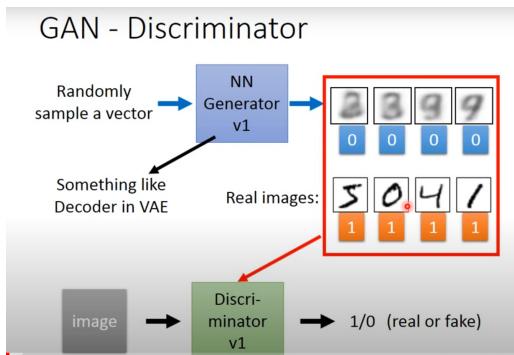
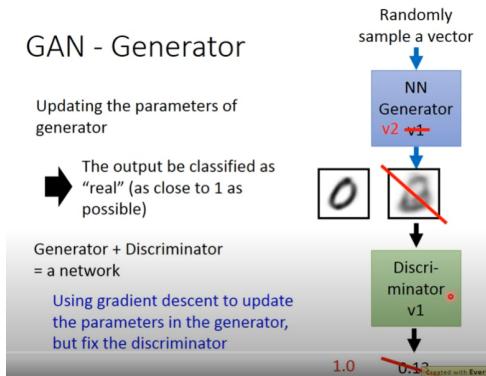


## GAN

flow to train Discriminator?



## How to train Generator?



## 数学原理：

## Maximum Likelihood Estimation

- Give a data distribution  $P_{\text{data}}(x)$
  - We have a distribution  $P_G(x; \theta)$  parameterized by  $\theta$ 
    - E.g.  $P_G(x; \theta)$  is a GMM,  $\theta$  are means and variances of Gaussian
    - We want to find  $\theta$  such that  $P_G(x; \theta)$  close  $\rightarrow P_{\text{data}}(x)$
  - Sample  $\{x^1, x^2, \dots, x^m\}$  from  $P_{\text{data}}(x)$
  - We can compute  $P_G(x^i; \theta)$
  - Likelihood of generating the samples:  $L = \prod_{i=1}^m P_G(x^i; \theta)$
  - Find  $\theta^*$  maximizing the likelihood.

$$\theta^* = \arg \max_{\theta} \sum_{i=1}^m p_G(x^i, \theta)$$

$$= \arg \max_{\theta} \sum_{i=1}^m \log p_{\theta}(x^i, \theta) \quad \{x^1, x^2, \dots, x^m\} \text{ from } p_{\text{data}}(x)$$

$$\approx \arg \max_{\theta} E_{x \sim P_{\text{data}}} [\log P_G(x; \theta)]$$

$$= \arg \max_Q \int_x P_{\text{data}}(x) \log P_Q(x^i, Q) dx$$

与无关，所以不能用

$$= \underset{\theta}{\operatorname{argmax}} \int p_{\text{data}}(x) \log p_G(x^i; \theta) dx - \underline{\int x p_{\text{data}}(x) \log p_{\text{data}}(x) dx}$$

$$= \arg \max_{\theta} \sum_x [P_{\text{data}}(x) \log \frac{P_{\theta}(x^i, \theta)}{P_{\text{data}}(x)}] dx$$

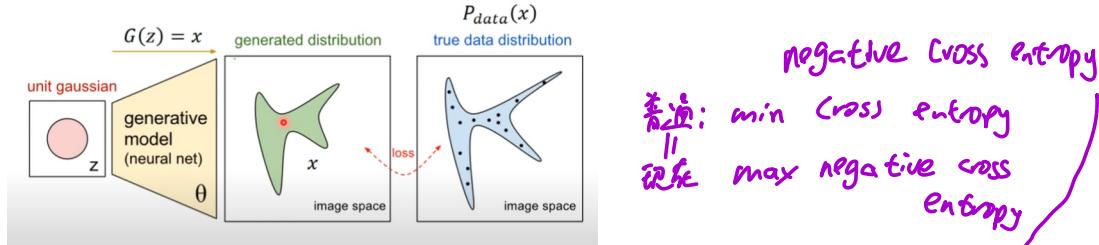
$$= \underset{\theta}{\operatorname{argmin}} \int_x \left[ P_{\text{data}}(x) \log \frac{P_{\text{data}}(x)}{P_G(x; \theta)} \right] dx$$

$$= \underset{\theta}{\operatorname{argmin}} \text{KL}(P_{\text{data}}(x) || P_G(x; \theta))$$

KL 故意选小，代表逼真分布  
越相近

Question: How to have a very general  $P_G(x; \theta)$ ?

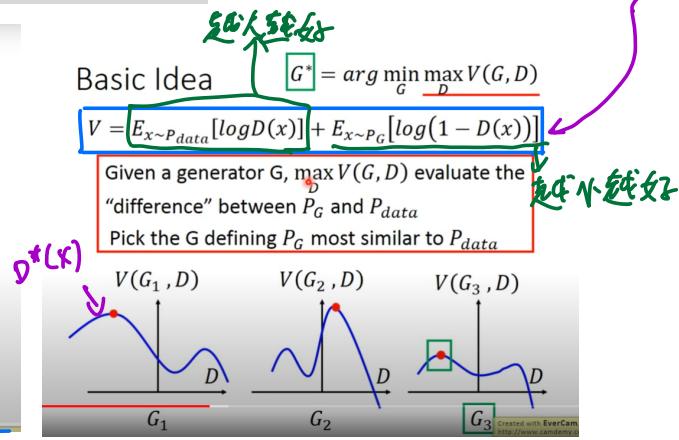
Now  $P_G(x; \theta)$  is a NN



## Basic Idea of GAN

- Generator G Hard to learn by maximum likelihood
  - G is a function, input z, output x
  - Given a prior distribution  $P_{\text{prior}}(z)$ , a probability distribution  $P_G(x)$  is defined by function G
- Discriminator D
  - D is a function, input x, output scalar
  - Evaluate the "difference" between  $P_G(x)$  and  $P_{\text{data}}(x)$
- There is a function  $V(G, D)$ .

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



Given  $G$ , maximize  $D$ .

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

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

$$= \int P_{\text{data}}(x) \log D(x) dx + \int P_G(x) \log(1 - D(x)) dx$$

$$\Rightarrow \int (P_{\text{data}}(x) \log D(x) + P_G(x) \log(1 - D(x))) dx$$

找  $D(x)$ . 全部卫最大

$$\Rightarrow a \log D(x) + b \log(1 - D(x)).$$

不容易得

$$D^*(x) = \frac{P_{\text{data}}(x)}{P_{\text{data}}(x) + P_G(x)}$$

Given  $D^*$ , minimize  $G$ :

$$\begin{aligned}
 V &= \mathbb{E}_{x \sim P_{\text{data}}} [\log D(x)] + \mathbb{E}_{x \sim P_G} [\log(1 - D(x))] \\
 &= \mathbb{E}_{x \sim P_{\text{data}}} \left[ \log \frac{P_{\text{data}}(x)}{P_{\text{data}}(x) + P_G(x)} \right] + \mathbb{E}_{x \sim P_G} \left[ \log \left( 1 - \frac{P_{\text{data}}(x)}{P_{\text{data}}(x) + P_G(x)} \right) \right] \\
 &= \int_x P_{\text{data}}(x) \left[ \log \frac{P_{\text{data}}(x)}{P_{\text{data}}(x) + P_G(x)} \right] dx + \int_x P_G(x) \left[ \log \left( 1 - \frac{P_{\text{data}}(x)}{P_{\text{data}}(x) + P_G(x)} \right) \right] dx \\
 &= \int_x P_{\text{data}}(x) \left[ \log \frac{P_{\text{data}}(x)}{(P_{\text{data}}(x) + P_G(x))/2} \right] dx + \int_x P_G(x) \left[ \log \left( 1 - \frac{P_{\text{data}}(x)}{(P_{\text{data}}(x) + P_G(x))/2} \right) \right] dx - 2\log 2 \\
 &= -2\log 2 + \text{KL} \left( P_{\text{data}}(x) \parallel \frac{P_{\text{data}}(x) + P_G(x)}{2} \right) + \text{KL} \left( P_G(x) \parallel \frac{P_{\text{data}}(x) + P_G(x)}{2} \right)
 \end{aligned}$$

**Remark:**

$$\text{JSD}(P||Q) = \frac{1}{2} D(P||M) + \frac{1}{2} D(Q||M) \quad \text{where } M = \frac{1}{2}(P+Q)$$

$$= -2\log 2 + 2\text{JSD}(P_{\text{data}}(x) \parallel P_G(x))$$

In the end .....

$$\begin{aligned}
 V &= \mathbb{E}_{x \sim P_{\text{data}}} [\log D(x)] \\
 &\quad + \mathbb{E}_{x \sim P_G} [\log(1 - D(x))]
 \end{aligned}$$

- Generator G, Discriminator D
- Looking for  $G^*$  such that

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

- Given  $G$ ,  $\max_D V(G, D)$

$$= -2\log 2 + 2\text{JSD}(P_{\text{data}}(x) \parallel P_G(x))$$

- What is the optimal G?

$$P_G(x) = P_{\text{data}}(x)$$

the value is related to JS divergence

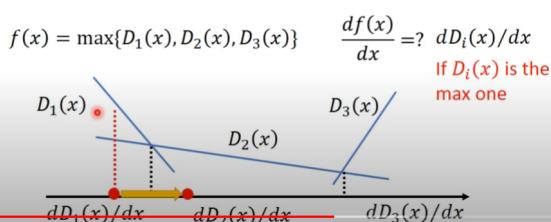
Question : **Ans** ?

Algorithm

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

- To find the best G minimizing the loss function  $L(G)$ ,

$$\theta_G \leftarrow \theta_G - \eta \partial L(G) / \partial \theta_G \quad \theta_G \text{ defines } G$$



Algorithm

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

- Given  $G_0$

- Find  $D_0^*$  maximizing  $V(G_0, D)$

$$V(G_0, D_0^*) \text{ is the JS divergence between } P_{\text{data}}(x) \text{ and } P_{G_0}(x)$$

- $\theta_G \leftarrow \theta_G - \eta \partial V(G, D_0^*) / \partial \theta_G \rightarrow$  Obtain  $G_1$  Decrease JS divergence(?)

- Find  $D_1^*$  maximizing  $V(G_1, D)$

$$V(G_1, D_1^*) \text{ is the JS divergence between } P_{\text{data}}(x) \text{ and } P_{G_1}(x)$$

- $\theta_G \leftarrow \theta_G - \eta \partial V(G, D_1^*) / \partial \theta_G \rightarrow$  Obtain  $G_2$  Decrease JS divergence(?)

.....

In practice ...

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

- Given  $G$ , how to compute  $\max_D V(G, D)$
- Sample  $\{x^1, x^2, \dots, x^m\}$  from  $P_{data}(x)$ , sample  $\{\tilde{x}^1, \tilde{x}^2, \dots, \tilde{x}^m\}$  from generator  $P_G(x)$

$$\text{Maximize } \tilde{V} = \frac{1}{m} \sum_{i=1}^m \log D(x^i) + \frac{1}{m} \sum_{i=1}^m \log(1 - D(\tilde{x}^i))$$

Binary Classifier

Output is  $D(x)$  Minimize Cross-entropy  
 If  $x$  is a positive example  $\rightarrow$  Minimize  $-\log D(x)$   
 If  $x$  is a negative example  $\rightarrow$  Minimize  $-\log(1-D(x))$

Binary Classifier

Output is  $f(x)$  Minimize Cross-entropy  
 If  $x$  is a positive example  $\rightarrow$  Minimize  $-\log f(x)$   
 If  $x$  is a negative example  $\rightarrow$  Minimize  $-\log(1-f(x))$

$D$  is a binary classifier (can be deep) with parameters  $\theta_d$

$\{x^1, x^2, \dots, x^m\}$  from  $P_{data}(x)$   $\rightarrow$  Positive examples

$\{\tilde{x}^1, \tilde{x}^2, \dots, \tilde{x}^m\}$  from  $P_G(x)$   $\rightarrow$  Negative examples

$$\text{Minimize } L = \frac{1}{m} \sum_{i=1}^m \log D(x^i) + \frac{1}{m} \sum_{i=1}^m \log(1 - D(\tilde{x}^i))$$

$$\text{Maximize } \tilde{V} = \frac{1}{m} \sum_{i=1}^m \log D(x^i) + \frac{1}{m} \sum_{i=1}^m \log(1 - D(\tilde{x}^i))$$

Created with EverCam

### Algorithm

- Initialize  $\theta_d$  for  $D$  and  $\theta_g$  for  $G$
- Can only find lower bound of  $\max_D V(G, D)$
- In each training iteration:
- 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
    - $\tilde{V} = \frac{1}{m} \sum_{i=1}^m \log D(x^i) + \frac{1}{m} \sum_{i=1}^m \log(1 - D(\tilde{x}^i))$
    - $\theta_d \leftarrow \theta_d + \eta \nabla \tilde{V}(\theta_d)$
  - 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
    - $\tilde{V} = \frac{1}{m} \sum_{i=1}^m \log D(x^i) + \frac{1}{m} \sum_{i=1}^m \log(1 - D(G(z^i)))$
    - $\theta_g \leftarrow \theta_g - \eta \nabla \tilde{V}(\theta_g)$

Learning  $D$   
 Repeat  $k$  times

Learning  $G$   
 Only Once

Implementation of GANs in TensorFlow 2.0

Background :

JS divergence is not suitable

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

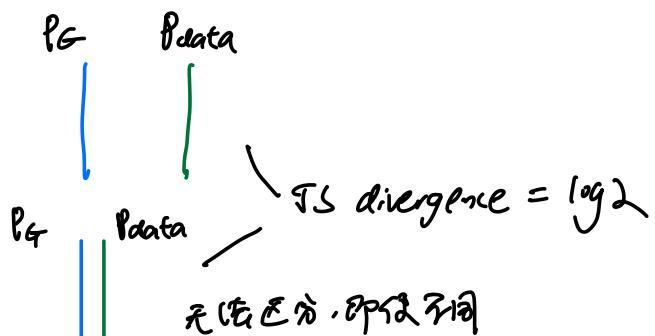
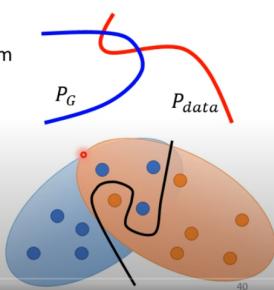
Both  $P_{data}$  and  $P_G$  are low-dim manifold in high-dim space.

The overlap can be ignored.

- 2. Sampling

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

If you do not have enough sampling



What's the problem of JS divergence?

JS divergence is always  $\log 2$  if two distribution do not overlap!

解决方法：

WGAN

WGAN

Evaluate Wasserstein distance between  $P_{data}$  and  $P_G$

spectral normalization

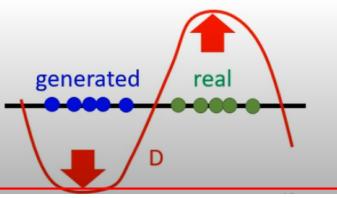
$$\max_{D \in 1\text{-Lipschitz}} \{E_{x \sim P_{data}}[D(x)] - E_{x \sim P_G}[D(x)]\}$$

D has to be smooth enough.

平滑限制了两个波峰相差不会很大

Without the constraint, the training of D will not converge.

Keeping the D smooth forces  $D(x)$  become  $\infty$  and  $-\infty$   
doesn't



## Evaluation of Generator (takeaways: '很多, 但因为没有多样性')

1. 用图片分类 / 识别系统挑出来的图像 .

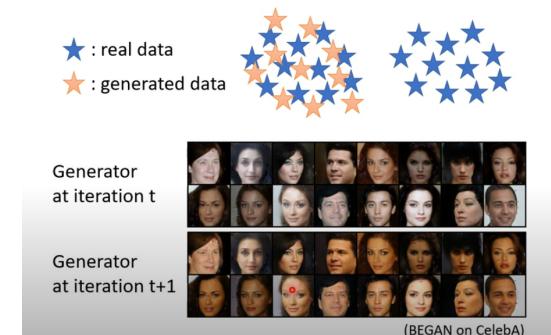
### 1. Drawback: Mode Collapse

#### Diversity - Mode Collapse

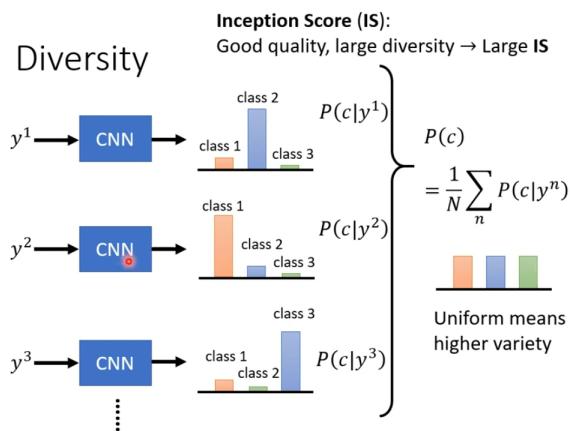


### mode Dropping

#### Diversity - Mode Dropping



2. 如何评价是高 diversity ?



## Conditional Generation (Conditional GAN)

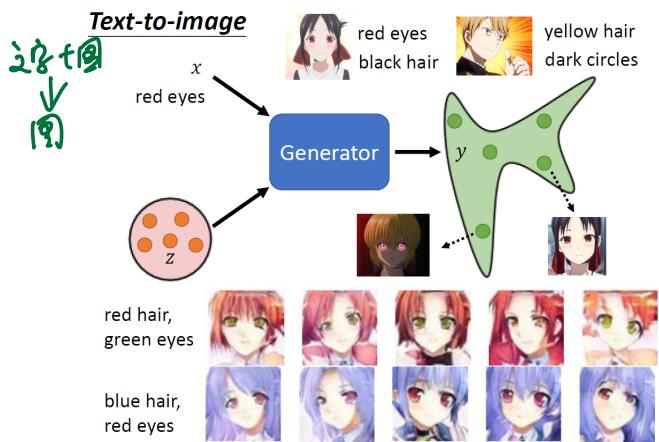
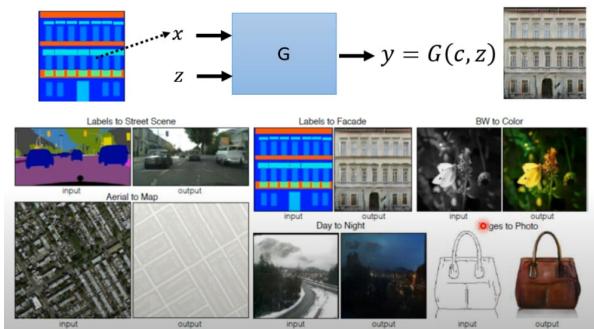


图 → 图

### Conditional GAN



## Learning from Unpaired Data

### Cycle GAN

