

Online Joint GlueX-EIC-PANDA Machine Learning Workshop

Generative Adversarial Networks

P. Jiang

Sep. 21-25, 2020

GSI Helmholtzzentrum für Schwerionenforschung GmbH

Institute of Modern Physics, Chinese Academy of Sciences

Many thanks to K. Goetzen, R. Kliemt, F. Nerling and K. Peters

Supported by China and Germany Postdoctoral Exchange Program

Outline



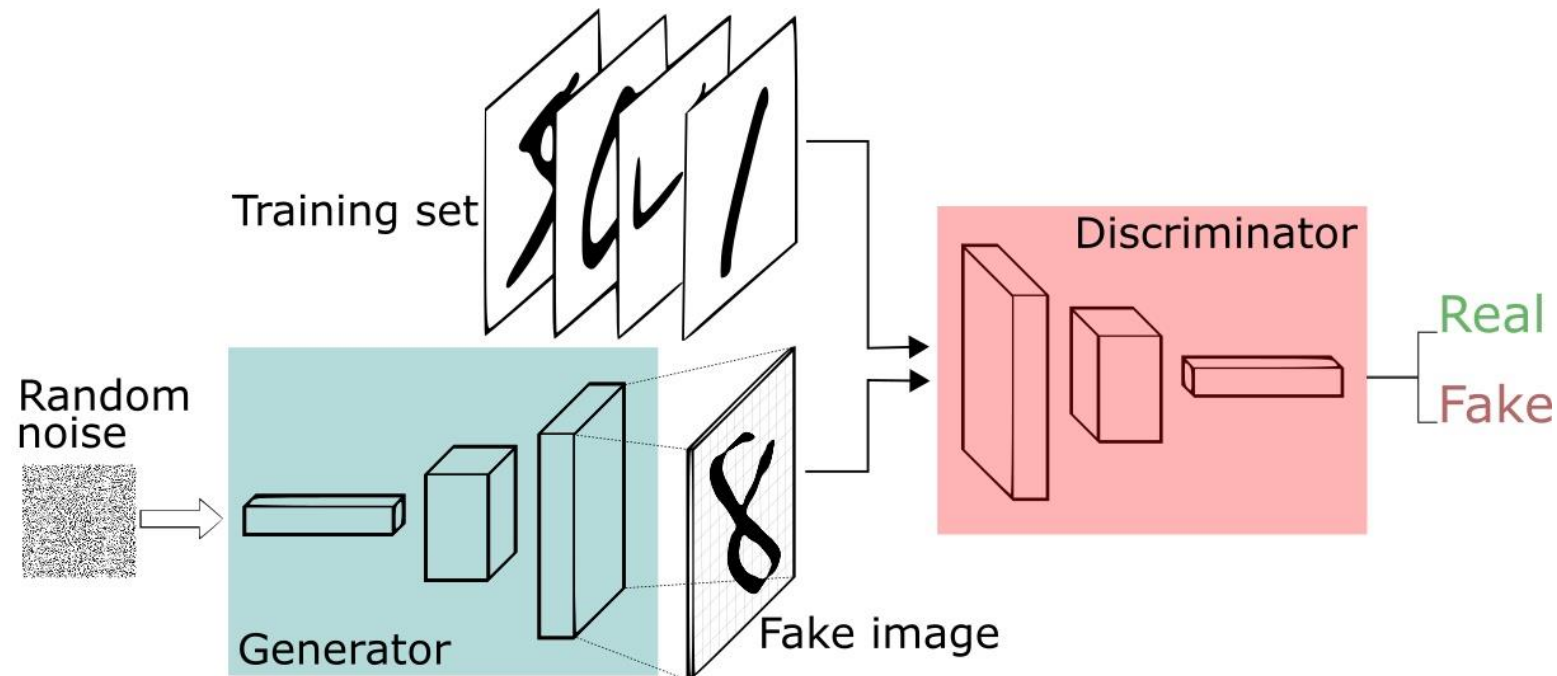
- Introduction
 - Prerequisites
 - Basic idea
 - Algorithm
 - Theory
 - WGAN
 - GANs
 - Evaluation
 - Tricks
 - Practice
- — What is Generative Adversarial Network (GAN)?
 - — What is the prior knowledge for GANs?
 - — How does a GAN work?
 - — How to train a GAN?
 - — What is the mathematics behind a GAN?
 - — Is there a more stable GAN?
 - — Are there other GANs?
 - — How good is a GAN?
 - — What are the tricks to make a GAN work?
 - — How to code the first line?

Introduction

— What is Generative Adversarial Network (GAN)?

Introduction

- The two neural networks, the generator and the discriminator, cooperate in an adversarial way to achieve the purpose of joint training.
- The generator generates fake samples while the discriminator discriminates the real samples and the fake samples.
- The generator is trained to “fool” the discriminator by generating the realistic looking fake samples.
- The discriminator is trained to discriminate the real samples and the fake samples.



Applications

- Image processing and computer vision
 - Super-resolution
 - Image synthesis and manipulation
 - Texture synthesis
 - Object detection
 - Video applications
 - Other image and vision applications
- Sequential data
 - Natural language processing
 - Music
 - Speech and Audio
- Other applications
 - Medical field
 - Data science



2014



2015



2016



2017

[arXiv:1802.07228](https://arxiv.org/abs/1802.07228)

[arXiv:2001.06937](https://arxiv.org/abs/2001.06937)

[arXiv 1703.10593](https://arxiv.org/abs/1703.10593)

Prerequisites

- What is the prior knowledge for GANs?

Binary classification

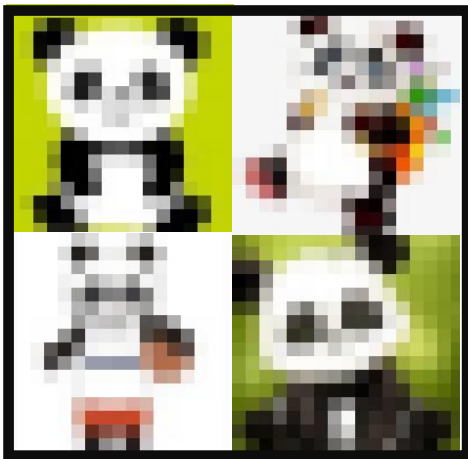
Real samples
from labeled
dataset

label=1

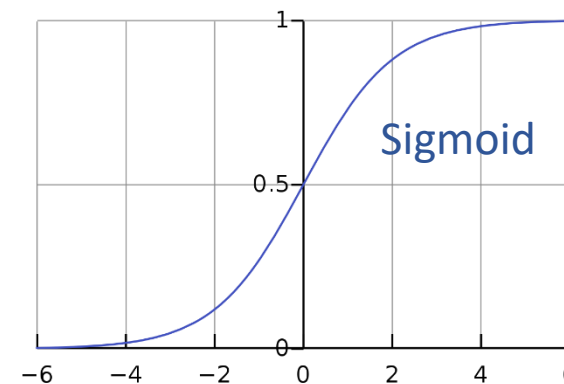
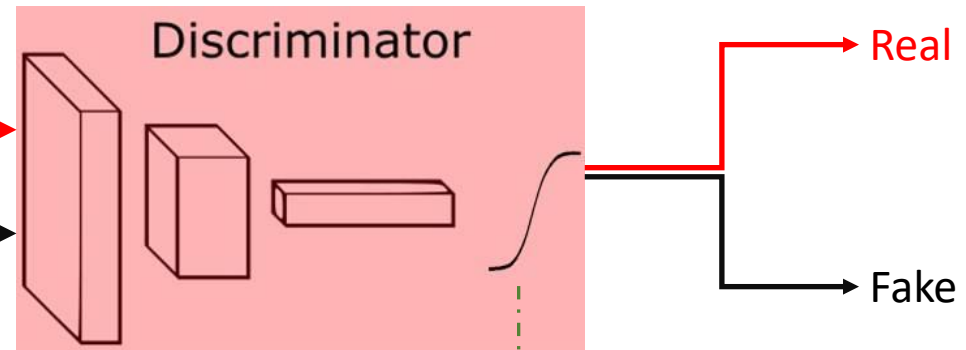


Fake samples
from labeled
dataset

label=0



$$\text{BCE loss} = - \sum (y \log(\hat{y}(\theta)) + (1-y) \log(1-\hat{y}(\theta)))$$



Binary cross entropy



- Given a sample, it can only be classified as 0 or 1. It follows the Bernoulli distribution.
- If the probability of the sample being classified as 1 is $\hat{y}(\theta)$, the probability of being classified as 0 is $1 - \hat{y}(\theta)$, where θ are neural network parameters.

- The probability mass function is $f(y; \hat{y}(\theta)) = \begin{cases} \hat{y}(\theta) & \text{if } y=1, \\ 1 - \hat{y}(\theta) & \text{if } y=0. \end{cases} = \hat{y}(\theta)^y (1 - \hat{y}(\theta))^{(1-y)} \quad y \in \{0,1\}$

- When we have multiple samples, we can construct the likelihood function $L(\theta) = \prod \hat{y}(\theta)^y (1 - \hat{y}(\theta))^{(1-y)}$

- By maximizing $L(\theta)$, update $\theta \leftarrow \underset{\theta}{\operatorname{argmax}} L(\theta)$

- Maximizing likelihood function is equivalent to minimizing binary cross entropy

$$\underset{\theta}{\operatorname{argmax}} L(\theta) = \underset{\theta}{\operatorname{argmax}} \prod \hat{y}(\theta)^y (1 - \hat{y}(\theta))^{(1-y)} = \underset{\theta}{\operatorname{argmax}} \log(\prod \hat{y}(\theta)^y (1 - \hat{y}(\theta))^{(1-y)})$$

$$= \underset{\theta}{\operatorname{argmax}} \sum (y \log(\hat{y}(\theta)) + (1-y) \log(1 - \hat{y}(\theta))) = \underset{\theta}{\operatorname{argmin}} [-\sum (y \log(\hat{y}(\theta)) + (1-y) \log(1 - \hat{y}(\theta)))]$$

$$= \underset{\theta}{\operatorname{argmin}} \operatorname{BCE}(y; \hat{y}(\theta))$$

- **BCE loss** = $-\sum (y \log(\hat{y}(\theta)) + (1-y) \log(1 - \hat{y}(\theta)))$

Basic idea

— How does a GAN work?

From binary classification to GAN

Real samples
from labeled
dataset

label=1



Fake samples
from labeled
dataset
from generator

label=0

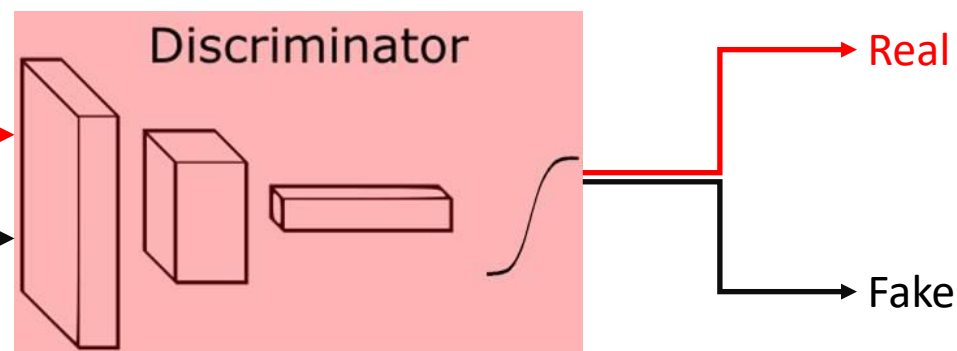


$$\text{BCE loss} = - \sum (y \log(\hat{y}(\theta)) + (1-y) \log(1-\hat{y}(\theta)))$$

$$\text{BCE loss} = - \sum (y_{\text{real}} \log(\hat{y}_{\text{real}}) + (1-y_{\text{fake}}) \log(1-\hat{y}_{\text{fake}}))$$

$$\text{BCE loss} = - \sum (\log(\hat{y}_{\text{real}}) + \log(1-\hat{y}_{\text{fake}}))$$

$$\text{BCE loss} = - \sum (\log(D(x_{\text{real}})) + \log(1-D(x_{\text{fake}})))$$



Encoder-Decoder



[mouth = 0.49
ear = -0.67
⋮
eyes = 0.44
nose = 1.80]



Encoder

My daddy has a big mouth...

About 0.49

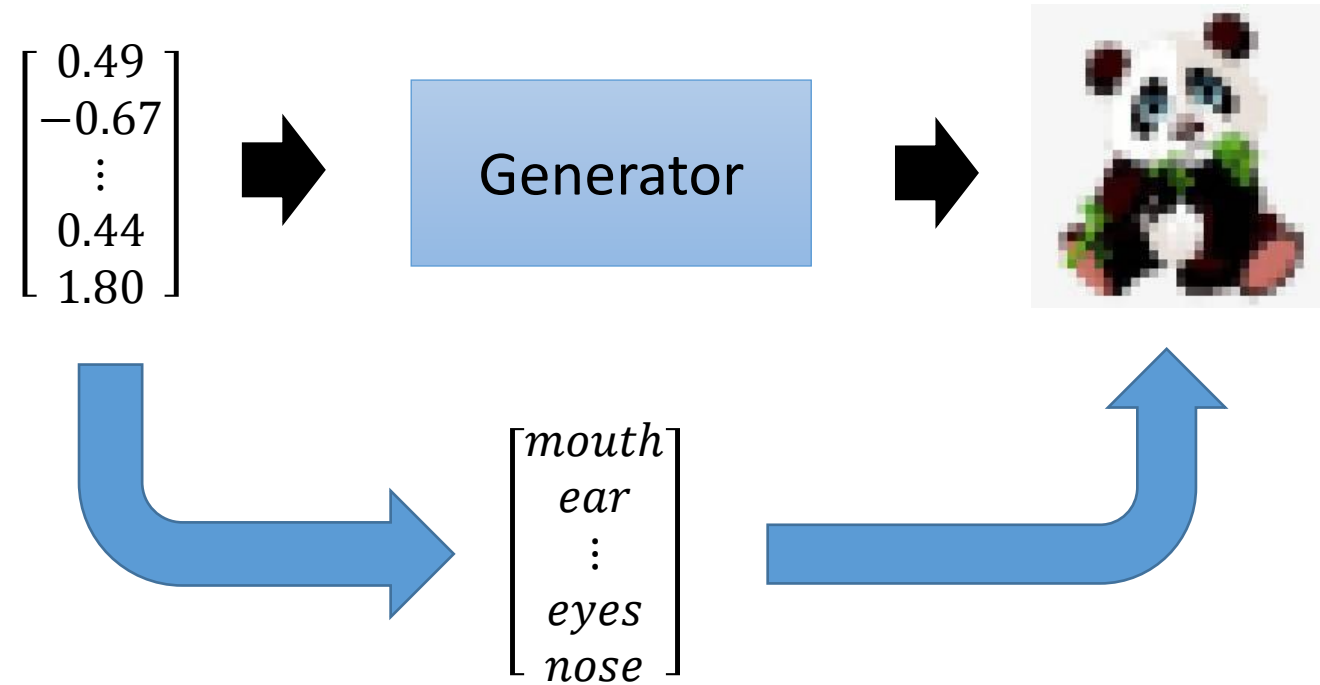
How big?

Wow! Huge...



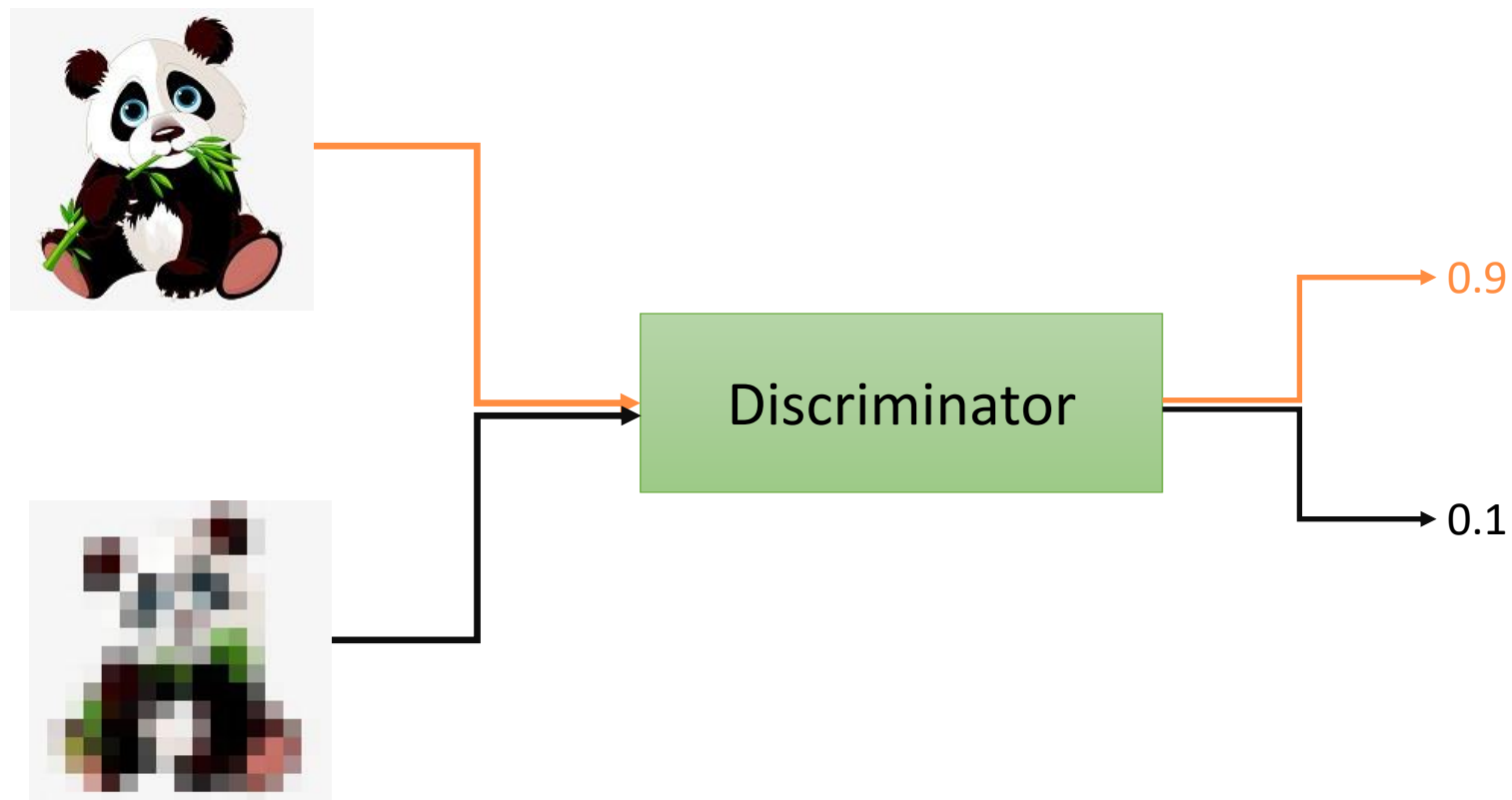
Decoder

Generator



Each dimension of input vector represents some characteristics.

Discriminator



Contest



Real money



Discriminator



Generator

Our slogan is to make the most realistic looking fake money (counterfeit money)...

Fake money



Currency detector



Algorithm

— How to train a GAN?

Step 1: Train D (Fix G)



real dataset

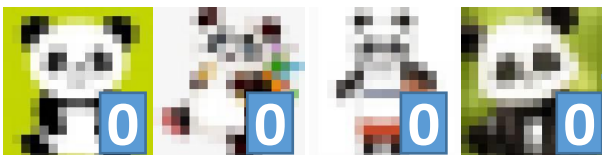


sample

real samples



fake samples



train



Gradient descent

~~0.5~~

0.2

G

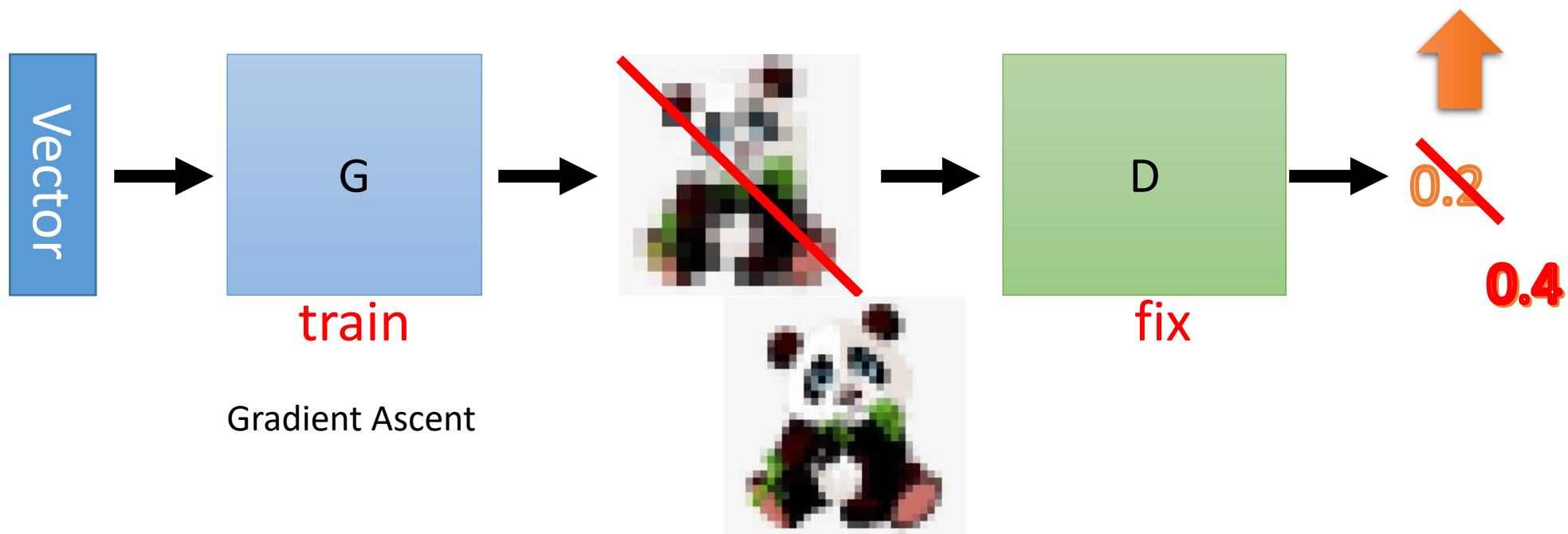
fix

randomly sampled



Step 2: Train G (Fix D)

Generator learns to “fool” the Discriminator.



Algorithm

- Initialize θ_d for D and θ_g for G

- In each training iteration:

- Sample m real samples $\{x^1, x^2, \dots, x^m\}$ from dataset
- Sample m noise samples $\{z^1, z^2, \dots, z^m\}$ from a distribution
- Generate m fake samples $\{\tilde{x}^1, \tilde{x}^2, \dots, \tilde{x}^m\}$, $\tilde{x}^i = G(z^i)$
- Update discriminator parameters θ_d to **minimize**

- $\bar{V}_d = -\left[\frac{1}{m} \sum_{i=1}^m \log D(x^i) + \frac{1}{m} \sum_{i=1}^m \log(1 - D(\tilde{x}^i))\right]$
- $\theta_d \leftarrow \theta_d - \eta \nabla_{\theta_d} \bar{V}_d(\theta_d)$

Train D

- Sample m noise samples $\{z^1, z^2, \dots, z^m\}$ from a distribution
- Update generator parameters θ_g to **maximize**

- ~~$\bar{V}_g = -\left[\frac{1}{m} \sum_{i=1}^m \log D(x^i) + \frac{1}{m} \sum_{i=1}^m \log(1 - D(\tilde{x}^i))\right]$~~ $\bar{V}_g = -\left[\frac{1}{m} \sum_{i=1}^m \log(1 - D(\tilde{x}^i))\right]$
- $\bar{V}_g = -\frac{1}{m} \sum_{i=1}^m \log(1 - D(G(z^i)))$
- $\theta_g \leftarrow \theta_g + \eta \nabla_{\theta_g} \bar{V}_g(\theta_g)$

Train G

Algorithm of GAN

- Initialize θ_d for D and θ_g for G
- In each training iteration:

- Sample m real samples $\{x^1, x^2, \dots, x^m\}$ from dataset
- Sample m noise samples $\{z^1, z^2, \dots, z^m\}$ from a distribution
- Generate m fake samples $\{\tilde{x}^1, \tilde{x}^2, \dots, \tilde{x}^m\}$, $\tilde{x}^i = G(z^i)$
- Update discriminator parameters θ_d to **minimize**
 - $\bar{V}_d = -\left(\frac{1}{m} \sum_{i=1}^m \log D(x^i) + \frac{1}{m} \sum_{i=1}^m \log(1 - D(\tilde{x}^i))\right)$
 - $\theta_d \leftarrow \theta_d - \eta \nabla_{\theta_d} \bar{V}_d(\theta_d)$

Train D

- Sample m noise samples $\{z^1, z^2, \dots, z^m\}$ from a distribution
- Update generator parameters θ_g to **minimize**
 - $V_g = \frac{1}{m} \sum_{i=1}^m \log(1 - D(G(z^i)))$
 - $\theta_g \leftarrow \theta_g - \eta \nabla_{\theta_g} V_g(\theta_g)$

Train G

$$\bar{V}_g = -\frac{1}{m} \sum_{i=1}^m \log(1 - D(G(z^i)))$$

Algorithm of GAN

- Initialize θ_d for D and θ_g for G
- In each training iteration:

- Sample m real samples $\{x^1, x^2, \dots, x^m\}$ from dataset
- Sample m noise samples $\{z^1, z^2, \dots, z^m\}$ from a distribution
- Generate m fake samples $\{\tilde{x}^1, \tilde{x}^2, \dots, \tilde{x}^m\}$, $\tilde{x}^i = G(z^i)$
- Update discriminator parameters θ_d to minimize
 - $\bar{V}_d = \text{BCELoss}$
 - $\theta_d \leftarrow \theta_d - \eta \nabla_{\theta_d} \bar{V}_d(\theta_d)$

Train D

- Sample m noise samples $\{z^1, z^2, \dots, z^m\}$ from a distribution
- Update generator parameters θ_g to minimize
 - $V_g = -\text{BCELoss}$
 - $\theta_g \leftarrow \theta_g - \eta \nabla_{\theta_g} V_g(\theta_g)$

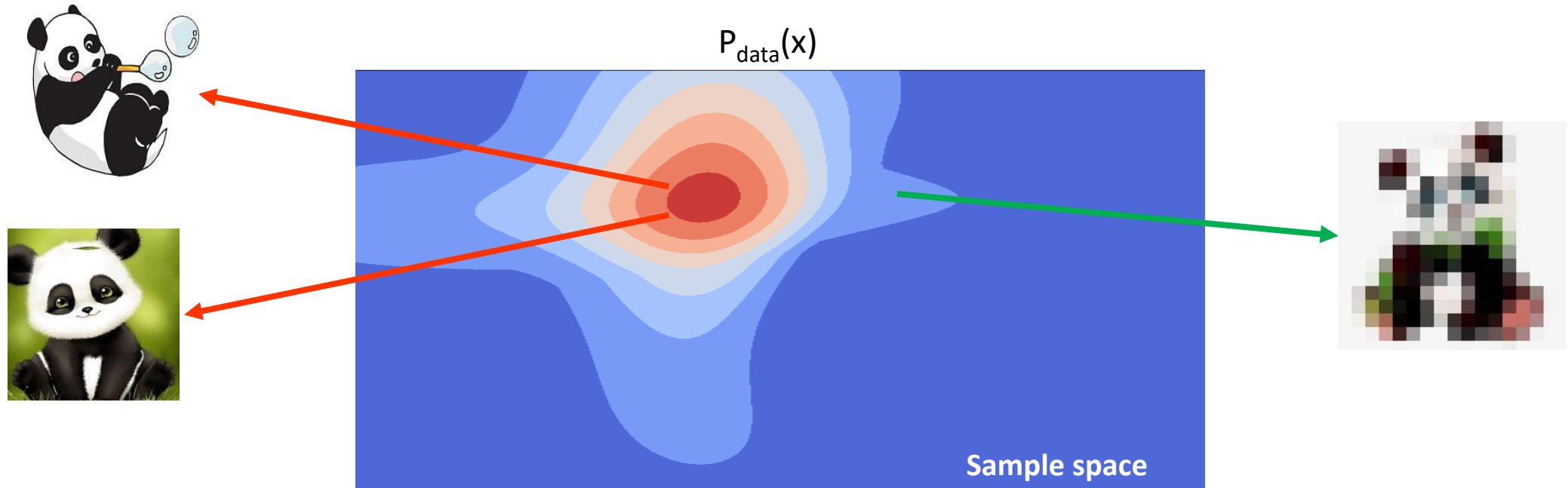
Train G

Theory

— What is the mathematics behind a GAN?

Probability distribution - $P_{\text{data}}(x)$

- Probability distribution of given data
 - $P_{\text{data}}(x)$
 - x : samples



Find $P_{\text{data}}(x)$ - Maximum likelihood estimation

- **Given** a data probability distribution $P_{\text{data}}(x)$
- Assume a distribution $P_G(x;\theta)$ parameterized by θ
- Aim: Find θ such that $P_G(x;\theta)$ close to $P_{\text{data}}(x)$

- Assume $P_G(x;\theta)$ is a model parameterized by θ
(e.g. $P_G(x;\theta)$ is a Gaussian Mixture Model, θ are means and variances of the Gaussians)

- Sample $\{x^1, x^2, \dots, x^m\}$ from $P_{\text{data}}(x)$
- Compute $P_G(x^i; \theta)$
- Likelihood function of the sampled samples

$$L = \prod_{i=1}^m P_G(x^i; \theta)$$

- Find θ^* by maximizing the likelihood function

Maximizing likelihood function is equivalent to minimizing Kullback-Leibler Divergence



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

$$\theta^* = \operatorname{argmax}_{\theta} \prod_{i=1}^m P_G(x^i; \theta)$$

$$= \operatorname{argmax}_{\theta} \log \prod_{i=1}^m P_G(x^i; \theta)$$

$$= \operatorname{argmax}_{\theta} \sum_{i=1}^m \log P_G(x^i; \theta)$$

$$\approx \operatorname{argmax}_{\theta} E_{x \sim P_{\text{data}}} [\log P_G(x; \theta)]$$

$$= \operatorname{argmax}_{\theta} \int_{\mathbf{x}} P_{\text{data}}(\mathbf{x}) \log P_G(\mathbf{x}; \theta) d\mathbf{x}$$

$$= \operatorname{argmin}_{\theta} \left[- \int_{\mathbf{x}} P_{\text{data}}(\mathbf{x}) \log P_G(\mathbf{x}; \theta) d\mathbf{x} \right]$$

$$= \operatorname{argmin}_{\theta} \left[\underbrace{- \int_{\mathbf{x}} P_{\text{data}}(\mathbf{x}) \log P_G(\mathbf{x}; \theta) d\mathbf{x}}_{\text{KL}(P_{\text{data}} || P_G)} + \int_{\mathbf{x}} P_{\text{data}}(\mathbf{x}) \log P_{\text{data}}(\mathbf{x}) d\mathbf{x} \right]$$

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

Maximizing likelihood function is equivalent to minimizing KLD

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

$$\theta^* = \operatorname{argmax}_{\theta} \prod_{i=1}^m P_G(x^i; \theta)$$

$$= \operatorname{argmax}_{\theta} \log \prod_{i=1}^m P_G(x^i; \theta)$$

$$= \operatorname{argmax}_{\theta} \sum_{i=1}^m \log P_G(x^i; \theta)$$

$$\approx \operatorname{argmax}_{\theta} E_{x \sim P_{\text{data}}} [\log P_G(x; \theta)]$$

$$= \operatorname{argmax}_{\theta} \int_{\mathbf{x}} P_{\text{data}}(\mathbf{x}) \log P_G(\mathbf{x}; \theta) d\mathbf{x}$$

$$= \operatorname{argmin}_{\theta} \left[- \int_{\mathbf{x}} P_{\text{data}}(\mathbf{x}) \log P_G(\mathbf{x}; \theta) d\mathbf{x} \right]$$

$$= \operatorname{argmin}_{\theta} \left[\underbrace{- \int_{\mathbf{x}} P_{\text{data}}(\mathbf{x}) \log P_G(\mathbf{x}; \theta) d\mathbf{x}}_{\text{Cross entropy}} + \underbrace{\int_{\mathbf{x}} P_{\text{data}}(\mathbf{x}) \log P_{\text{data}}(\mathbf{x}) d\mathbf{x}}_{\text{Entropy}} \right]$$

$$= \operatorname{argmin}_{\theta} \text{KL}(P_{\text{data}} \parallel P_G)$$

Cross entropy
 $H(P_{\text{data}}(x), P_G(x; \theta))$

Entropy
 $-H(P_{\text{data}}(x))$

- Point to Point

- p – norm distance = $(\sum_{i=1}^n |x_i - y_i|^p)^{1/p}$

- Distribution to Distribution

- Kullback–Leibler divergence

- $D_{\text{KL}}(P||Q) = \sum_{x \in X} P(x) \log\left(\frac{P(x)}{Q(x)}\right)$

- $D_{\text{KL}}(P||Q) = \int_{-\infty}^{\infty} P(x) \log\left(\frac{P(x)}{Q(x)}\right) dx$

- $D_{\text{KL}}(P||Q) = \sum_{x \in X} P(x) \log\left(\frac{P(x)}{Q(x)}\right)$
 $= -\sum_{x \in X} P(x) \log(Q(x)) + \sum_{x \in X} P(x) \log(P(x))$
 $= H(P, Q) - H(P)$

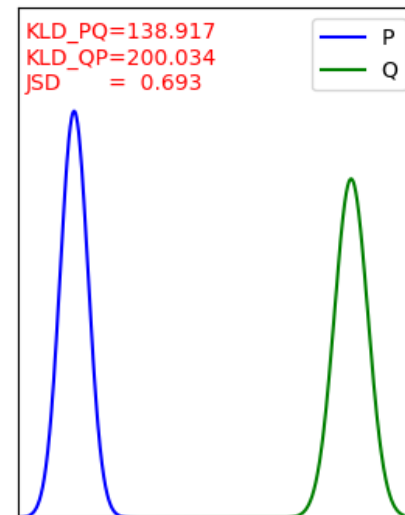
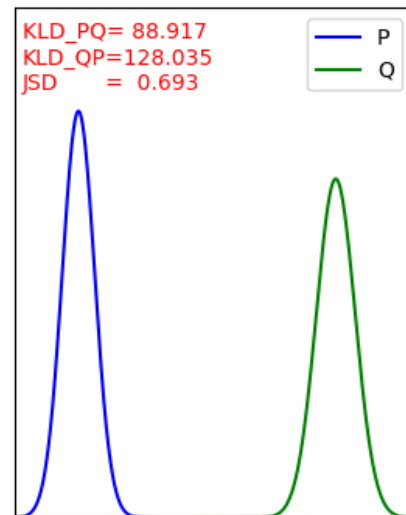
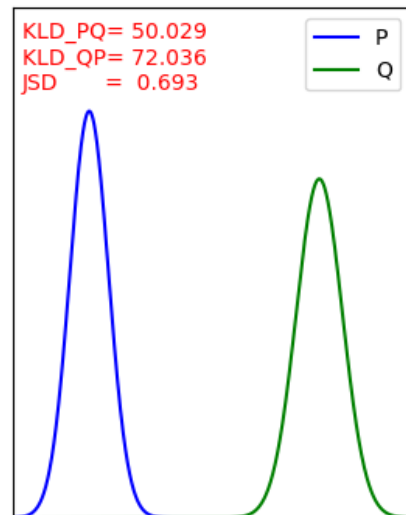
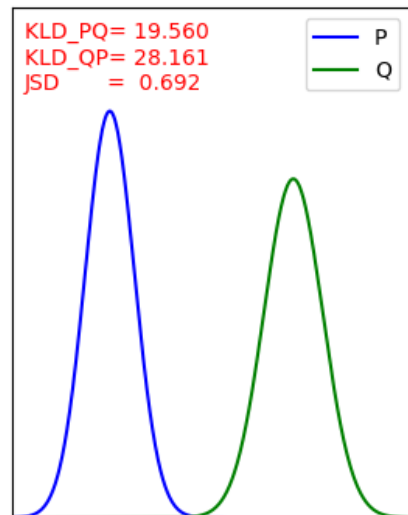
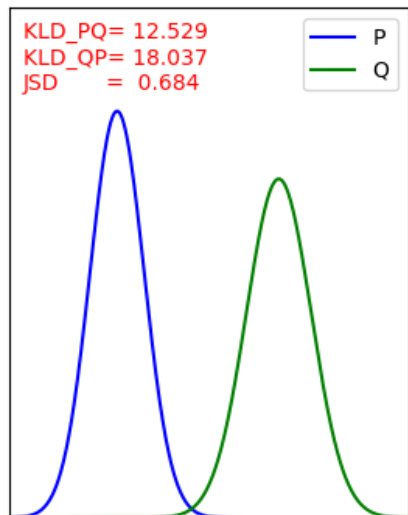
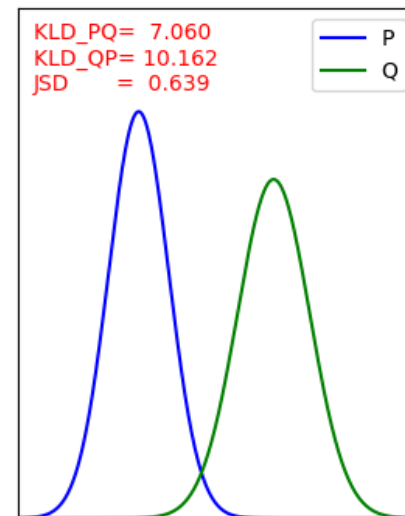
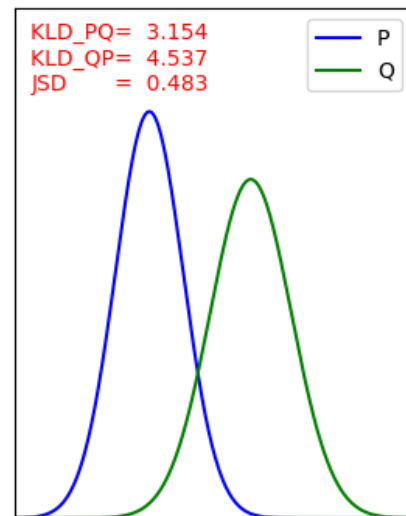
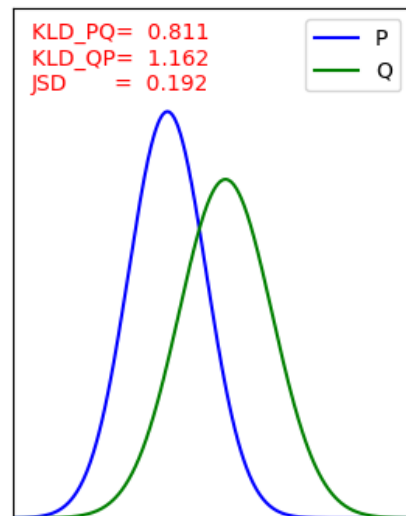
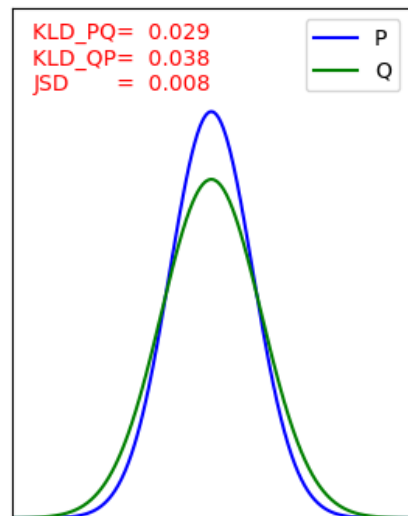
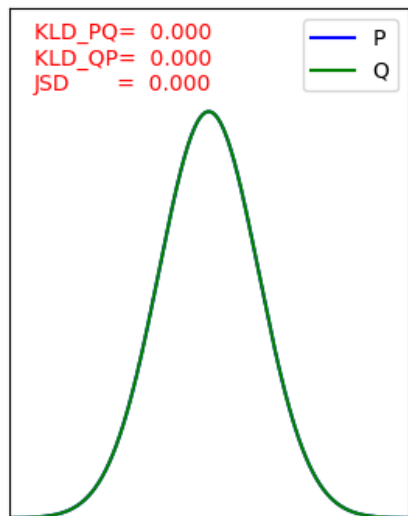
- $D_{\text{KL}}(P||Q) \neq D_{\text{KL}}(Q||P)$

- Jensen–Shannon divergence

- $\text{JSD}(P||Q) = \frac{1}{2} D_{\text{KL}}(P||M) + \frac{1}{2} D_{\text{KL}}(Q||M)$, where $M = \frac{1}{2}(P + Q)$

- $\text{JSD}(P||Q) = \text{JSD}(Q||P)$

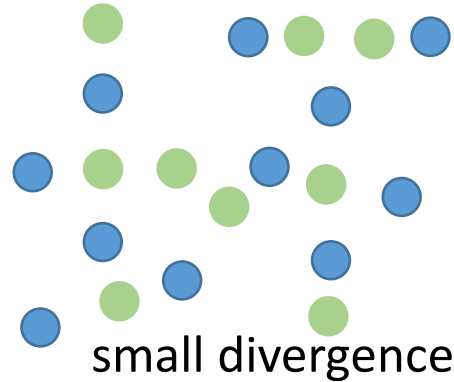
KLD & JSD



Divergence and discriminator

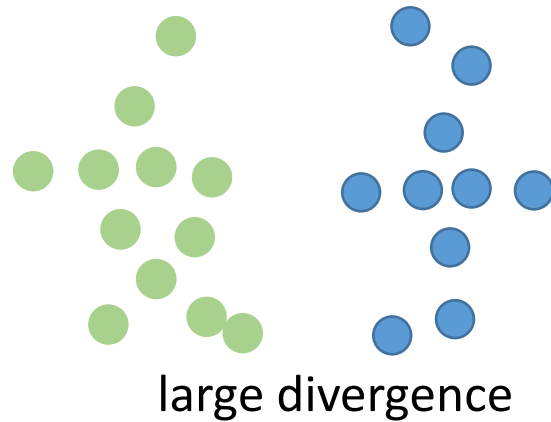
● : data sampled from P_{data}

● : data sampled from P_G



Discriminator

hard to discriminate

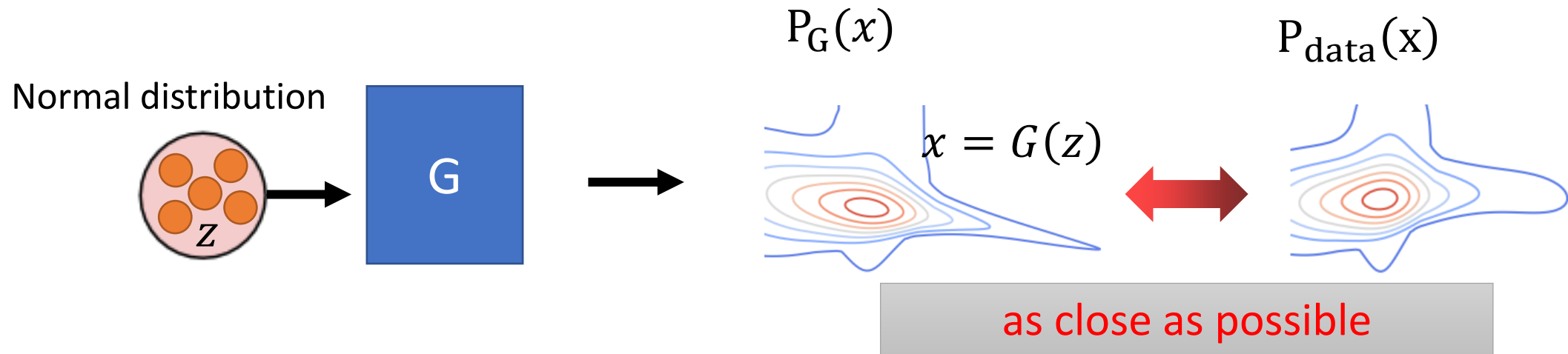


Discriminator

easy to discriminate

Generator

- A generator G is a network which defines a probability distribution P_G



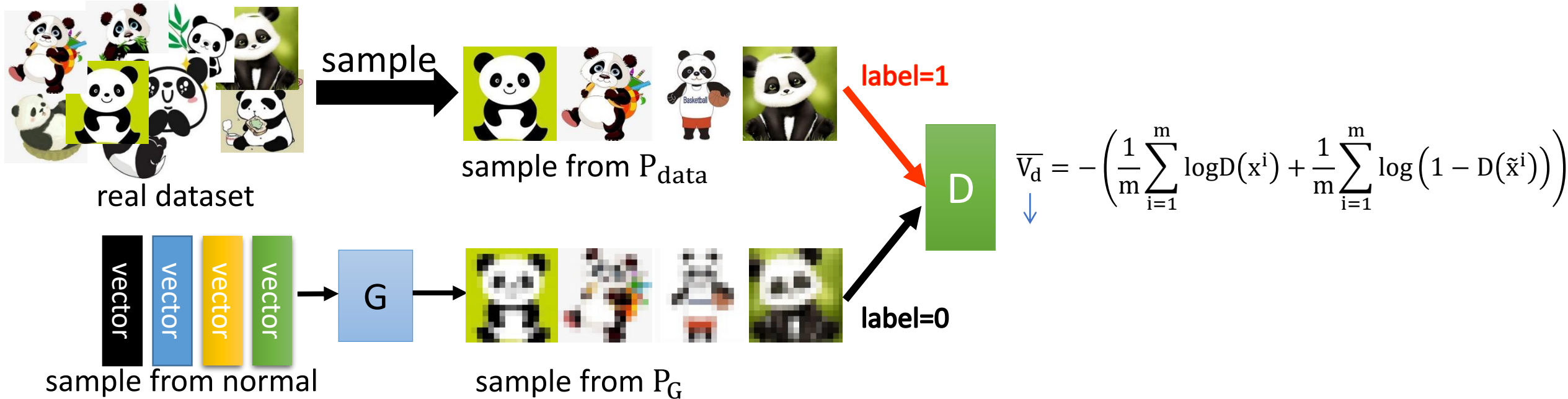
$$\text{Train } G: G^* = \arg \min_G \text{Div}(P_G, P_{\text{data}})$$

Divergence between distributions P_G and P_{data}

Discriminator



- Although we do not know the distributions of P_G and P_{data} , we can sample from them.



- Objective Function for D

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

↑ ↑ G is fixed ↓
0

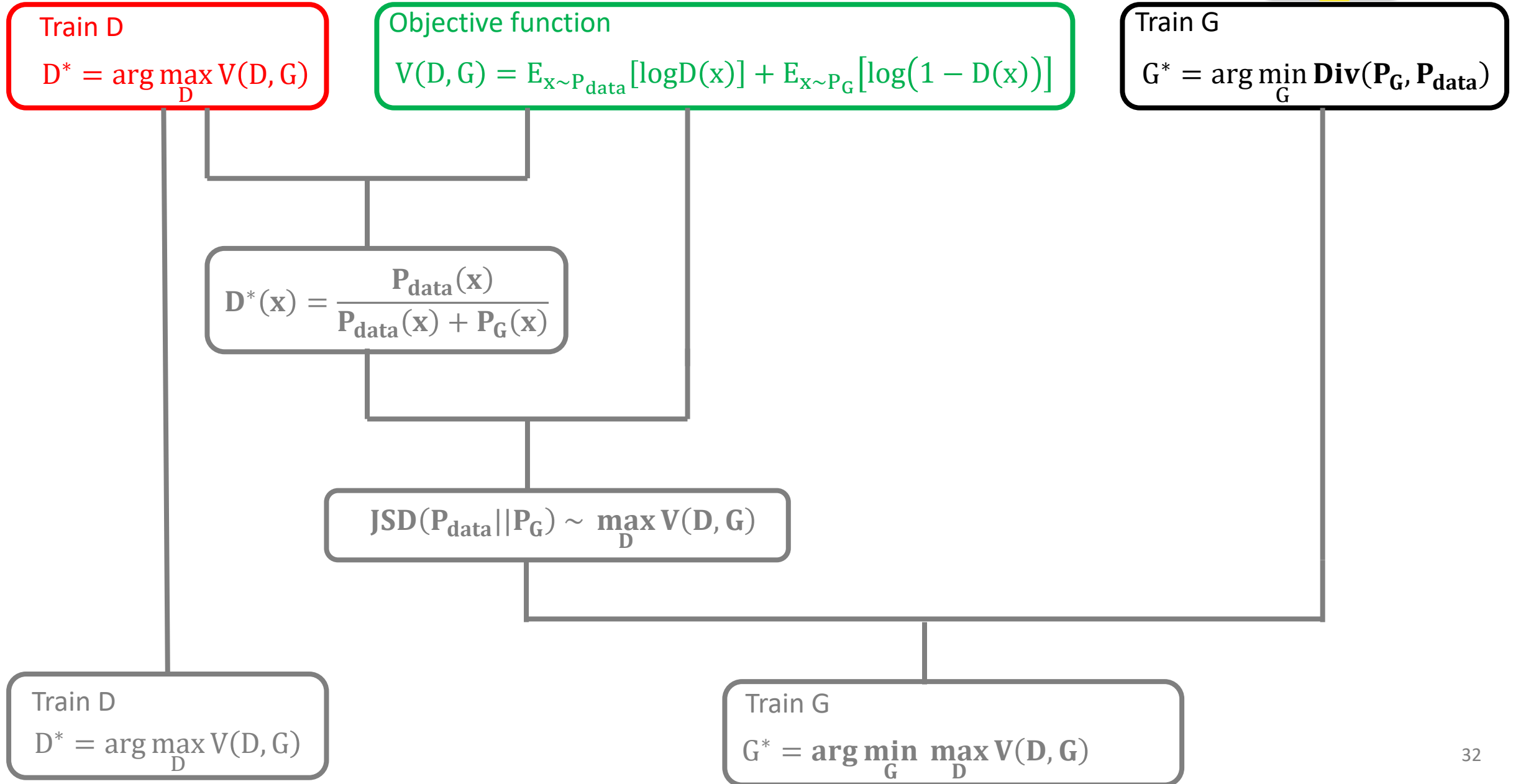
- Train D: $D^* = \arg \max_D V(D, G)$

G^* & D^*



- $G^* = \arg \min_G \text{Div}(P_G, P_{\text{data}})$
- $D^* = \arg \max_D V(D, G)$
- $V(D, G) = E_{x \sim P_{\text{data}}}[\log D(x)] + E_{x \sim P_G}[\log(1 - D(x))]$

Family tree



D*



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

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

- Given G, assuming that D(x) can be any function, the optimal D* maximize

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

- $f(D) = a \log(D) + b \log(1 - D)$

- $\frac{df(D)}{dD} = 0 \implies D^* = \frac{a}{a+b}$

- $D^*(x) = \frac{P_{\text{data}}(x)}{P_{\text{data}}(x) + P_G(x)} \quad D^*(x) \in (0,1)$

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



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

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

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

$$= \mathbf{E}_{\mathbf{x} \sim P_{\text{data}}} \left[\log \frac{P_{\text{data}}(\mathbf{x})}{P_{\text{data}}(\mathbf{x}) + P_G(\mathbf{x})} \right] + \mathbf{E}_{\mathbf{x} \sim P_G} \left[\log \frac{P_G(\mathbf{x})}{P_{\text{data}}(\mathbf{x}) + P_G(\mathbf{x})} \right]$$

$$= \int_{\mathbf{x}} P_{\text{data}}(\mathbf{x}) \log \frac{P_{\text{data}}(\mathbf{x})}{P_{\text{data}}(\mathbf{x}) + P_G(\mathbf{x})} d\mathbf{x} + \int_{\mathbf{x}} P_G(\mathbf{x}) \log \frac{P_G(\mathbf{x})}{P_{\text{data}}(\mathbf{x}) + P_G(\mathbf{x})} d\mathbf{x}$$

$$= \int_{\mathbf{x}} P_{\text{data}}(\mathbf{x}) \log \frac{\frac{1}{2} P_{\text{data}}(\mathbf{x})}{(P_{\text{data}}(\mathbf{x}) + P_G(\mathbf{x})) / 2} d\mathbf{x} + \int_{\mathbf{x}} P_G(\mathbf{x}) \log \frac{\frac{1}{2} P_G(\mathbf{x})}{(P_{\text{data}}(\mathbf{x}) + P_G(\mathbf{x})) / 2} d\mathbf{x}$$

$$= -2 \log 2 + \int_{\mathbf{x}} P_{\text{data}}(\mathbf{x}) \log \frac{P_{\text{data}}(\mathbf{x})}{(P_{\text{data}}(\mathbf{x}) + P_G(\mathbf{x})) / 2} d\mathbf{x} + \int_{\mathbf{x}} P_G(\mathbf{x}) \log \frac{P_G(\mathbf{x})}{(P_{\text{data}}(\mathbf{x}) + P_G(\mathbf{x})) / 2} d\mathbf{x}$$

$$= -2 \log 2 + \text{KL} \left(P_{\text{data}} \parallel \frac{P_{\text{data}} + P_G}{2} \right) + \text{KL} \left(P_G \parallel \frac{P_{\text{data}} + P_G}{2} \right)$$

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

G^* & D^* - Div $\rightarrow \max_D V(D, G)$

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

$$G^* = \arg \min_G \text{Div}(P_G, P_{\text{data}})$$

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

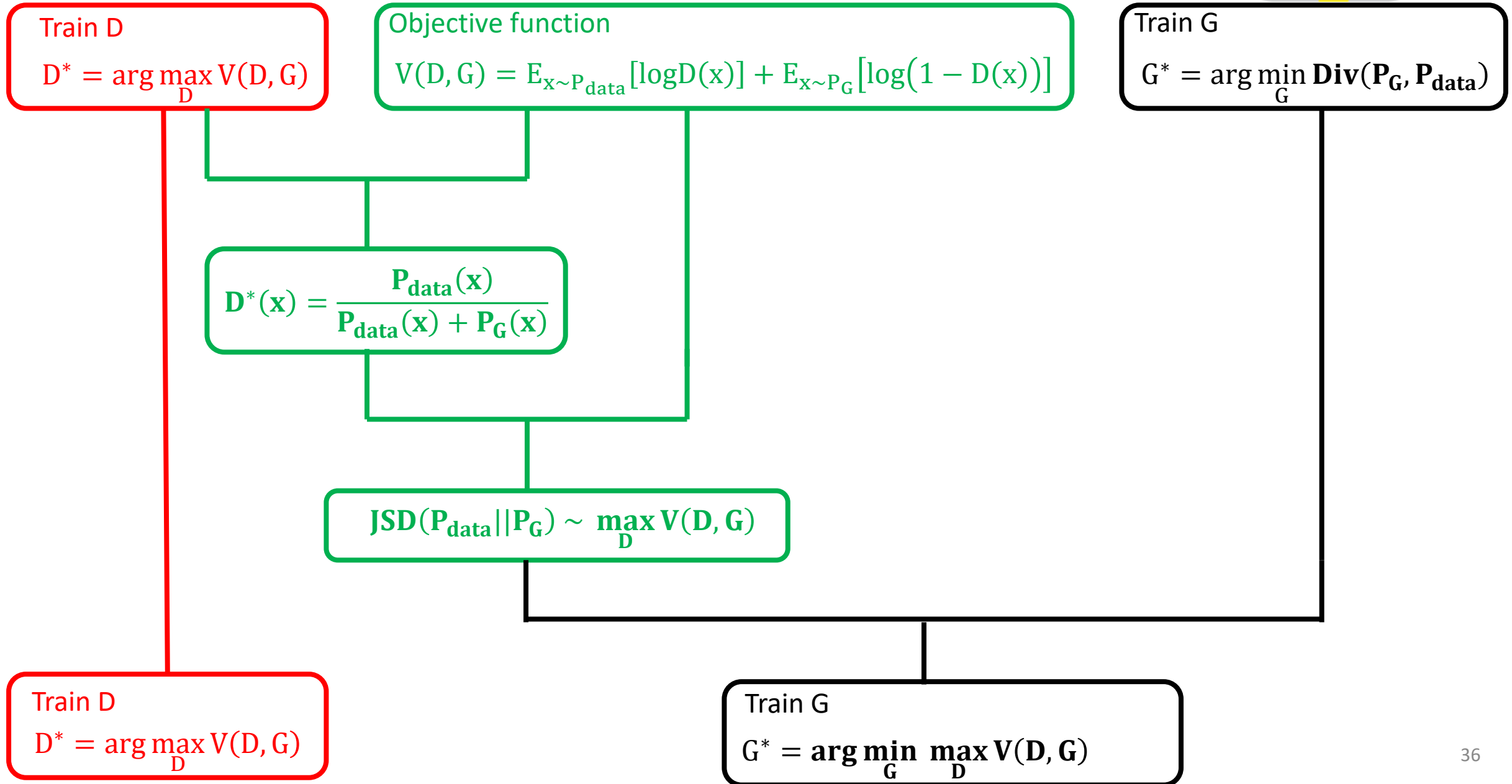


$$\text{JSD}(P_{\text{data}} || P_G) \sim \max_D V(D, G)$$

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

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

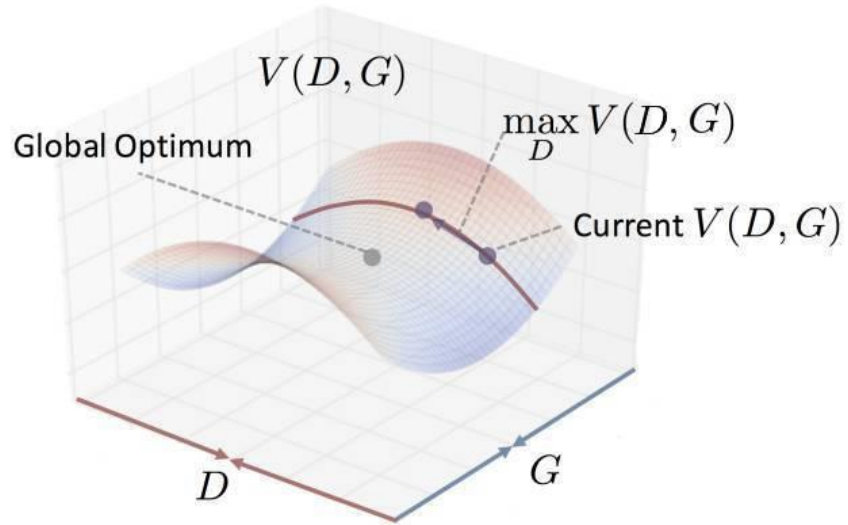
Family tree



D* & G* - Algorithm

Step 1. Train D

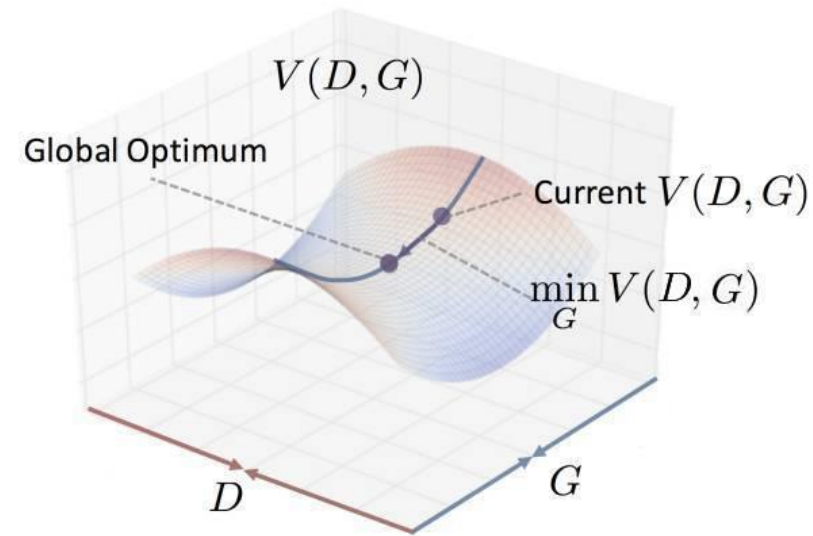
- $D^* = \arg \max_D V(D, G)$



- Fix G , and update D

Step 2. Train G

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

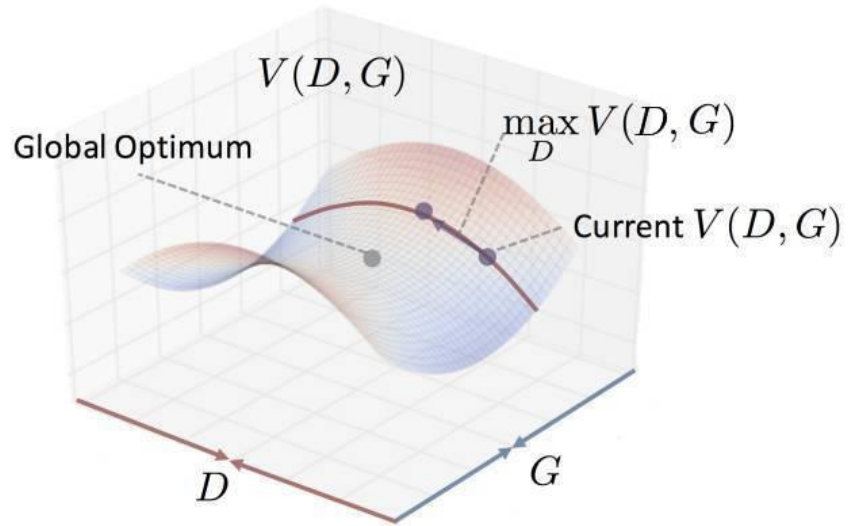


- Fix D , and update G

D* & G* - Algorithm

Step 1. Train D

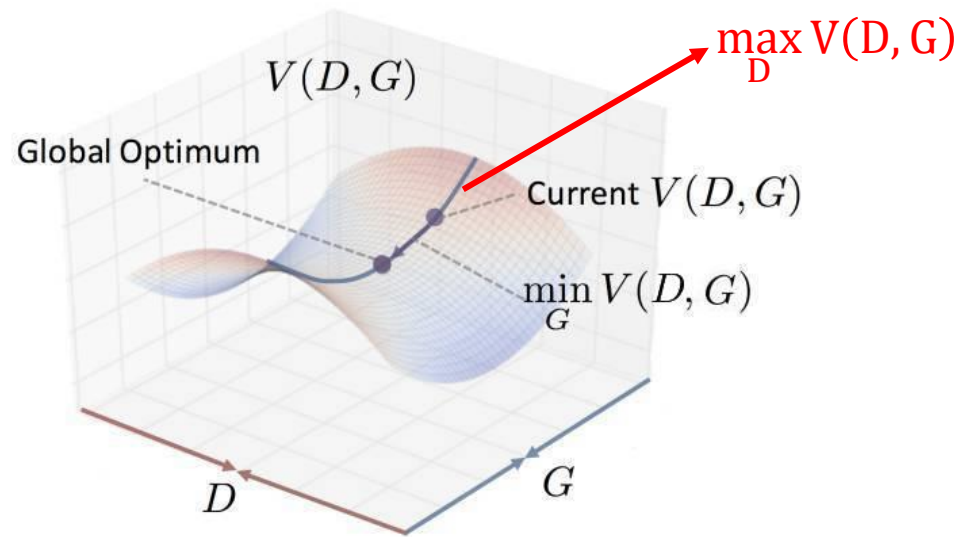
- $D^* = \arg \max_D V(D, G)$



- Fix G , and update D

Step 2. Train G

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

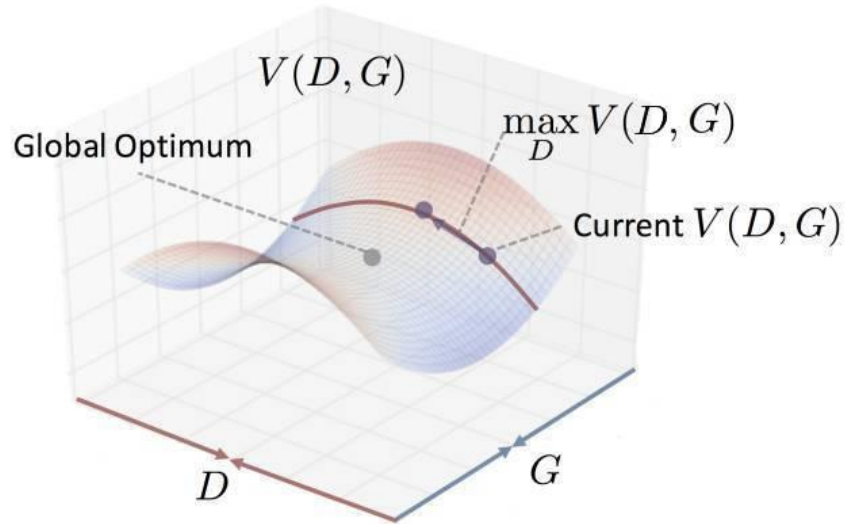


- Fix D , and update G

D* & G* - Algorithm

Step 1. Train D

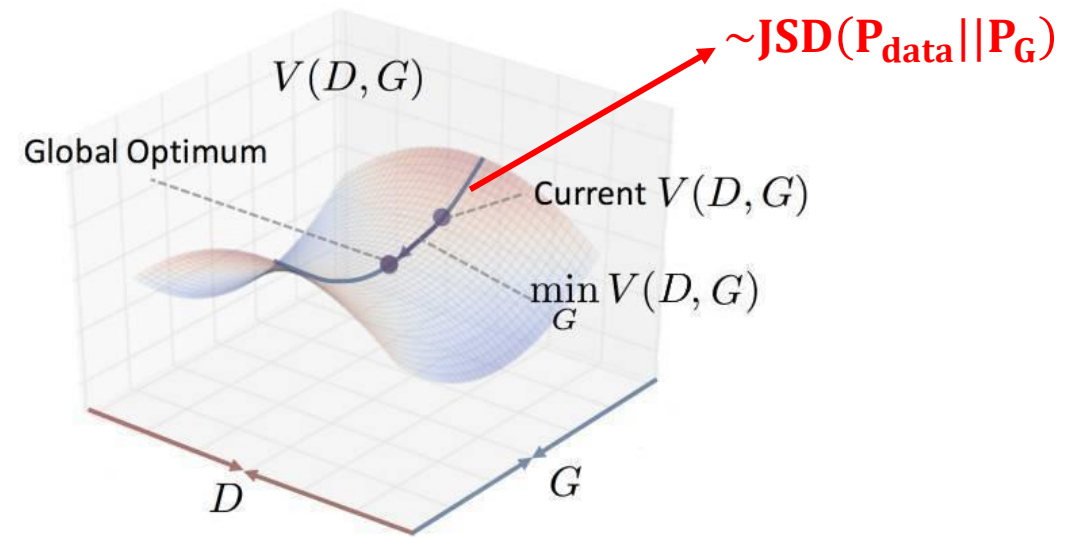
- $D^* = \arg \max_D V(D, G)$



- Fix G , and update D

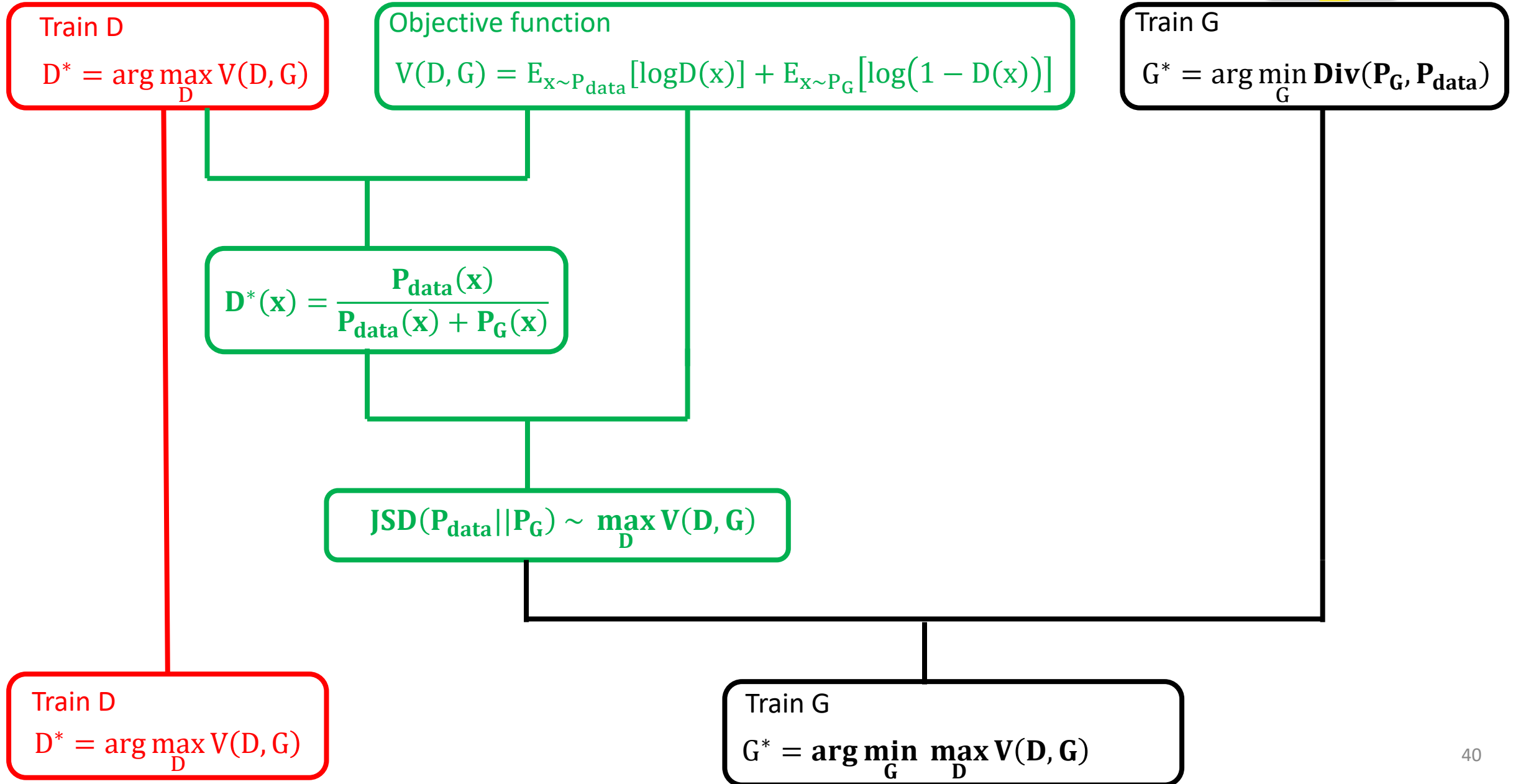
Step 2. Train G

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

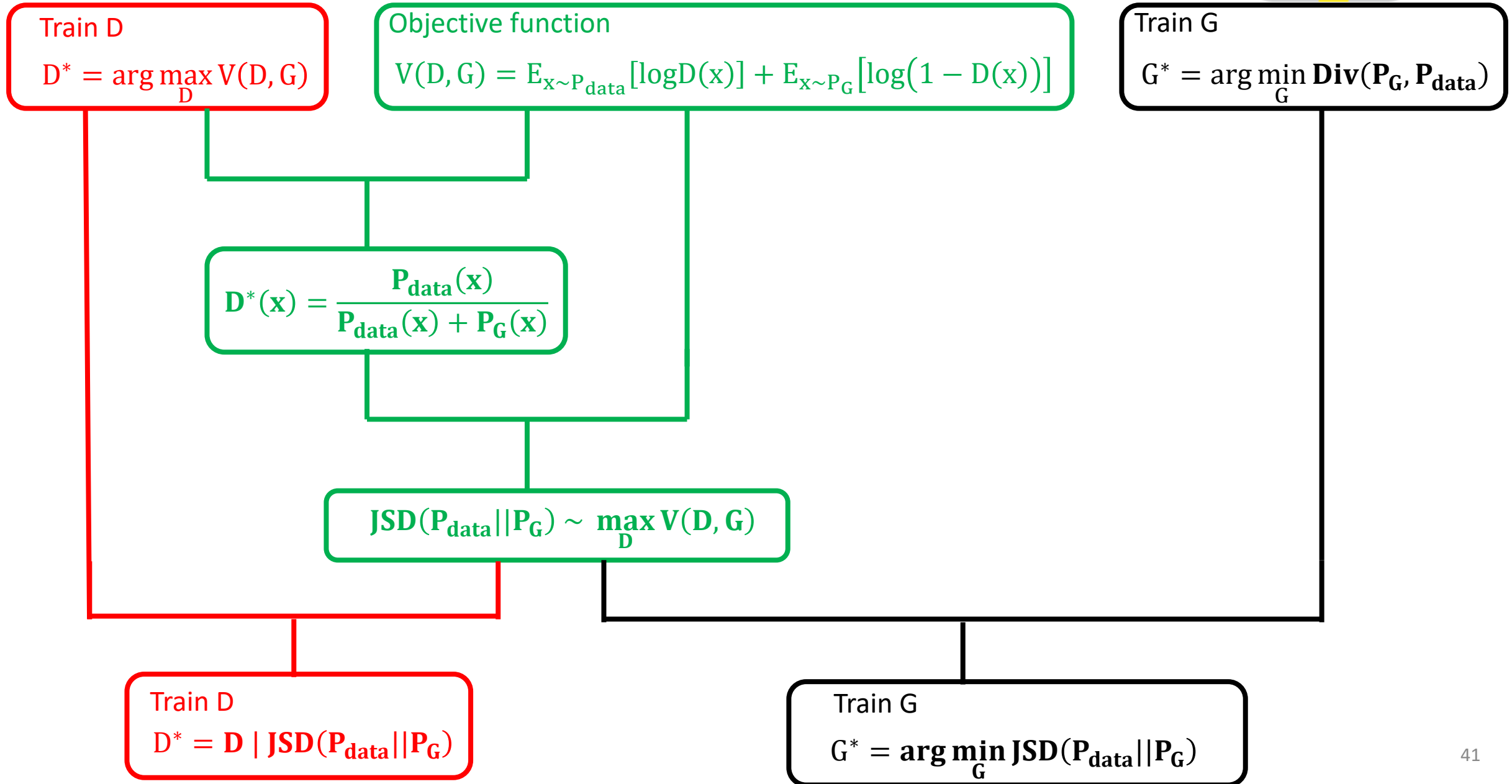


- Fix D , and update G

Family tree



Family tree



D* & G* - Algorithm

Step 1. Train D

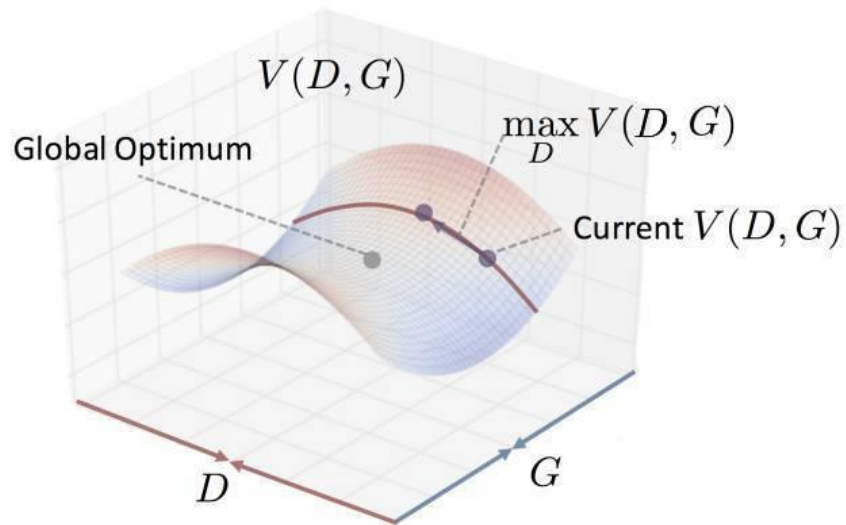
- $D^* = \mathbf{D} \mid \text{JSD}(\mathbf{P}_{\text{data}} \parallel \mathbf{P}_G)$

- Evaluate $\text{JSD}(\mathbf{P}_{\text{data}} \parallel \mathbf{P}_G)$

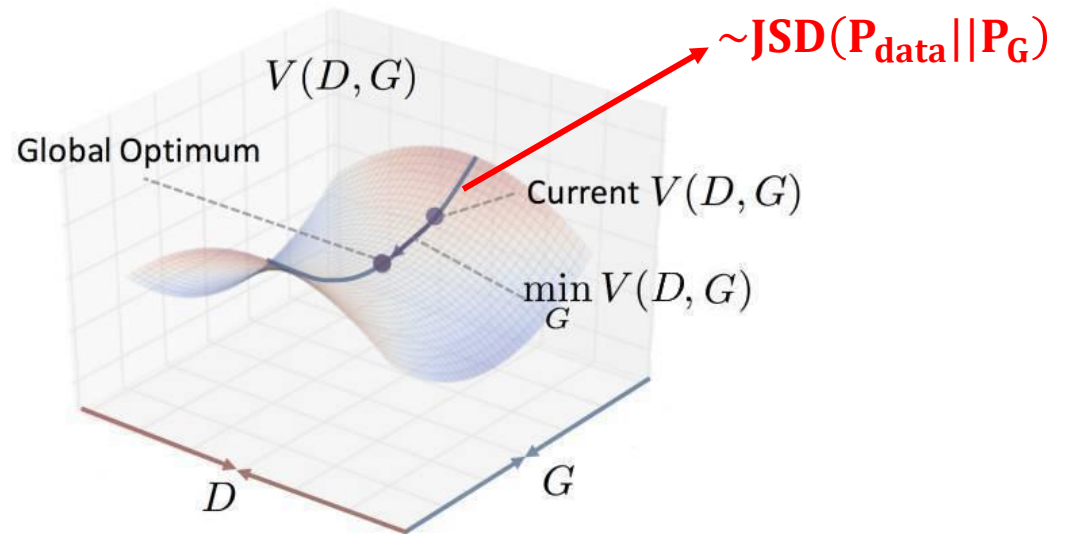
Step 2. Train G

- $G^* = \arg \min_G \text{JSD}(\mathbf{P}_{\text{data}} \parallel \mathbf{P}_G)$

- Minimize $\text{JSD}(\mathbf{P}_{\text{data}} \parallel \mathbf{P}_G)$



- Fix G , and update D



- Fix D , and update G

Weights update

D^* - Evaluate $JSD(P_{data} || P_G)$

- $D^* = D \mid JSD(P_{data} || P_G)$
- $D^* = \arg \max_D V(D, G)$
- $\theta_d \leftarrow \theta_d + \eta \nabla_{\theta_d} V(D, G)$
- Repeat n times to evaluate $JSD(P_{data} || P_G)$ correctly

G^* - Minimize $JSD(P_{data} || P_G)$

- $G^* = \arg \min_G JSD(P_{data} || P_G)$
- $G^* = \arg \min_G \max_D V(D, G)$
 $= \arg \min_G V(D^*, G)$
- $\theta_g \leftarrow \theta_g - \eta \nabla_{\theta_g} V(D, G)$
- Run once to change P_G slightly

Weights update



D^* - Evaluate $JSD(P_{\text{data}} || P_G)$

- $D^* = D \mid JSD(P_{\text{data}} || P_G)$
- $D^* = \arg \max_D V(D, G)$
- $\theta_d \leftarrow \theta_d + \eta \nabla_{\theta_d} V(D, G)$

- $V = E_{x \sim P_{\text{data}}} [\log D(x)] + E_{x \sim P_G} [\log(1 - D(x))]$
- Sample $\{x^1, x^2, \dots, x^m\}$ from $P_{\text{data}}(x)$
- Sample $\{\tilde{x}^1, \tilde{x}^2, \dots, \tilde{x}^m\}$ from generator $P_G(x)$

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

G^* - Minimize $JSD(P_{\text{data}} || P_G)$

- $G^* = \arg \min_G JSD(P_{\text{data}} || P_G)$
- $G^* = \arg \min_G \max_D V(D, G)$
 $= \arg \min_G V(D^*, G)$
- $\theta_g \leftarrow \theta_g - \eta \nabla_{\theta_g} V(D, G)$
- $V = E_{x \sim P_{\text{data}}} [\log D(x)] + E_{x \sim P_G} [\log(1 - D(x))]$
- ~~Sample $\{x^1, x^2, \dots, x^m\}$ from $P_{\text{data}}(x)$~~
- $\{z^1, z^2, \dots, z^m\}$ from the prior $P_{\text{prior}}(z)$
- ~~$V_g = \frac{1}{m} \sum_{i=1}^m \log D(x^i) + \frac{1}{m} \sum_{i=1}^m \log(1 - D(G(z)))$~~
- $V_g = \frac{1}{m} \sum_{i=1}^m \log(1 - D(G(z)))$

Weights update

D^* - Evaluate $JSD(P_{data} || P_G)$

- $D^* = D \mid JSD(P_{data} || P_G)$
- $D^* = \arg \max_D V(D, G)$

- $\theta_d \leftarrow \theta_d + \eta \nabla_{\theta_d} V(D, G)$
- $V_d = \frac{1}{m} \sum_{i=1}^m \log D(x^i) + \frac{1}{m} \sum_{i=1}^m \log(1 - D(\tilde{x}^i))$



- $\theta_d \leftarrow \theta_d - \eta \nabla_{\theta_d} V(D, G)$
- $\bar{V}_d = - \left[\frac{1}{m} \sum_{i=1}^m \log D(x^i) + \frac{1}{m} \sum_{i=1}^m \log(1 - D(\tilde{x}^i)) \right]$
- $\bar{V}_d = \text{BCE Loss}$

G^* - Minimize $JSD(P_{data} || P_G)$

- $G^* = \arg \min_G JSD(P_{data} || P_G)$
- $G^* = \arg \min_G \max_D V(D, G)$
- $= \arg \min_G V(D^*, G)$

- $\theta_g \leftarrow \theta_g - \eta \nabla_{\theta_g} V(D, G)$
- $V_g = \frac{1}{m} \sum_{i=1}^m \log(1 - D(G(z)))$



- $\theta_g \leftarrow \theta_g - \eta \nabla_{\theta_g} V(D, G)$
- $V_g = \frac{1}{m} \sum_{i=1}^m \log(1 - D(G(z)))$
- $V_g = - \text{BCE Loss}$

Algorithm of GAN

- Initialize θ_d for D and θ_g for G
- In each training iteration:

- Sample m real samples $\{x^1, x^2, \dots, x^m\}$ from dataset
- Sample m noise samples $\{z^1, z^2, \dots, z^m\}$ from a distribution
- Generate m fake samples $\{\tilde{x}^1, \tilde{x}^2, \dots, \tilde{x}^m\}$, $\tilde{x}^i = G(z^i)$
- Update discriminator parameters θ_d to minimize

$$\bar{V}_d = - \left(\frac{1}{m} \sum_{i=1}^m \log D(x^i) + \frac{1}{m} \sum_{i=1}^m \log (1 - D(\tilde{x}^i)) \right)$$

$\bar{V}_d = \text{BCE Loss}$

$$\theta_d \leftarrow \theta_d - \eta \nabla_{\theta_d} \bar{V}_d(\theta_d)$$

Train D
Repeat n times

- Sample m noise samples $\{z^1, z^2, \dots, z^m\}$ from a distribution
- Update generator parameters θ_g to minimize

$$V_g = \frac{1}{m} \sum_{i=1}^m \log(1 - D(G(z)))$$

$V_g = - \text{BCE Loss}$

$$\theta_g \leftarrow \theta_g - \eta \nabla_{\theta_g} V_g(\theta_g)$$

Train G
Run once

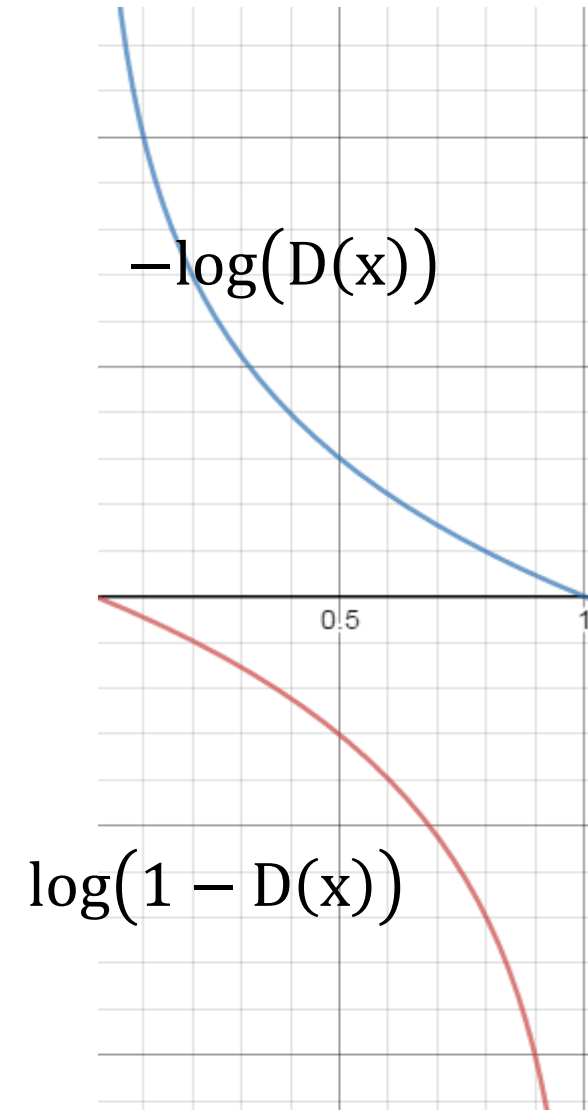
MMGAN & NSGAN

- MMGAN(Minimax Gan)

- $V_g = E_{x \sim P_G} [\log(1 - D(x))]$

- NSGAN(Non-saturating GAN)

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



Algorithm of NSGAN

- Initialize θ_d for D and θ_g for G
- In each training iteration:
 - Sample m real samples $\{x^1, x^2, \dots, x^m\}$ from dataset
 - Sample m noise samples $\{z^1, z^2, \dots, z^m\}$ from a distribution
 - Generate m fake samples $\{\tilde{x}^1, \tilde{x}^2, \dots, \tilde{x}^m\}$, $\tilde{x}^i = G(z^i)$
 - Update discriminator parameters θ_d to minimize
 - $\bar{V}_d = -\left(\frac{1}{m} \sum_{i=1}^m \log D(x^i) + \frac{1}{m} \sum_{i=1}^m \log(1 - D(\tilde{x}^i))\right)$
 - $\theta_d \leftarrow \theta_d - \eta \nabla_{\theta_d} \bar{V}_d(\theta_d)$

Train D
Repeat n times

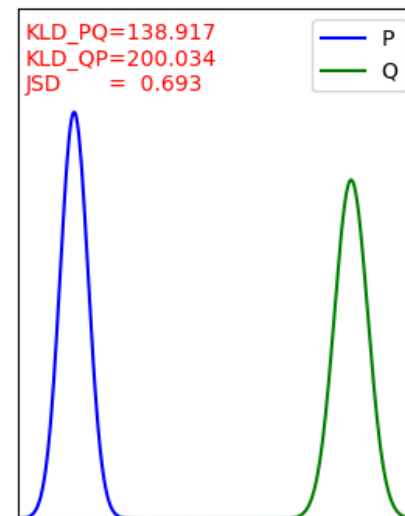
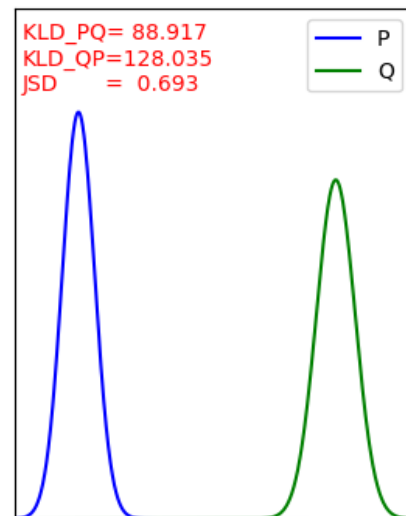
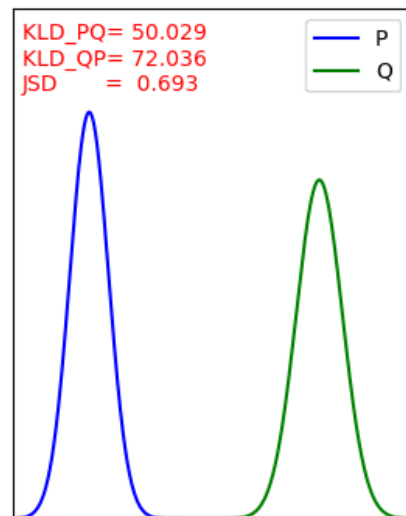
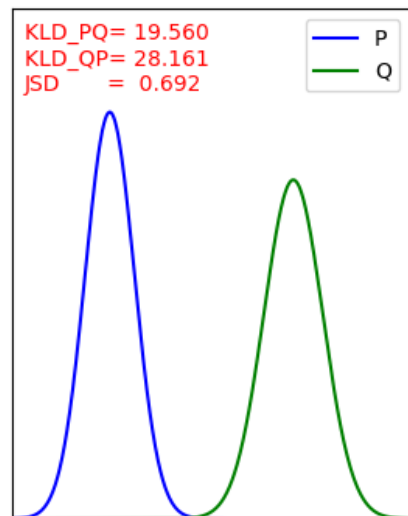
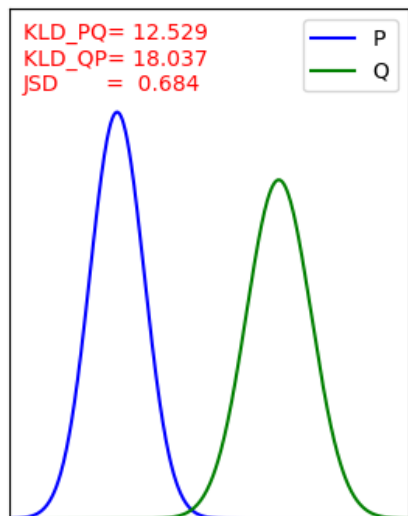
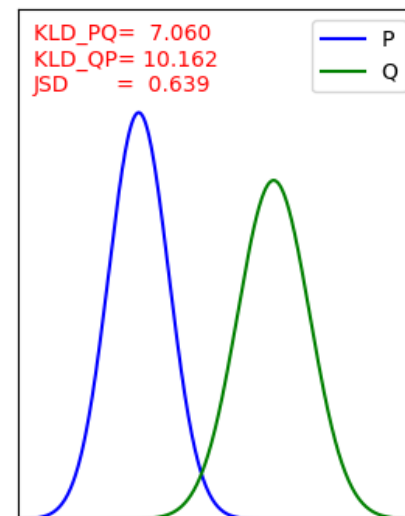
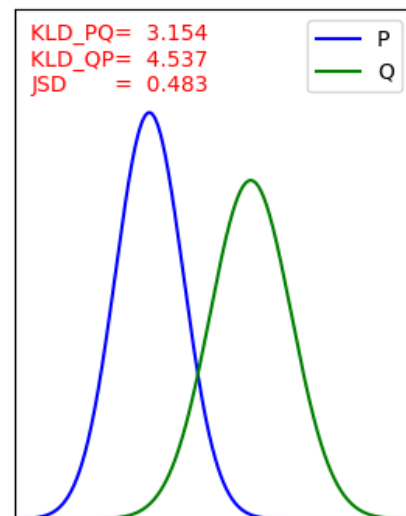
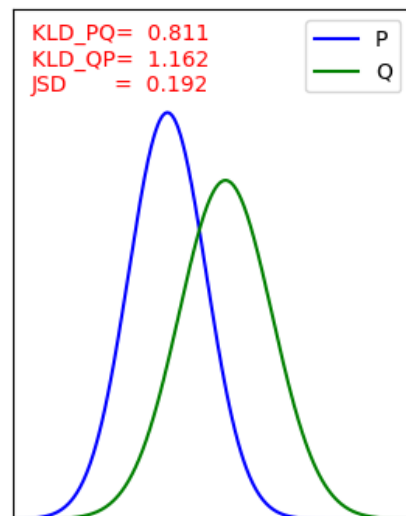
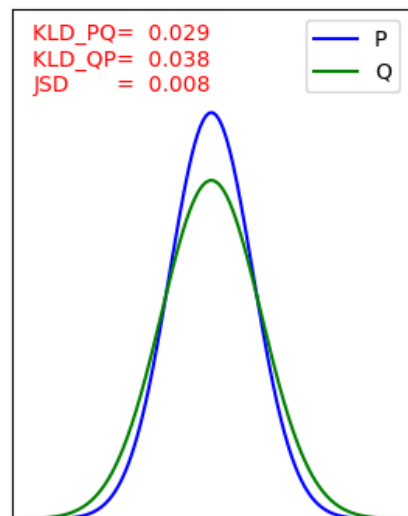
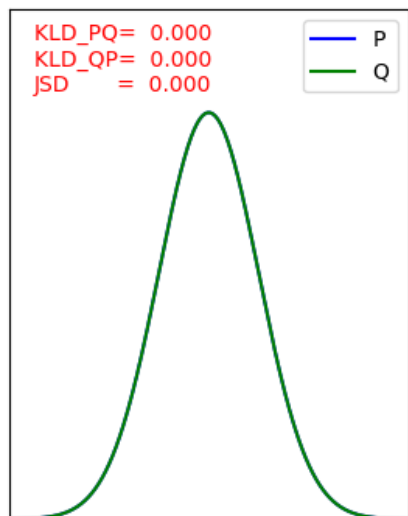
- Sample m noise samples $\{z^1, z^2, \dots, z^m\}$ from a distribution
- Update generator parameters θ_g to minimize
 - $V_g = \frac{1}{m} \sum_{i=1}^m [-\log(D(G(z)))]$
 - $\theta_g \leftarrow \theta_g - \eta \nabla_{\theta_g} V_g(\theta_g)$

Train G
Run once

WGAN

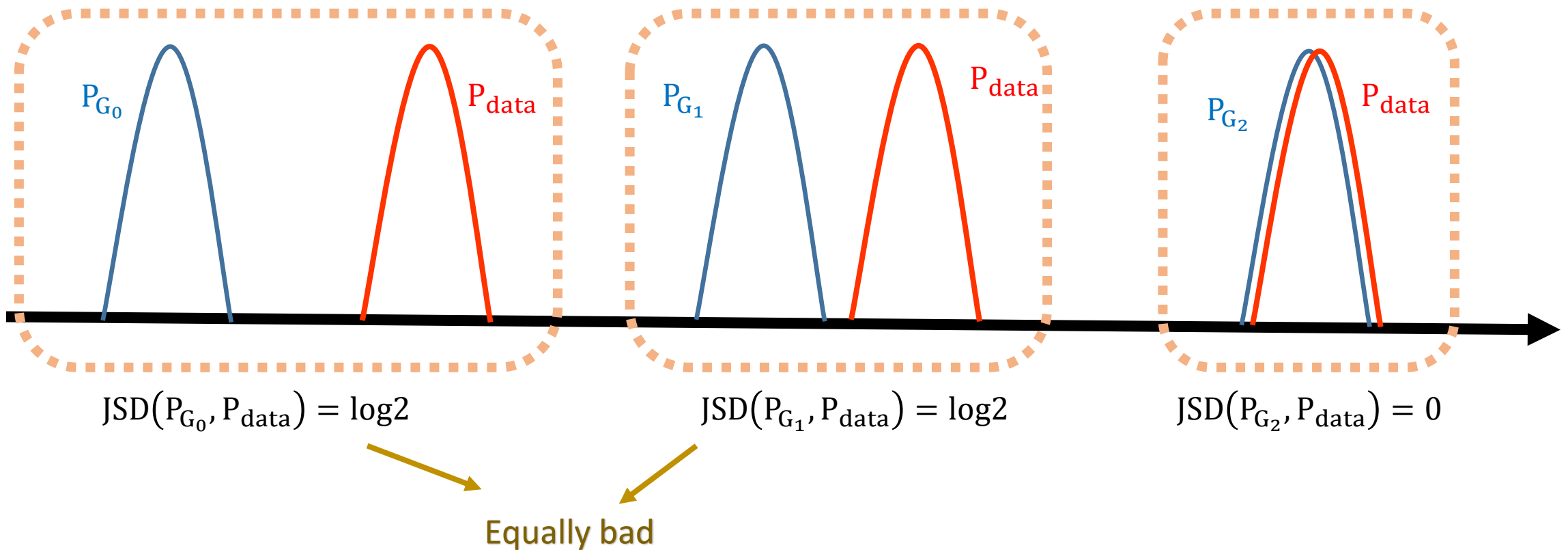
— Is there a more stable GAN?

KLD & JSD



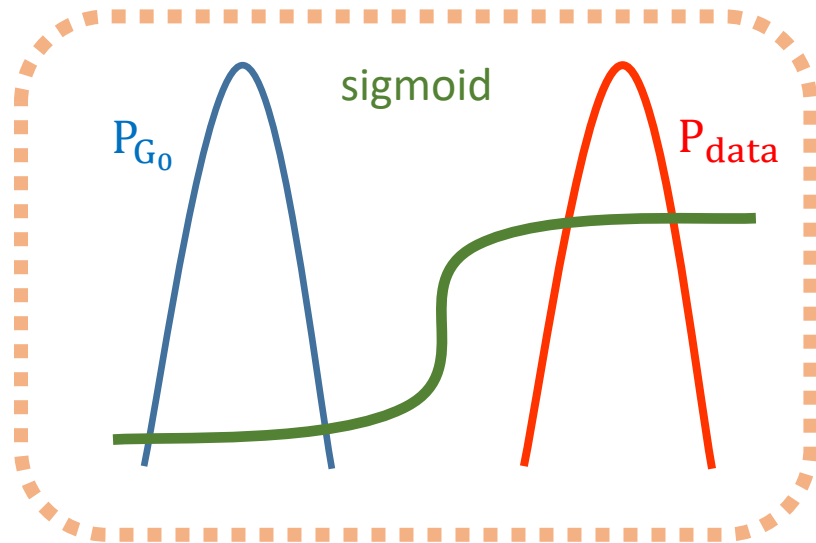
Jensen–Shannon divergence

- $JSD(P||Q) = \frac{1}{2} \text{KLD}(P||M) + \frac{1}{2} \text{KLD}(Q||M)$, where $M = \frac{1}{2}(P + Q)$
- $D_{\text{KL}}(P||Q) = \int_{-\infty}^{\infty} P(x) \log\left(\frac{P(x)}{Q(x)}\right) dx$
- if P_{data} and P_G are not overlapped, $JSD(P||Q) = \log 2$

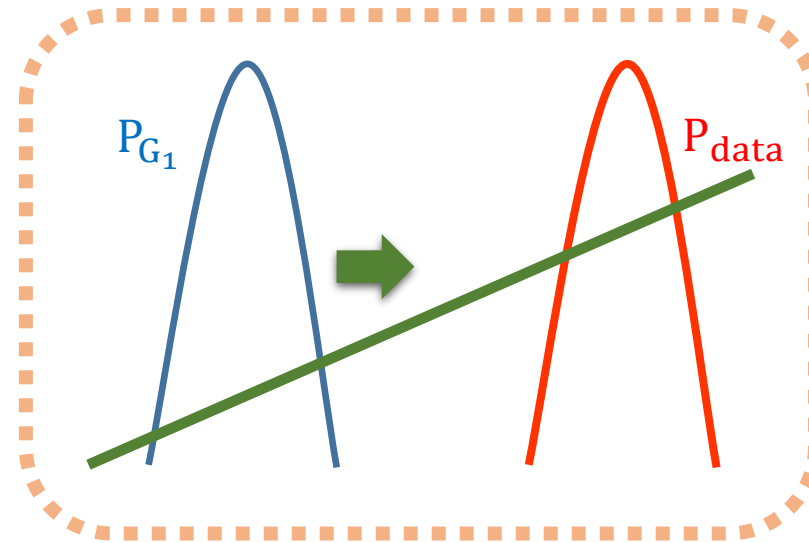


Solution: WGAN(Wasserstein GAN)

- Replace binary classification with linear regression
- Weight clipping
 - Force the parameters w between c and $-c$: if $w > c$, $w = c$; if $w < -c$, $w = -c$



P_{G_0} does not move



Algorithm of GAN

- Initialize θ_d for D and θ_g for G
- In each training iteration:

- Sample m real samples $\{x^1, x^2, \dots, x^m\}$ from dataset
- Sample m noise samples $\{z^1, z^2, \dots, z^m\}$ from a distribution
- Generate m fake samples $\{\tilde{x}^1, \tilde{x}^2, \dots, \tilde{x}^m\}$, $\tilde{x}^i = G(z^i)$
- Update discriminator parameters θ_d to minimize
 - $\bar{V}_d = -\left(\frac{1}{m} \sum_{i=1}^m \log D(x^i) + \frac{1}{m} \sum_{i=1}^m \log(1 - D(\tilde{x}^i))\right)$
 - $\theta_d \leftarrow \theta_d - \eta \nabla_{\theta_d} \bar{V}_d(\theta_d)$

- Sample m noise samples $\{z^1, z^2, \dots, z^m\}$ from a prior distribution
- Update generator parameters θ_g to minimize
 - $V_g = \frac{1}{m} \sum_{i=1}^m [-\log(D(G(z)))]$
 - $\theta_g \leftarrow \theta_g - \eta \nabla_{\theta_g} V_g(\theta_g)$

Algorithm of WGAN

- Initialize θ_d for D and θ_g for G
- In each training iteration:

- Sample m real samples $\{x^1, x^2, \dots, x^m\}$ from dataset
- Sample m noise samples $\{z^1, z^2, \dots, z^m\}$ from a distribution
- Generate m fake samples $\{\tilde{x}^1, \tilde{x}^2, \dots, \tilde{x}^m\}$, $\tilde{x}^i = G(z^i)$
- Update discriminator parameters θ_d to minimize
 - $\bar{V}_d = -\left[\frac{1}{m}\sum_{i=1}^m D(x^i) - \frac{1}{m}\sum_{i=1}^m D(\tilde{x}^i)\right]$
 - $\theta_d \leftarrow \theta_d - \eta \nabla_{\theta_d} \bar{V}_d(\theta_d)$

- Sample m noise samples $\{z^1, z^2, \dots, z^m\}$ from a prior distribution
- Update generator parameters θ_g to minimize
 - $V_g = \frac{1}{m}\sum_{i=1}^m -D(G(z^i))$
 - $\theta_g \leftarrow \theta_g - \eta \nabla_{\theta_g} V_g(\theta_g)$

No sigmoid for the output of D

Weight clipping

GANs

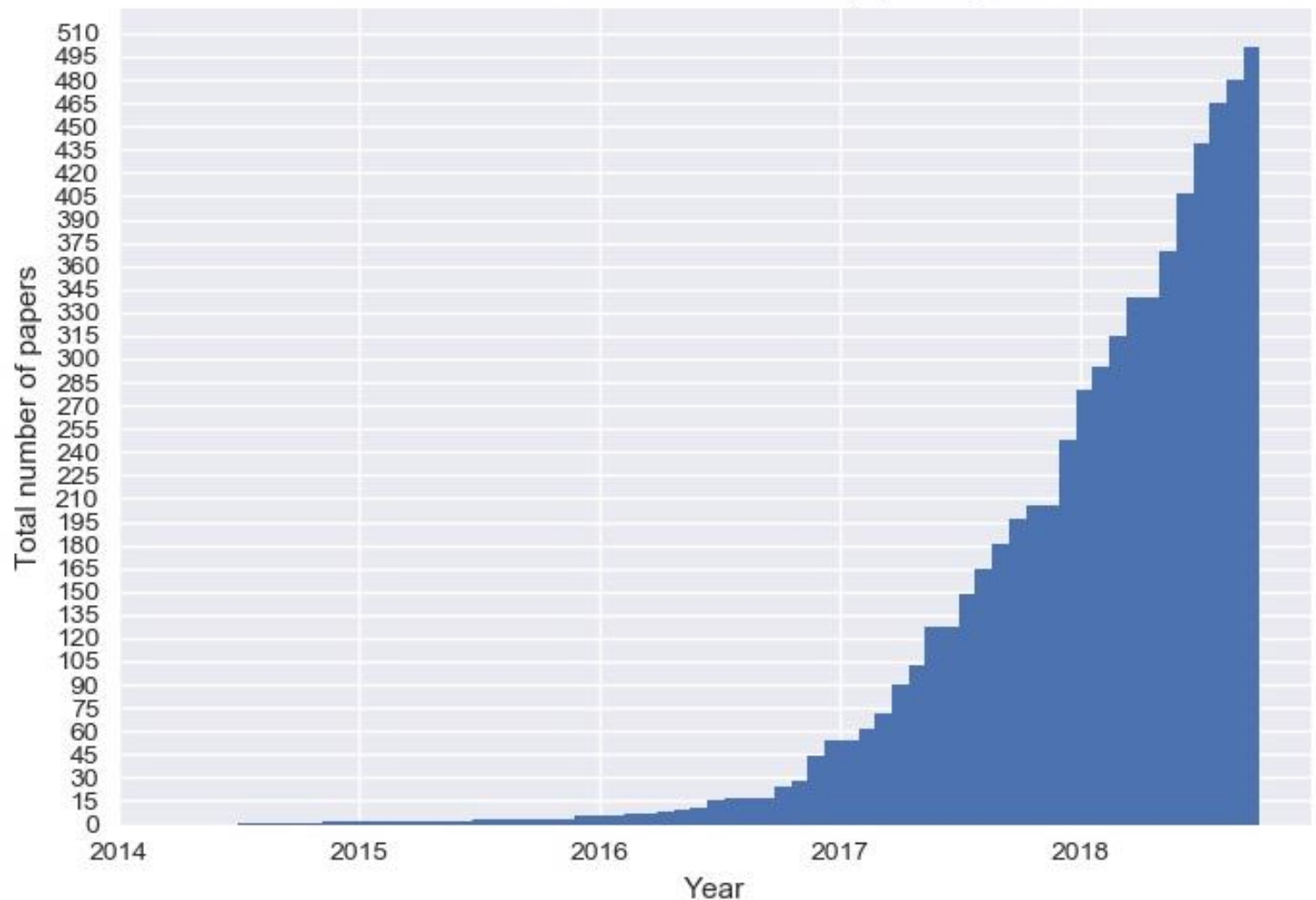
— Are there other GANs?

The GAN Zoo



- modified GAN-CLS - Generate the corresponding image from text description using modified Algorithm
- ModularGAN - Modular Generative Adversarial Networks
- MolGAN - MolGAN: An implicit generative model for small molecular graphs
- MPM-GAN - Message Passing Multi-Agent GANs
- MS-GAN - Temporal Coherency based Criteria for Predicting Video Frames using Deep Multi-Adversarial Networks
- MTGAN - MTGAN: Speaker Verification through Multitasking Triplet Generative Adversarial Net
- MuseGAN - MuseGAN: Symbolic-domain Music Generation and Accompaniment with Multi-tr Generative Adversarial Networks
- MV-BiGAN - Multi-view Generative Adversarial Networks
- N2RPP - N2RPP: An Adversarial Network to Rebuild Plantar Pressure for ACLD Patients
- NAN - Understanding Humans in Crowded Scenes: Deep Nested Adversarial Learning and A N Multi-Human Parsing
- NCE-GAN - Dihedral angle prediction using generative adversarial networks
- ND-GAN - Novelty Detection with GAN
- NetGAN - NetGAN: Generating Graphs via Random Walks
- OGAN - One-Class Adversarial Nets for Fraud Detection
- OptionGAN - OptionGAN: Learning Joint Reward-Policy Options using Generative Adversarial I Reinforcement Learning
- ORGAN - Objective-Reinforced Generative Adversarial Networks (ORGAN) for Sequence Gener
- ORGAN - 3D Reconstruction of Incomplete Archaeological Objects Using a Generative Adversa
- OT-GAN - Improving GANs Using Optimal Transport
- PacGAN - PacGAN: The power of two samples in generative adversarial networks
- PAN - Perceptual Adversarial Networks for Image-to-Image Transformation
- PassGAN - PassGAN: A Deep Learning Approach for Password Guessing
- PD-WGAN - Primal-Dual Wasserstein GAN
- Perceptual GAN - Perceptual Generative Adversarial Networks for Small Object Detection
- PGAN - Probabilistic Generative Adversarial Networks
- PGD-GAN - Solving Linear Inverse Problems Using GAN Priors: An Algorithm with Provable Gu

Cumulative number of named GAN papers by month

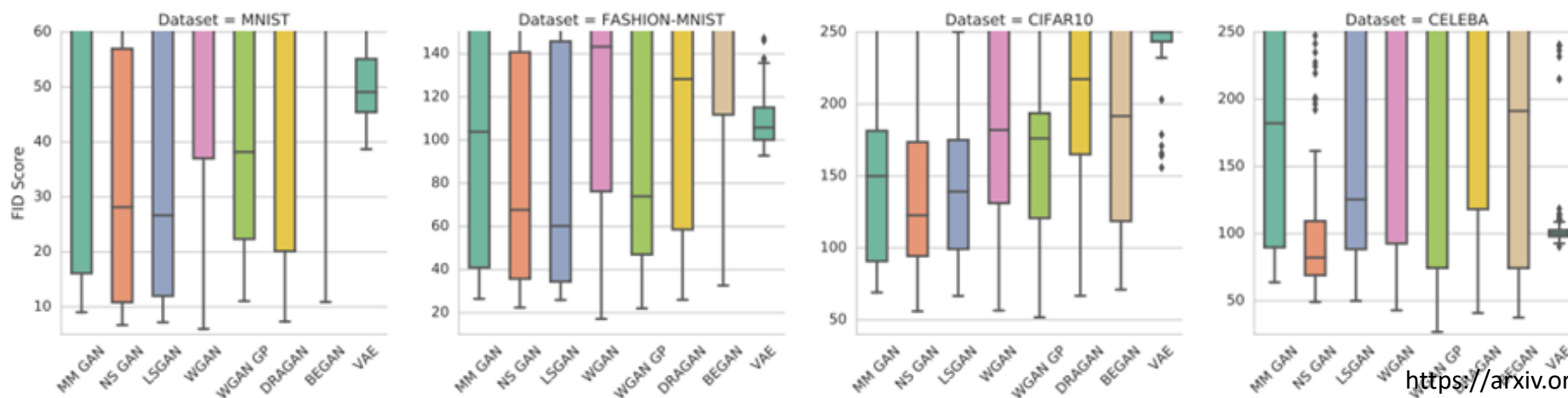


<https://github.com/hindupuravinash/the-gan-zoo>

GANs & Performance



GAN	DISCRIMINATOR LOSS	GENERATOR LOSS
MM GAN	$\mathcal{L}_D^{\text{GAN}} = -\mathbb{E}_{x \sim p_d} [\log(D(x))] + \mathbb{E}_{\hat{x} \sim p_g} [\log(1 - D(\hat{x}))]$	$\mathcal{L}_G^{\text{GAN}} = -\mathcal{L}_D^{\text{GAN}}$
NS GAN	$\mathcal{L}_D^{\text{NSGAN}} = \mathcal{L}_D^{\text{GAN}}$	$\mathcal{L}_G^{\text{NSGAN}} = \mathbb{E}_{\hat{x} \sim p_g} [\log(D(\hat{x}))]$
WGAN	$\mathcal{L}_D^{\text{WGAN}} = -\mathbb{E}_{x \sim p_d} [D(x)] + \mathbb{E}_{\hat{x} \sim p_g} [D(\hat{x})]$	$\mathcal{L}_G^{\text{WGAN}} = \mathcal{L}_D^{\text{WGAN}}$
WGAN GP	$\mathcal{L}_D^{\text{WGAN}} = \mathcal{L}_D^{\text{WGAN}} + \lambda \mathbb{E}_{\hat{x} \sim p_g} [(\ \nabla D(\alpha x + (1 - \alpha)\hat{x})\ _2 - 1)^2]$	$\mathcal{L}_G^{\text{WGAN}} = -\mathbb{E}_{\hat{x} \sim p_g} [D(\hat{x})]$
LS GAN	$\mathcal{L}_D^{\text{LSGAN}} = -\mathbb{E}_{x \sim p_d} [(D(x) - 1)^2] + \mathbb{E}_{\hat{x} \sim p_g} [D(\hat{x})^2]$	$\mathcal{L}_G^{\text{LSGAN}} = -\mathbb{E}_{\hat{x} \sim p_g} [(D(\hat{x}) - 1)^2]$
DRAGAN	$\mathcal{L}_D^{\text{DRAGAN}} = \mathcal{L}_D^{\text{GAN}} + \lambda \mathbb{E}_{\hat{x} \sim p_d + \mathcal{N}(0, c)} [(\ \nabla D(\hat{x})\ _2 - 1)^2]$	$\mathcal{L}_G^{\text{DRAGAN}} = -\mathcal{L}_D^{\text{NSGAN}}$
BEGAN	$\mathcal{L}_D^{\text{BEGAN}} = \mathbb{E}_{x \sim p_d} [\ x - \text{AE}(x)\ _1] - k_t \mathbb{E}_{\hat{x} \sim p_g} [\ \hat{x} - \text{AE}(\hat{x})\ _1]$	$\mathcal{L}_G^{\text{BEGAN}} = \mathbb{E}_{\hat{x} \sim p_g} [\ \hat{x} - \text{AE}(\hat{x})\ _1]$



The smaller, the better

<https://arxiv.org/abs/1711.10337>
<https://arxiv.org/pdf/1706.08500>

Evaluation

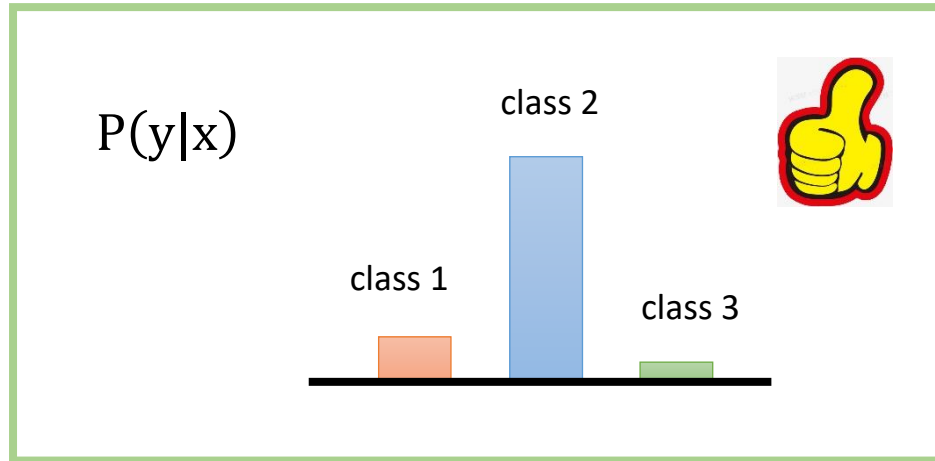
— How good is a GAN?

Inception Score

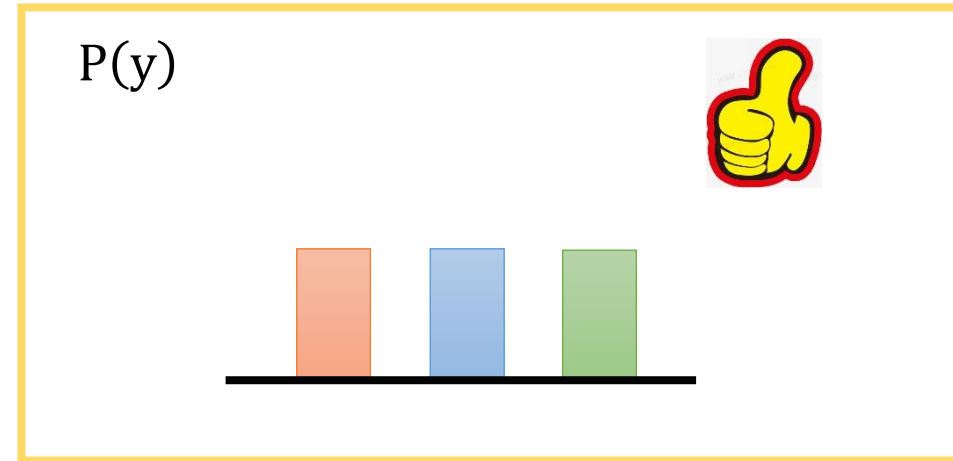
x: samples
y: classes



Quality



Diversity



$$\text{Inception Score} = \underbrace{\sum_{x,y} P(y|x) \log P(y|x)}_{\text{Quality}} - \underbrace{\sum_y P(y) \log P(y)}_{\text{Diversity}}$$

Negative entropy of $P(y|x)$

Entropy of $P(y)$

Tricks

— What are the tricks to make a GAN work?

Tips and tricks to make GANs work



- Normalize the samples between -1 and 1 and use Tanh as the last layer of the generator
- Sample from a Gaussian distribution rather than a Uniform distribution as the input of generator
- Replace Relu with LeakyReLU both in discriminator and generator
- Use soft and noisy labels by replacing label 0 with a random number between 0.0 and 0.3 and replacing label 1 with a random number between 0.8 and 1.2
- Increase discriminator and generator network capacity with more filters and more layers
- Add artificial noise to the real and generated samples before feeding them into the discriminator, and decay over time
- Feed real samples and generated samples into the discriminator separately in different batches
- Feed the discriminator with all recent generated samples to avoid overfitting
- Try different loss functions, though no single loss function always performs best among all datasets
- Stabilize training with Batch Normalization
- Adjust the learning rate to ease model collapse
-

[arXiv:1609.05158](https://arxiv.org/abs/1609.05158) [arXiv:1609.04468](https://arxiv.org/abs/1609.04468)
<https://github.com/soumith/ganhacks>

Practice

— How to code the first line?

GAN toy

https://github.com/pp-jiang/GANs4Joint_GlueX_EIC_PANDA_ML_Workshop.git



```
import numpy as np; import matplotlib.pyplot as plt
```

```
import torch; from torch.autograd import Variable; import torch.nn as nn; import torch.nn.functional as F
```

```
Tensor = torch.FloatTensor
```

```
numSample=10; numEpoch=10000; num_DG=3; dimSample = 32; dimLatent=10; batchSize=128; bceLoss = torch.nn.BCELoss()
```

```
G = nn.Sequential(nn.Linear(dimLatent, 128),nn.LeakyReLU(0.2, inplace=True),nn.Linear(128, 256),nn.LeakyReLU(0.2, inplace=True),  
nn.Linear(256, 512),nn.LeakyReLU(0.2, inplace=True),nn.Linear(512, 1024),nn.LeakyReLU(0.2, inplace=True),nn.Linear(1024,  
dimSample),nn.Tanh())
```

```
D = nn.Sequential(nn.Linear(dimSample, 512),nn.LeakyReLU(0.2, inplace=True),nn.Linear(512, 256),nn.LeakyReLU(0.2,  
inplace=True),nn.Linear(256, 1),nn.Sigmoid())
```

```
dataSet=Variable(torch.sin(torch.mm(torch.linspace(1,1.2,numSample).unsqueeze(1),torch.linspace(-  
2.5,3.5,dimSample).unsqueeze(0))),requires_grad=False)
```

```
optimG = torch.optim.Adam(G.parameters(),lr=0.00001); optimD = torch.optim.SGD(D.parameters(),lr=0.001) #optimD =  
torch.optim.Adam(D.parameters(),lr=0.00001)
```

```
lossDList=[]; lossGList=[]
```

GAN toy - Training



```
for epoch in range(numEpoch):
```

```
    xReal=dataset[torch.randint(0,numSample,(batchSize,)),:] # xReal
```

```
    z = torch.normal(0, 1, size=(batchSize, dimLatent), requires_grad=False, out=None) # latent noise
```

```
    yReal= Variable(Tensor(batchSize, 1).fill_(1.0), requires_grad=False) # real labels
```

```
    yFake= Variable(Tensor(batchSize, 1).fill_(0.0), requires_grad=False) # fake labels
```

Train Discriminator

```
    optimD.zero_grad()
```

```
    lossD=(bceLoss(D(G(z)),yFake)+bceLoss(D(xReal),yReal))/2. # lossG = BCE Loss
```

```
    lossD.backward(); optimD.step()
```

Train Generator

```
    if epoch%num_DG==0: # Repeat num_DG times training D and then run once training G
```

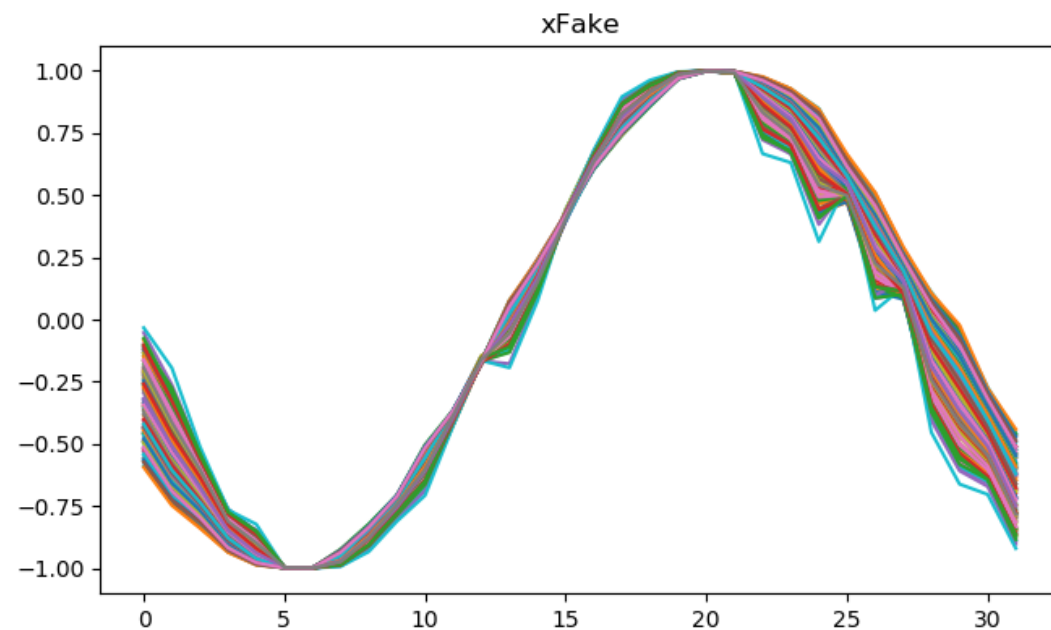
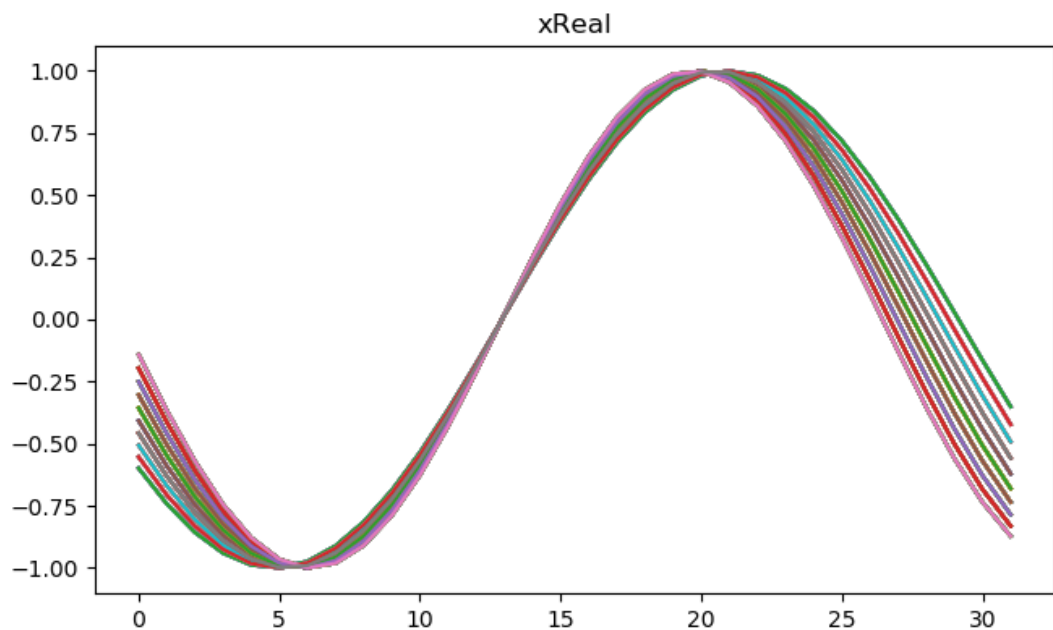
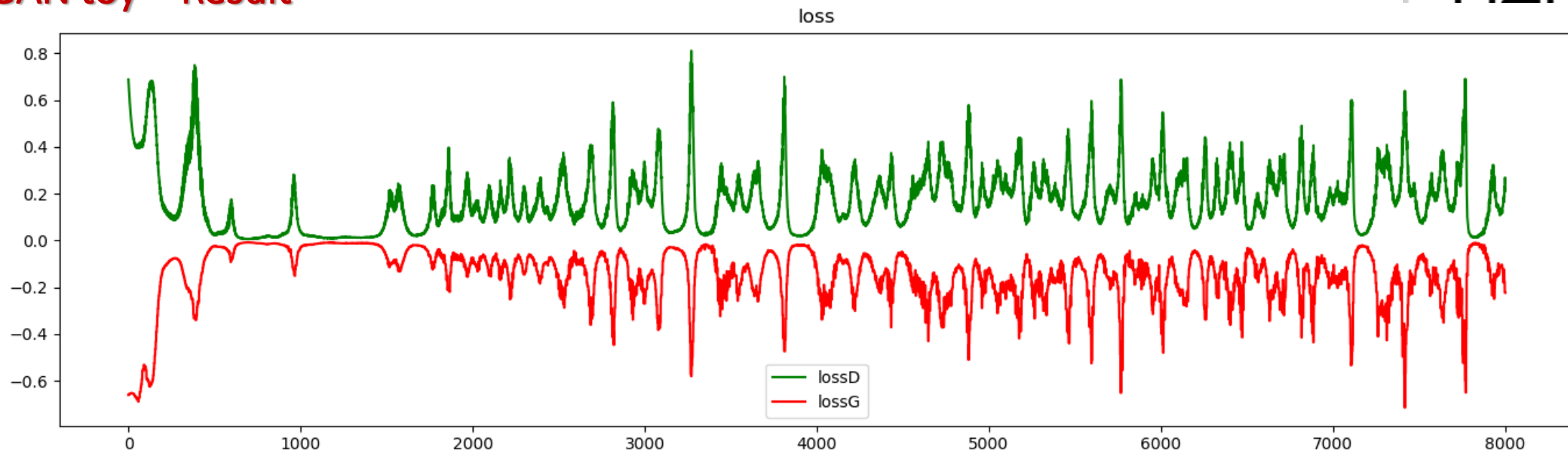
```
        optimG.zero_grad()
```

```
        lossG=-bceLoss(D(G(z)),yFake) # lossG = -BCE Loss
```

```
        lossG.backward(); optimG.step()
```

```
    lossDListot.append(lossD.item()); lossGList.append(lossG.item())
```


GAN toy – Result



Thank you for your attention!

