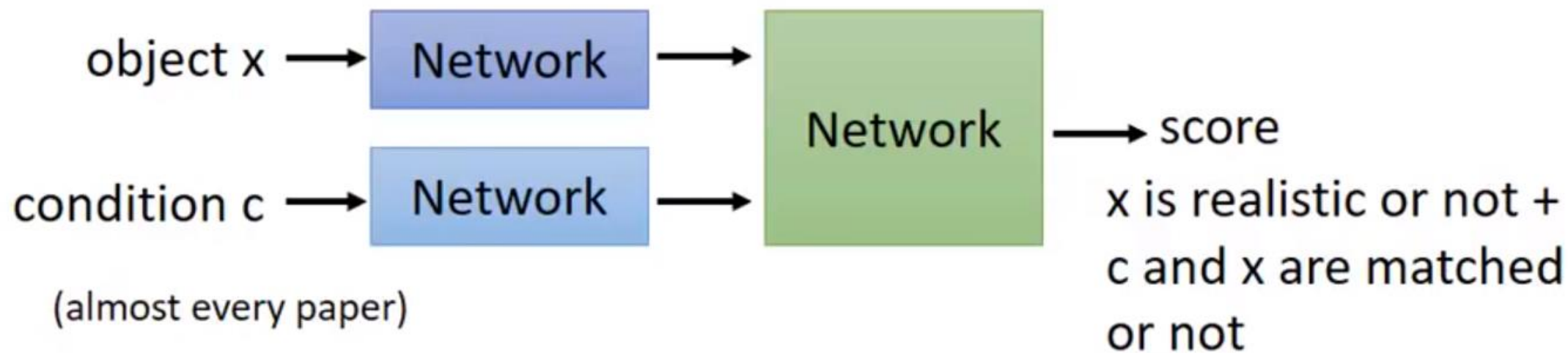
输出：生成数据 $x_1$

判别器

输入：生成数据 $x_1$ ，真实  
数据 $x_2$ ，标签向量 $c$

输出：标量Scalar

- In each training iteration:
  - Sample  $m$  positive examples  $\{(c^1, x^1), (c^2, x^2), \dots, (c^m, x^m)\}$  from database
  - Sample  $m$  noise samples  $\{z^1, z^2, \dots, z^m\}$  from a distribution
  - Obtaining generated data  $\{\tilde{x}^1, \tilde{x}^2, \dots, \tilde{x}^m\}$ ,  $\tilde{x}^i = G(c^i, z^i)$
  - Sample  $m$  objects  $\{\hat{x}^1, \hat{x}^2, \dots, \hat{x}^m\}$  from database
  - Update discriminator parameters  $\theta_d$  to maximize
    - $\tilde{V} = \frac{1}{m} \sum_{i=1}^m \log D(c^i, x^i)$   
 $+ \frac{1}{m} \sum_{i=1}^m \log (1 - D(c^i, \tilde{x}^i)) + \frac{1}{m} \sum_{i=1}^m \log (1 - D(c^i, \hat{x}^i))$
    - $\theta_d \leftarrow \theta_d + \eta \nabla \tilde{V}(\theta_d)$
  - Sample  $m$  noise samples  $\{z^1, z^2, \dots, z^m\}$  from a distribution
  - Sample  $m$  conditions  $\{c^1, c^2, \dots, c^m\}$  from a database
  - Update generator parameters  $\theta_g$  to maximize
    - $\tilde{V} = \frac{1}{m} \sum_{i=1}^m \log (D(G(c^i, z^i)))$ ,  $\theta_g \leftarrow \theta_g - \eta \nabla \tilde{V}(\theta_g)$





*paired data*



blue eyes  
red hair  
short hair

Collecting anime faces  
and the description of its  
characteristics

red hair,  
green eyes



blue hair,  
red eyes



# f-GAN

为什么出现f-GAN?

$$G^* = \arg \min_G \max_D V(G, D)$$

$$V = E_{x \sim P_{data}} [\log D(x)] + E_{x \sim P_G} [\log(1 - D(x))]$$

$$G^* = \arg \min_G D_f(P_{data} || P_G)$$

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

$$= \arg \min_G \max_D V(G, D)$$



# f-divergence

用来衡量P，Q两分布之间的差异

$$D_f(P||Q) = \int_x q(x) f\left(\frac{p(x)}{q(x)}\right) dx$$

性质：凸函数且 $f(1)=0$

# 共轭函数

若一个函数 $f$ 为凸函数，则它存在一个共轭函数

$$f^*(t) = \max_{x \in \text{dom}(f)} \{xt - f(x)\}$$

此处省略一万个字!!!!!!!!!!!!!!!!!!!!

$$f^*(t) = \sup_{x \in \text{dom}(f)} \{xt - f(x)\} \longleftrightarrow f(x) = \max_{t \in \text{dom}(f^*)} \{xt - f^*(t)\}$$

$$f^*(t) = \sup_{x \in \text{dom}(f)} \{xt - f(x)\} \longleftrightarrow f(x) = \max_{t \in \text{dom}(f^*)} \{xt - f^*(t)\}$$

$$\begin{aligned} D_f(P||Q) &= \int_x q(x) f\left(\frac{p(x)}{q(x)}\right) dx \\ &= \int_x q(x) \left( \max_{t \in \text{dom}(f^*)} \left\{ \frac{p(x)}{q(x)} t - f^*(t) \right\} \right) dx \\ &\approx \max_D \int_x p(x) D(x) dx - \int_x q(x) f^*(D(x)) dx \end{aligned}$$

假设有一个D，输入是x，输出是t

$$\begin{aligned} D_f(P||Q) &\geq \int_x q(x) \left( \frac{p(x)}{q(x)} \underline{D(x)} - f^*(\underline{D(x)}) \right) dx \\ &= \int_x p(x) D(x) dx - \int_x q(x) f^*(D(x)) dx \end{aligned}$$

$$\begin{aligned}
D_f(P||Q) &\approx \max_D \int_x p(x)D(x)dx - \int_x q(x)f^*(D(x))dx \\
&= \max_D \{ \underbrace{E_{x \sim P}[D(x)]}_{\text{Samples from P}} - \underbrace{E_{x \sim Q}[f^*(D(x))]}_{\text{Samples from Q}} \}
\end{aligned}$$

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

$$G^* = \arg \min_G D_f(P_{data}||P_G)$$

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

$$= \arg \min_G \max_D V(G, D)$$

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

Name	$D_f(P  Q)$	Generator $f(u)$
Total variation	$\frac{1}{2} \int  p(x) - q(x)  \, dx$	$\frac{1}{2} u - 1 $
Kullback-Leibler	$\int p(x) \log \frac{p(x)}{q(x)} \, dx$	$u \log u$
Reverse Kullback-Leibler	$\int q(x) \log \frac{q(x)}{p(x)} \, dx$	$-\log u$
Pearson $\chi^2$	$\int \frac{(q(x)-p(x))^2}{p(x)} \, dx$	$(u - 1)^2$
Neyman $\chi^2$	$\int \frac{(p(x)-q(x))^2}{q(x)} \, dx$	$\frac{(1-u)^2}{u}$
Squared Hellinger	$\int \left( \sqrt{p(x)} - \sqrt{q(x)} \right)^2 \, dx$	$(\sqrt{u} - 1)^2$
Jeffrey	$\int (p(x) - q(x)) \log \left( \frac{p(x)}{q(x)} \right) \, dx$	$(u - 1) \log u$
Jensen-Shannon	$\frac{1}{2} \int p(x) \log \frac{2p(x)}{p(x)+q(x)} + q(x) \log \frac{2q(x)}{p(x)+q(x)} \, dx$	$-(u + 1) \log \frac{1+u}{2} + u \log u$
Jensen-Shannon-weighted	$\int p(x) \pi \log \frac{p(x)}{\pi p(x)+(1-\pi)q(x)} + (1 - \pi)q(x) \log \frac{q(x)}{\pi p(x)+(1-\pi)q(x)} \, dx$	$\pi u \log u - (1 - \pi + \pi u) \log(1 - \pi + \pi u)$
GAN	$\int p(x) \log \frac{2p(x)}{p(x)+q(x)} + q(x) \log \frac{2q(x)}{p(x)+q(x)} \, dx - \log(4)$	$u \log u - (u + 1) \log(u + 1)$

Using the f-divergence  
you like ☺

<https://arxiv.org/pdf/1606.00709.pdf>

Name	Conjugate $f^*(t)$
Total variation	$t$
Kullback-Leibler (KL)	$\exp(t - 1)$
Reverse KL	$-1 - \log(-t)$
Pearson $\chi^2$	$\frac{1}{4}t^2 + t$
Neyman $\chi^2$	$2 - 2\sqrt{1 - t}$
Squared Hellinger	$\frac{t}{1-t}$
Jeffrey	$W(e^{1-t}) + \frac{1}{W(e^{1-t})} + t - 2$
Jensen-Shannon	$-\log(2 - \exp(t))$
Jensen-Shannon-weighted	$(1 - \pi) \log \frac{1-\pi}{1-\pi e^{t/\pi}}$
GAN	$-\log(1 - \exp(t))$

# Wasserstein GAN

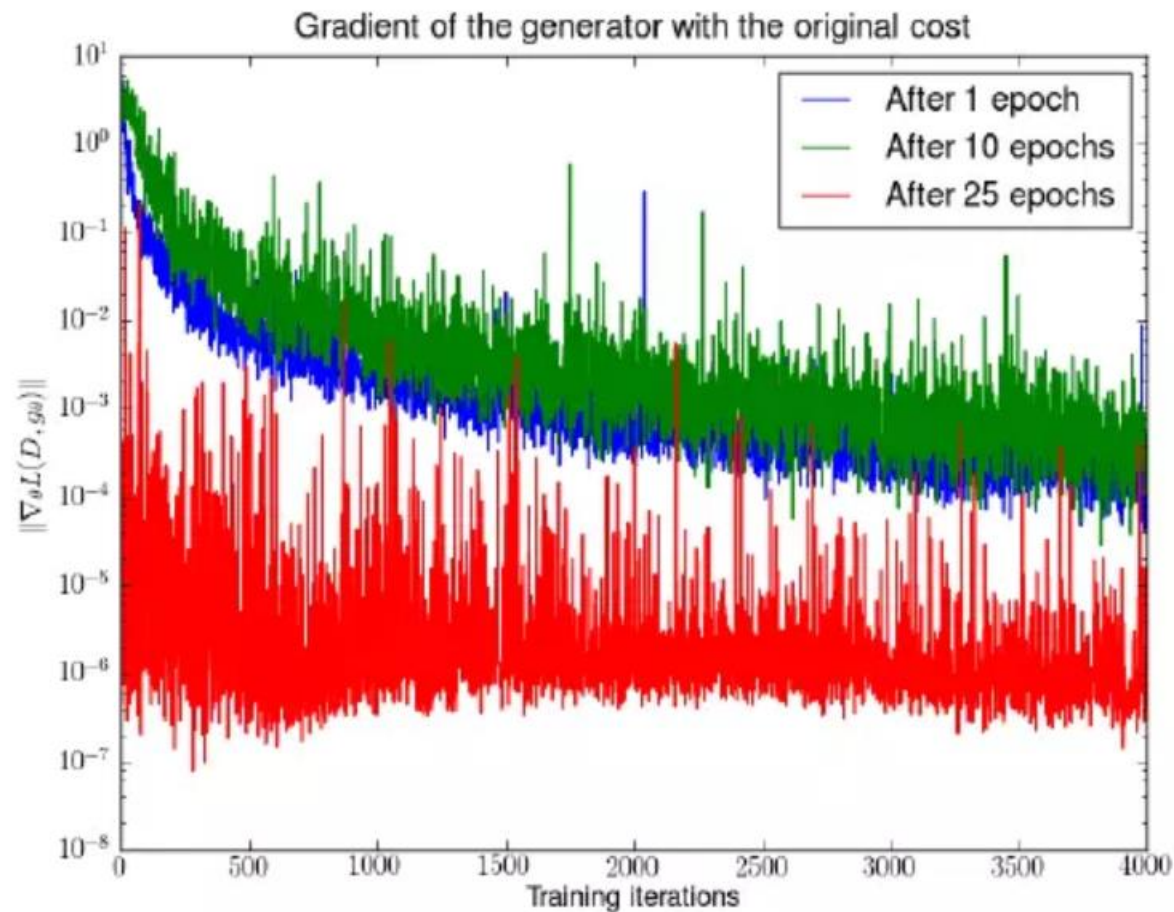
实质是什么？

将JS散度变成EM距离（ Wasserstein 距离）

# JS散度的问题

$$\max_D V(G, D) = -2\log 2 + 2JSD(P_{data}(x) || P_G(x))$$

梯度消失的问题



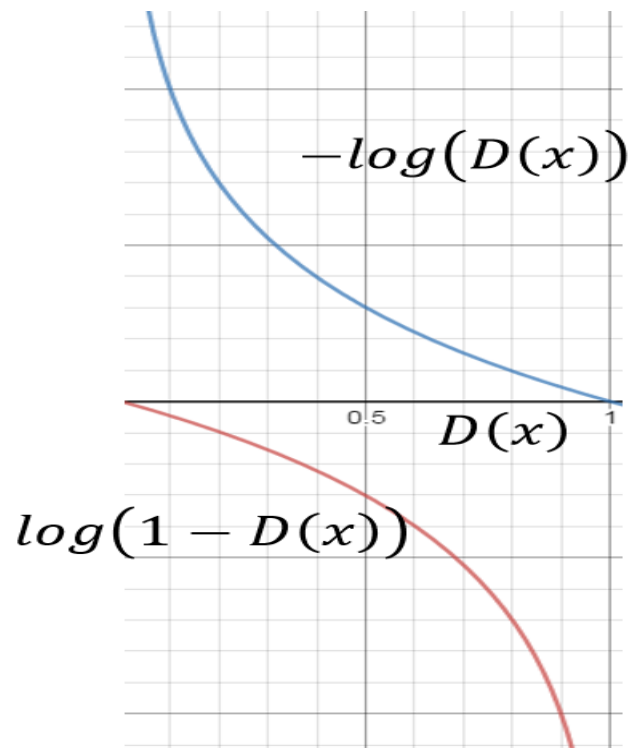


$$V = \cancel{E_{x \sim P_{data}} [\log D(x)]} \\ + E_{x \sim P_G} [\log(1 - D(x))]$$

Slow at the beginning

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

Real implementation:  
label  $x$  from  $P_G$  as positive



$$\mathbb{E}_{z \sim p(z)} [-\nabla_{\theta} \log D^*(g_{\theta}(z)) |_{\theta=\theta_0}] = \nabla_{\theta} [KL(\mathbb{P}_{g_{\theta}} \parallel \mathbb{P}_r) - 2JSD(\mathbb{P}_{g_{\theta}} \parallel \mathbb{P}_r)] |_{\theta=\theta_0}$$

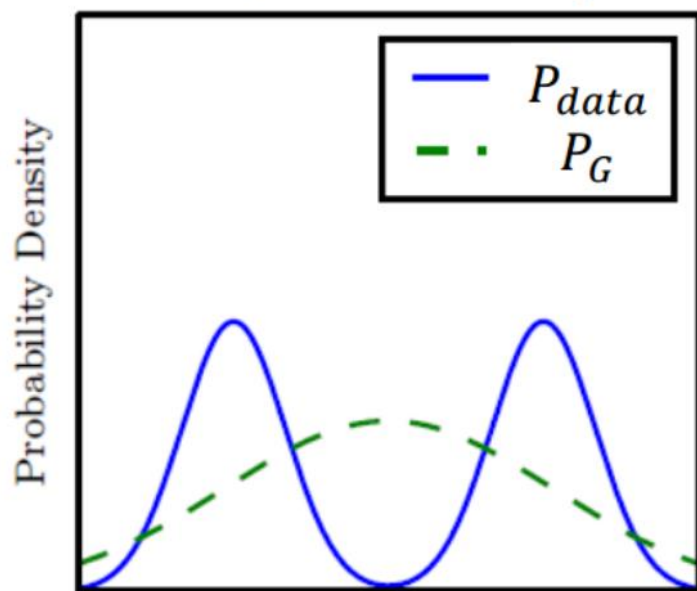
1. 第一项的KL散度会被最小化，这会带来严重的mode collapse问题。
2. 第二项意味着最大化真实数据分布和生成数据分布之间的JS散度，也就是让两者差异化更大，这显然违背了最初的优化目标，算是一种缺陷
3. 训练不稳定

# Mode Collapse

Generator会生成大量高质量却缺乏多样性的样本

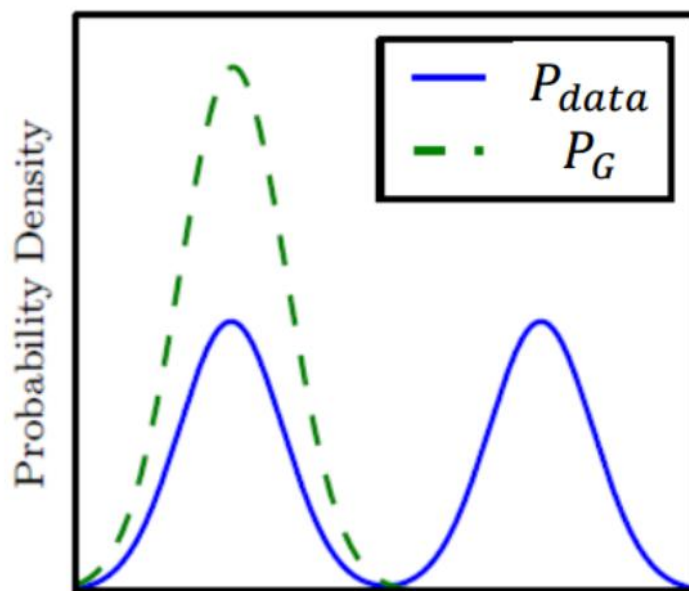
$$\nabla_{\theta} [KL(\mathbb{P}_{g_{\theta}} || \mathbb{P}_r)]$$

$$KL = \int P_{data} \log \frac{P_{data}}{P_G} dx$$



Maximum likelihood  
(minimize  $KL(P_{data} || P_G)$ )

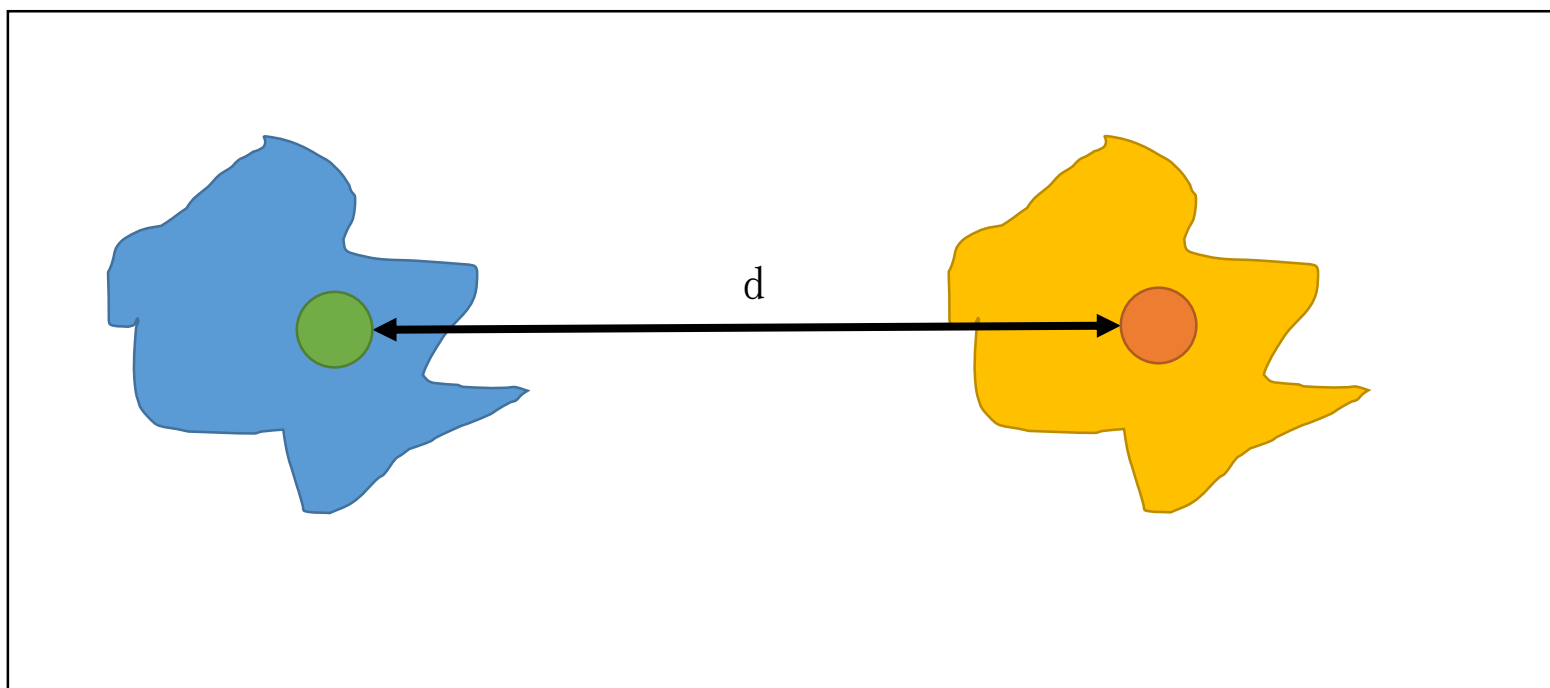
$$\text{Reverse KL} = \int P_G \log \frac{P_G}{P_{data}} dx$$



Minimize  $KL(P_G || P_{data})$   
(reverse KL)

# EM距离 earth mover ( Wasserstein 距离 )

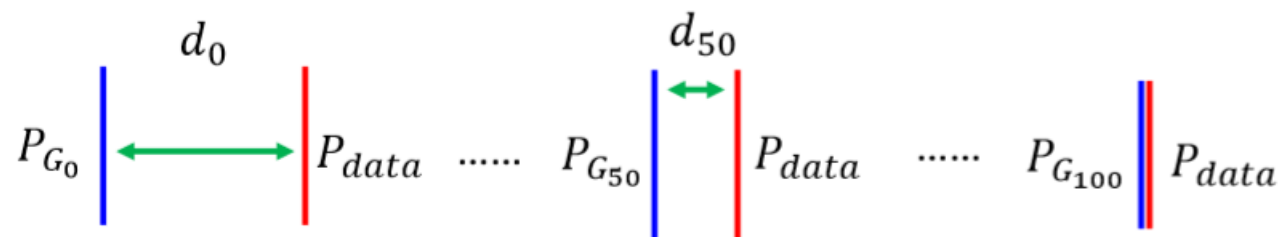
P, Q分布很远几乎无重叠的情况，仍能反映两个分布的远近



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



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



$$JS(P_{G_0}, P_{data}) \\ = \log 2$$

$$JS(P_{G_{50}}, P_{data}) \\ = \log 2$$

$$JS(P_{G_{100}}, P_{data}) \\ = 0$$

$$W(P_{G_0}, P_{data}) \\ = d_0$$

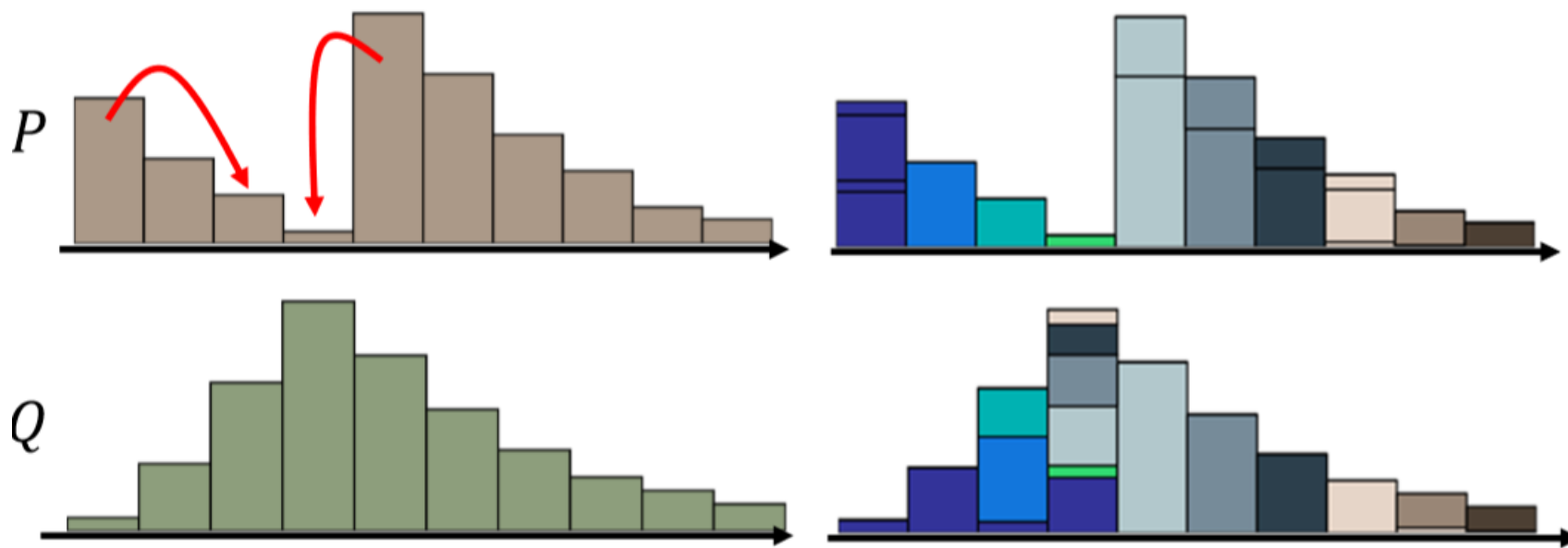
$$W(P_{G_{50}}, P_{data}) \\ = d_{50}$$

$$W(P_{G_{100}}, P_{data}) \\ = 0$$

# 形象表示

## Earth Mover's Distance

Best “moving plans”  
of this example



# 定义式

$$\mathcal{W}[p, q] = \inf_{\gamma \in \Pi[p, q]} \iint \gamma(\mathbf{x}, \mathbf{y}) \underline{d(\mathbf{x}, \mathbf{y})} d\mathbf{x} d\mathbf{y}$$

$d(\mathbf{x}, \mathbf{y})$ ，它不一定是距离，其准确含义应该是一个成本函数，代表着从  $\mathbf{x}$  运输到  $\mathbf{y}$  的成本

# 定义式

$$\mathcal{W}[p, q] = \inf_{\gamma \in \Pi[p, q]} \iint \gamma(\mathbf{x}, \mathbf{y}) d(\mathbf{x}, \mathbf{y}) d\mathbf{x} d\mathbf{y}$$

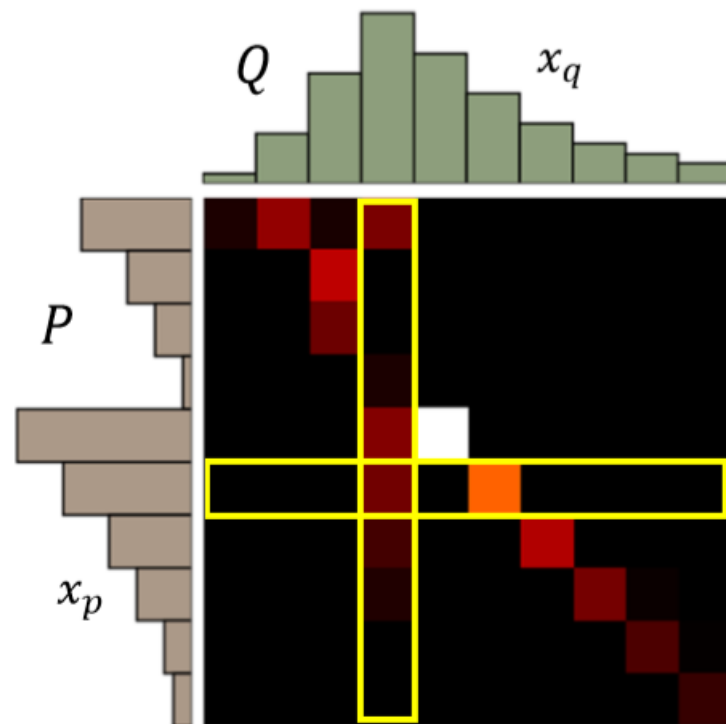
$$\int \gamma(\mathbf{x}, \mathbf{y}) d\mathbf{y} = p(\mathbf{x}) \quad \text{且} \quad \int \gamma(\mathbf{x}, \mathbf{y}) d\mathbf{x} = q(\mathbf{y})$$

$\gamma$  是一个联合分布，它的边缘分布就是原来的  $p$  和  $q$



$$\mathcal{W}[p, q] = \inf_{\gamma \in \Pi[p, q]} \iint \gamma(\mathbf{x}, \mathbf{y}) d(\mathbf{x}, \mathbf{y}) d\mathbf{x} d\mathbf{y}$$

下确界，简单来说就是取最小，也就是说，要从所有的运输方案中，找出总运输成本  $\iint \gamma(\mathbf{x}, \mathbf{y}) d(\mathbf{x}, \mathbf{y}) d\mathbf{x} d\mathbf{y}$  最小的方案，这个方案的成本，就是我们要算的  $\mathcal{W}[p, q]$



A "moving plan" is a matrix

The value of the element is the amount of earth from one position to another.

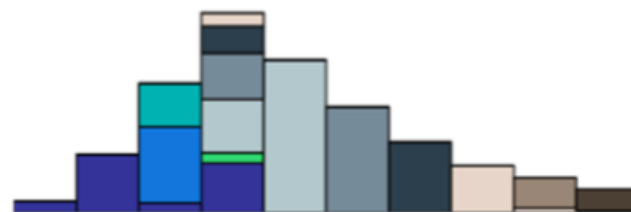
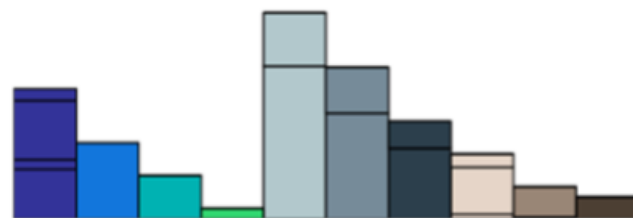
Average distance of a plan  $\gamma$ :

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

Earth Mover's Distance:

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

The best plan



$$\mathcal{W}[p, q] = \inf_{\gamma \in \Pi[p, q]} \iint \gamma(\mathbf{x}, \mathbf{y}) d(\mathbf{x}, \mathbf{y}) d\mathbf{x} d\mathbf{y}$$

其实就是最小化

$$\iint \gamma(\mathbf{x}, \mathbf{y}) \underline{d(\mathbf{x}, \mathbf{y})} d\mathbf{x} d\mathbf{y}$$

$d(\mathbf{x}, \mathbf{y})$  是定值

需要满足约束:  $\int \gamma(\mathbf{x}, \mathbf{y}) d\mathbf{y} = p(\mathbf{x}), \quad \int \gamma(\mathbf{x}, \mathbf{y}) d\mathbf{x} = q(\mathbf{y}), \quad \gamma(\mathbf{x}, \mathbf{y}) \geq 0$

积分只是求和的极限形式,  
所以我们可以把  $\gamma(\mathbf{x}, \mathbf{y})$  和  
 $d(\mathbf{x}, \mathbf{y})$  离散化

相当于就是将  $\Gamma$  和  $D$  对应  
位置相乘, 然后求和, 这  
不就是内积  $\langle \Gamma, D \rangle$  了吗

$$\Gamma = \begin{pmatrix} \gamma(\mathbf{x}_1, \mathbf{y}_1) \\ \gamma(\mathbf{x}_1, \mathbf{y}_2) \\ \vdots \\ \gamma(\mathbf{x}_2, \mathbf{y}_1) \\ \gamma(\mathbf{x}_2, \mathbf{y}_2) \\ \vdots \\ \vdots \\ \gamma(\mathbf{x}_n, \mathbf{y}_1) \\ \gamma(\mathbf{x}_n, \mathbf{y}_2) \\ \vdots \\ \vdots \end{pmatrix}, \quad D = \begin{pmatrix} d(\mathbf{x}_1, \mathbf{y}_1) \\ d(\mathbf{x}_1, \mathbf{y}_2) \\ \vdots \\ d(\mathbf{x}_2, \mathbf{y}_1) \\ d(\mathbf{x}_2, \mathbf{y}_2) \\ \vdots \\ \vdots \\ d(\mathbf{x}_n, \mathbf{y}_1) \\ d(\mathbf{x}_n, \mathbf{y}_2) \\ \vdots \\ \vdots \end{pmatrix}$$

$$\int \gamma(\mathbf{x}, \mathbf{y}) d\mathbf{y} = p(\mathbf{x}), \quad \int \gamma(\mathbf{x}, \mathbf{y}) d\mathbf{x} = q(\mathbf{y}), \quad \gamma(\mathbf{x}, \mathbf{y}) \geq 0$$

约束条件也可以写成矩阵形式  $A\Gamma = \mathbf{b}$

$$\underbrace{\begin{pmatrix} 1 & 1 & \dots & 0 & 0 & \dots & \dots & 0 & 0 & \dots & \dots \\ 0 & 0 & \dots & 1 & 1 & \dots & \dots & 0 & 0 & \dots & \dots \\ \vdots & \vdots & \ddots & \vdots & \vdots & \ddots & \ddots & \vdots & \vdots & \ddots & \ddots \\ 0 & 0 & \dots & 0 & 0 & \dots & \dots & 1 & 1 & \dots & \dots \\ \vdots & \vdots & \ddots & \vdots & \vdots & \ddots & \ddots & \vdots & \vdots & \ddots & \ddots \\ \hline 1 & 0 & \dots & 1 & 0 & \dots & \dots & 1 & 0 & \dots & \dots \\ 0 & 1 & \dots & 0 & 1 & \dots & \dots & 0 & 1 & \dots & \dots \\ \vdots & \vdots & \ddots & \vdots & \vdots & \ddots & \ddots & \vdots & \vdots & \ddots & \ddots \\ 0 & 0 & \dots & 0 & 0 & \dots & \dots & 0 & 0 & \dots & \dots \\ \vdots & \vdots & \ddots & \vdots & \vdots & \ddots & \ddots & \vdots & \vdots & \ddots & \ddots \end{pmatrix}}_A \underbrace{\begin{pmatrix} \gamma(\mathbf{x}_1, \mathbf{y}_1) \\ \gamma(\mathbf{x}_1, \mathbf{y}_2) \\ \vdots \\ \gamma(\mathbf{x}_2, \mathbf{y}_1) \\ \gamma(\mathbf{x}_2, \mathbf{y}_2) \\ \vdots \\ \vdots \\ \vdots \\ \gamma(\mathbf{x}_n, \mathbf{y}_1) \\ \gamma(\mathbf{x}_n, \mathbf{y}_2) \\ \vdots \\ \vdots \end{pmatrix}}_{\Gamma} = \underbrace{\begin{pmatrix} p(\mathbf{x}_1) \\ p(\mathbf{x}_2) \\ \vdots \\ p(\mathbf{x}_n) \\ \vdots \\ \vdots \\ \vdots \\ q(\mathbf{y}_1) \\ q(\mathbf{y}_2) \\ \vdots \\ q(\mathbf{y}_n) \\ \vdots \end{pmatrix}}_{\mathbf{b}} \longrightarrow \min_{\Gamma} \{ \langle \Gamma, D \rangle \mid A\Gamma = \mathbf{b}, \Gamma \geq 0 \}$$

线性约束下的线性函数最小值

# 线性约束下的线性函数最小值

$$\min_x \{c^\top x \mid Ax = b, x \geq 0\}$$

弱对偶形式

设置最小值在  $x^*$  取到

$$\text{两边乘以一个 } y^\top \in R^m \longrightarrow y^\top Ax^* \in y^\top b$$

$$\begin{array}{ccc} \text{此时假设 } y^\top A \leq c^\top & \longrightarrow & y^\top Ax^* \leq c^\top x^* \\ & & \nwarrow \\ & & y^\top b \leq c^\top x^* \end{array} \quad \text{因为 } x^* > 0$$

$$\text{在条件下 } y^\top A \leq c^\top \longrightarrow \max_y \{b^\top y \mid A^\top y \leq c\} \leq \min_x \{c^\top x \mid Ax = b, x \geq 0\}$$

现在我们将原来的最小值问题变成了一个最大值问题，这便有了对偶的味道。  
从应用角度，其实弱对偶形式给出的下界都已经够用了，  
因为深度学习中的问题都很复杂，能有一个近似的目标去优化都已经很不错了。

$$\max_{\mathbf{y}} \{ \mathbf{b}^\top \mathbf{y} \mid \mathbf{A}^\top \mathbf{y} \leq \mathbf{c} \} \leq \min_{\mathbf{x}} \{ \mathbf{c}^\top \mathbf{x} \mid \mathbf{A} \mathbf{x} = \mathbf{b}, \mathbf{x} \geq 0 \}$$

$$\min_{\mathbf{\Gamma}} \{ \langle \mathbf{\Gamma}, \mathbf{D} \rangle \mid \mathbf{A} \mathbf{\Gamma} = \mathbf{b}, \mathbf{\Gamma} \geq 0 \} = \max_{\mathbf{F}} \{ \langle \mathbf{b}, \mathbf{F} \rangle \mid \mathbf{A}^\top \mathbf{F} \leq \mathbf{D} \}$$

b 是由两部分拼起来的，所以我們也可以把 F 类似地写成

$$\underbrace{\begin{pmatrix} 1 & 1 & \dots & 0 & 0 & \dots & \dots & 0 & 0 & \dots & \dots \\ 0 & 0 & \dots & 1 & 1 & \dots & \dots & 0 & 0 & \dots & \dots \\ \vdots & \vdots & \ddots & \vdots & \vdots & \ddots & \ddots & \vdots & \vdots & \ddots & \ddots \\ 0 & 0 & \dots & 0 & 0 & \dots & \dots & 1 & 1 & \dots & \dots \\ \vdots & \vdots & \ddots & \vdots & \vdots & \ddots & \ddots & \vdots & \vdots & \ddots & \ddots \\ \hline 1 & 0 & \dots & 1 & 0 & \dots & \dots & 1 & 0 & \dots & \dots \\ 0 & 1 & \dots & 0 & 1 & \dots & \dots & 0 & 1 & \dots & \dots \\ \vdots & \vdots & \ddots & \vdots & \vdots & \ddots & \ddots & \vdots & \vdots & \ddots & \ddots \\ 0 & 0 & \dots & 0 & 0 & \dots & \dots & 0 & 0 & \dots & \dots \\ \vdots & \vdots & \ddots & \vdots & \vdots & \ddots & \ddots & \vdots & \vdots & \ddots & \ddots \end{pmatrix}}_{\mathbf{A}} \underbrace{\begin{pmatrix} \gamma(\mathbf{x}_1, \mathbf{y}_1) \\ \gamma(\mathbf{x}_1, \mathbf{y}_2) \\ \vdots \\ \gamma(\mathbf{x}_2, \mathbf{y}_1) \\ \gamma(\mathbf{x}_2, \mathbf{y}_2) \\ \vdots \\ \gamma(\mathbf{x}_n, \mathbf{y}_1) \\ \gamma(\mathbf{x}_n, \mathbf{y}_2) \\ \vdots \end{pmatrix}}_{\mathbf{\Gamma}} = \underbrace{\begin{pmatrix} p(\mathbf{x}_1) \\ p(\mathbf{x}_2) \\ \vdots \\ p(\mathbf{x}_n) \\ \vdots \\ q(\mathbf{y}_1) \\ q(\mathbf{y}_2) \\ \vdots \\ q(\mathbf{y}_n) \\ \vdots \end{pmatrix}}_{\mathbf{b}} \quad \mathbf{F} = \begin{pmatrix} f(\mathbf{x}_1) \\ f(\mathbf{x}_2) \\ \vdots \\ f(\mathbf{x}_n) \\ \vdots \\ g(\mathbf{y}_1) \\ g(\mathbf{y}_2) \\ \vdots \\ g(\mathbf{y}_n) \\ \vdots \end{pmatrix}$$

现在  $\langle \mathbf{b}, \mathbf{F} \rangle$   
是什么？

$$\langle \mathbf{b}, \mathbf{F} \rangle = \sum_n p(\mathbf{x}_n) f(\mathbf{x}_n) + \sum_n q(\mathbf{x}_n) g(\mathbf{x}_n)$$

$$\langle \mathbf{b}, \mathbf{F} \rangle = \sum_n p(\mathbf{x}_n) f(\mathbf{x}_n) + \sum_n q(\mathbf{x}_n) g(\mathbf{x}_n) \longrightarrow \langle \mathbf{b}, \mathbf{F} \rangle = \int [p(\mathbf{x}) f(\mathbf{x}) + q(\mathbf{x}) g(\mathbf{x})] d\mathbf{x}$$

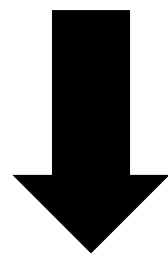
约束条件:

$$\underbrace{\begin{pmatrix} 1 & 0 & \dots & 0 & \dots & | & 1 & 0 & \dots & 0 & \dots \\ 1 & 0 & \dots & 0 & \dots & | & 0 & 1 & \dots & 0 & \dots \\ \vdots & \vdots & \ddots & \vdots & \ddots & | & \vdots & \vdots & \ddots & \vdots & \ddots \\ \hline 0 & 1 & \dots & 0 & \dots & | & 1 & 0 & \dots & 0 & \dots \\ 0 & 1 & \dots & 0 & \dots & | & 0 & 1 & \dots & 0 & \dots \\ \vdots & \vdots & \ddots & \vdots & \ddots & | & \vdots & \vdots & \ddots & \vdots & \ddots \\ \hline \vdots & \vdots & \ddots & \vdots & \ddots & | & \vdots & \vdots & \ddots & \vdots & \ddots \\ \hline 0 & 0 & \dots & 1 & \dots & | & 1 & 0 & \dots & 0 & \dots \\ 0 & 0 & \dots & 1 & \dots & | & 0 & 1 & \dots & 0 & \dots \\ \vdots & \vdots & \ddots & \vdots & \ddots & | & \vdots & \vdots & \ddots & \vdots & \ddots \\ \hline \vdots & \vdots & \ddots & \vdots & \ddots & | & \vdots & \vdots & \ddots & \vdots & \ddots \end{pmatrix}}_{\mathbf{A}^\top} \underbrace{\begin{pmatrix} f(\mathbf{x}_1) \\ f(\mathbf{x}_2) \\ \vdots \\ f(\mathbf{x}_n) \\ \vdots \\ g(\mathbf{y}_1) \\ g(\mathbf{y}_2) \\ \vdots \\ g(\mathbf{y}_n) \\ \vdots \end{pmatrix}}_{\mathbf{F}} \leq \underbrace{\begin{pmatrix} d(\mathbf{x}_1, \mathbf{y}_1) \\ d(\mathbf{x}_1, \mathbf{y}_2) \\ \vdots \\ d(\mathbf{x}_2, \mathbf{y}_1) \\ d(\mathbf{x}_2, \mathbf{y}_2) \\ \vdots \\ \vdots \\ d(\mathbf{x}_n, \mathbf{y}_1) \\ d(\mathbf{x}_n, \mathbf{y}_2) \\ \vdots \\ \vdots \end{pmatrix}}_{\mathbf{D}}$$

实际上就是:  $\forall i, j, f(\mathbf{x}_i) + g(\mathbf{y}_j) \leq d(\mathbf{x}_i, \mathbf{y}_j) \longrightarrow \forall \mathbf{x}, \mathbf{y}, f(\mathbf{x}) + g(\mathbf{y}) \leq d(\mathbf{x}, \mathbf{y})$



$$\mathcal{W}[p, q] = \inf_{\gamma \in \Pi[p, q]} \iint \gamma(\mathbf{x}, \mathbf{y}) d(\mathbf{x}, \mathbf{y}) d\mathbf{x} d\mathbf{y}$$



$$\mathcal{W}[p, q] = \max_{f, g} \left\{ \int [p(\mathbf{x})f(\mathbf{x}) + q(\mathbf{x})g(\mathbf{x})] d\mathbf{x} \mid f(\mathbf{x}) + g(\mathbf{y}) \leq d(\mathbf{x}, \mathbf{y}) \right\}$$

$$\mathcal{W}[p, q] = \max_{f, g} \left\{ \underbrace{\int [p(\mathbf{x})f(\mathbf{x}) + q(\mathbf{x})g(\mathbf{x})] d\mathbf{x}}_{\downarrow} \left| \underbrace{f(\mathbf{x}) + g(\mathbf{y}) \leq d(\mathbf{x}, \mathbf{y})}_{\downarrow} \right. \right\}$$

$$\begin{aligned} p(\mathbf{x})f(\mathbf{x}) + q(\mathbf{x})g(\mathbf{x}) &\leq p(\mathbf{x})f(\mathbf{x}) + q(\mathbf{x})[-f(\mathbf{x})] \\ &= p(\mathbf{x})f(\mathbf{x}) - q(\mathbf{x})f(\mathbf{x}) \end{aligned} \quad \longleftarrow f(\mathbf{x}) + g(\mathbf{x}) \leq d(\mathbf{x}, \mathbf{x}) = 0$$

$$\mathcal{W}[p, q] = \max_f \left\{ \int [p(\mathbf{x})f(\mathbf{x}) - q(\mathbf{x})f(\mathbf{x})] d\mathbf{x} \left| f(\mathbf{x}) - f(\mathbf{y}) \leq d(\mathbf{x}, \mathbf{y}) \right. \right\}$$

$$\mathcal{W}[p, q] = \max_{f, \|f\|_L \leq 1} \mathbb{E}_{\mathbf{x} \sim p(\mathbf{x})}[f(\mathbf{x})] - \mathbb{E}_{\mathbf{x} \sim q(\mathbf{x})}[f(\mathbf{x})] \implies \min_G \max_{f, \|f\|_L \leq 1} \mathbb{E}_{\mathbf{x} \sim p(\mathbf{x})}[f(\mathbf{x})] - \mathbb{E}_{\mathbf{z} \sim q(\mathbf{z})}[f(G(\mathbf{z}))]$$

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

### **Lipschitz Function**

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

Output  
change

Input  
change

K=1 for "1 - Lipschitz"

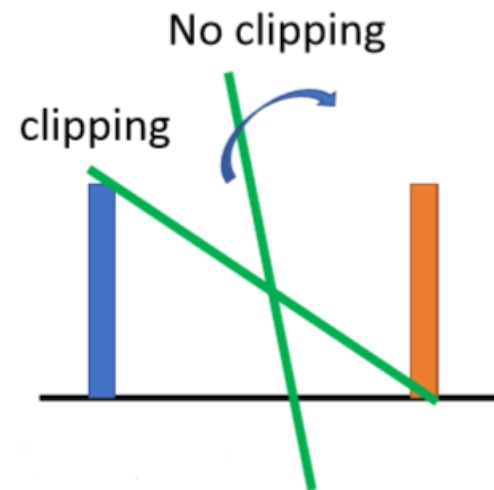
Do not change fast

# 问题

在计算梯度的时候，怎么解决约束问题？

(1) weight clipping

强制把更新后的参数固定在 $[-c, c]$



if  $w > c$ , then  $w=c$ ; if  $w < -c$ , then  $w=-c$

## Algorithm of WGAN

- In each training iteration:

No sigmoid for the output of D

Learning  
D

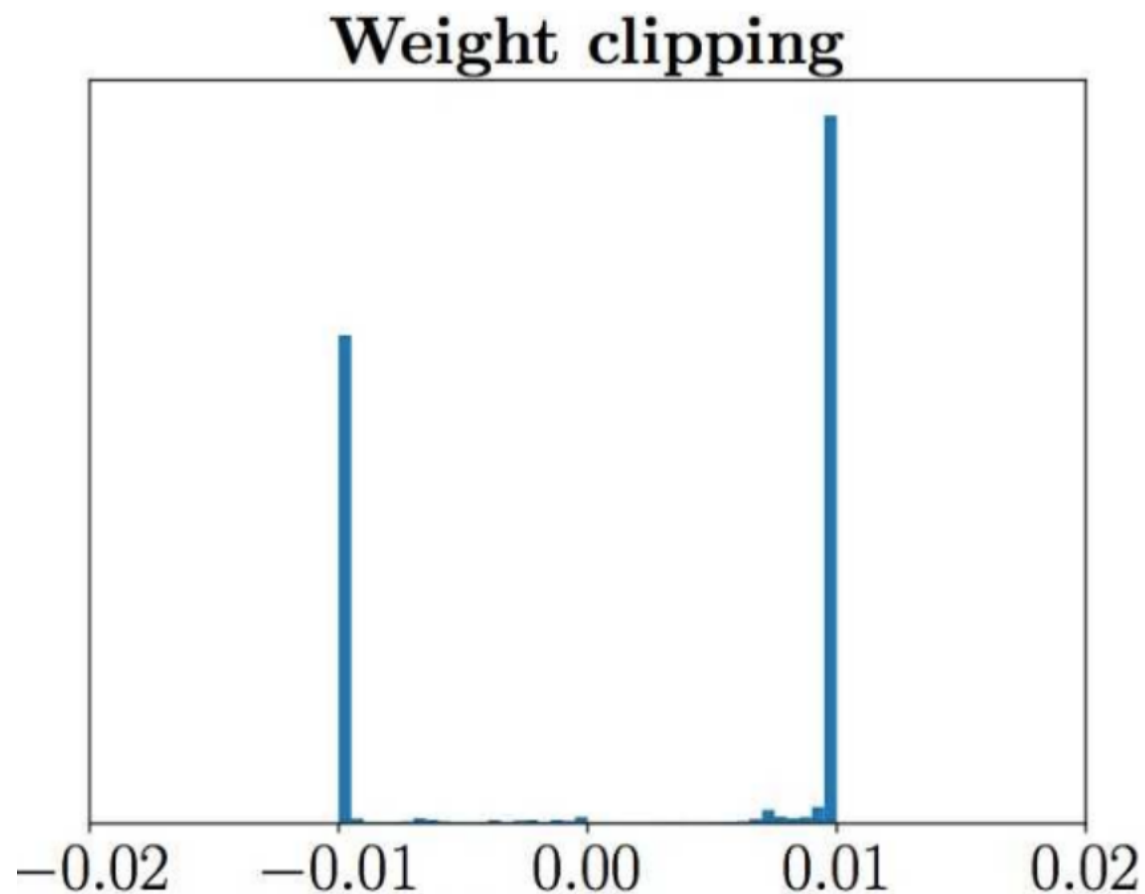
Repeat  
k times

- Sample  $m$  examples  $\{x^1, x^2, \dots, x^m\}$  from data distribution  $P_{data}(x)$
- Sample  $m$  noise samples  $\{z^1, z^2, \dots, z^m\}$  from the prior  $P_{prior}(z)$
- Obtaining generated data  $\{\tilde{x}^1, \tilde{x}^2, \dots, \tilde{x}^m\}$ ,  $\tilde{x}^i = G(z^i)$
- 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  $\{z^1, z^2, \dots, z^m\}$  from the prior  $P_{prior}(z)$
- Update generator parameters  $\theta_g$  to minimize
  - $\tilde{V} = \frac{1}{m} \sum_{i=1}^m \log D(x^i) - \frac{1}{m} \sum_{i=1}^m D(G(z^i))$
  - $\theta_g \leftarrow \theta_g - \eta \nabla \tilde{V}(\theta_g)$



很容易一不小心就梯度消失或者梯度爆炸。原因是判别器是一个多层网络，如果把 clipping threshold 设得稍微小了一点，每经过一层网络，梯度就变小一点点，多层之后就会指数衰减；反之，如果设得稍微大了一点，每经过一层网络，梯度变大一点点，多层之后就会指数爆炸。只有设得不大不小，才能让生成器获得恰到好处的回传梯度，然而在实际应用中这个平衡区域可能很狭窄，就会给调参工作带来麻烦。

## (2) WGAN-GP (gradient penalty)

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



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

让梯度尽可能的接近1



不能全部x都算，要采样

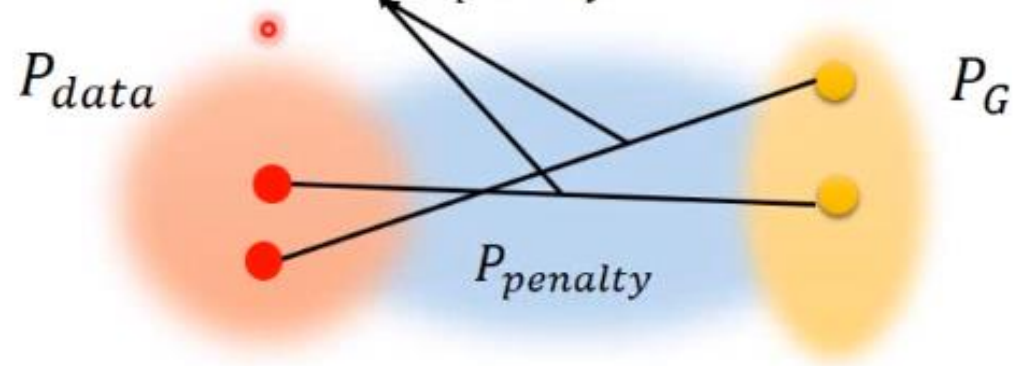
$$- \lambda E_{x \sim P_{penalty}} [\max(0, \|\nabla_x D(x)\| - 1)]$$



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



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



谢谢！