# Lecture 4

# **Generative Adversarial Networks**

6.S978 Deep Generative Models

Kaiming He
Fall 2024, EECS, MIT

# Overview

- Generative Adversarial Networks (GAN)

- Wasserstein GAN (W-GAN)

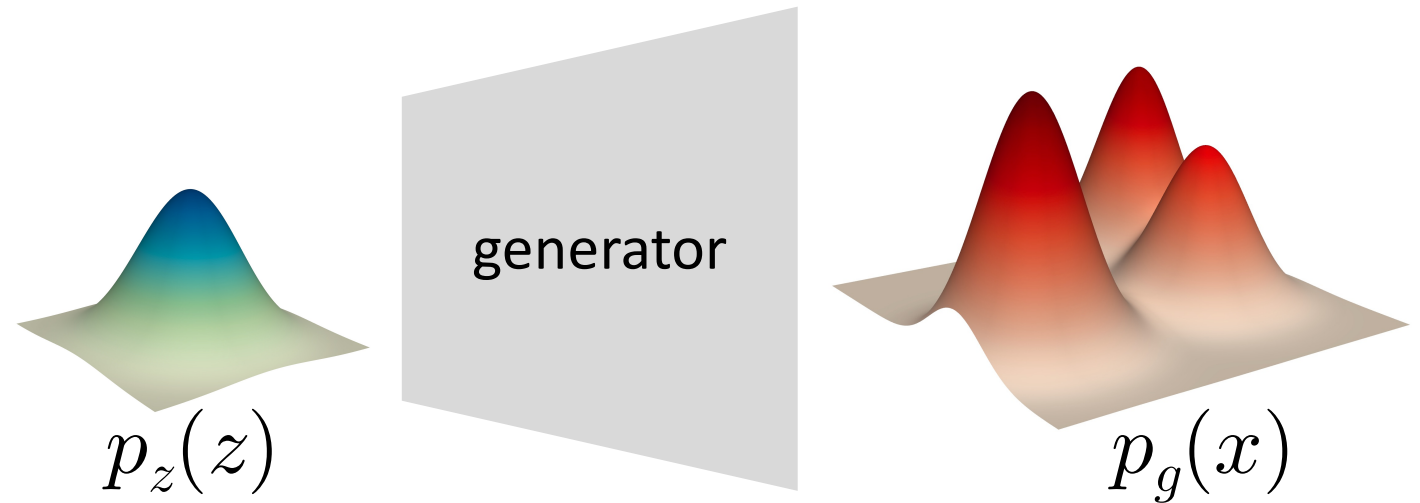- Adversary as a Loss Function

# Generative Adversarial Networks (GAN)

# Introduction

- "**Generative**"
  - "Discriminative" was dominant back then

- "**Adversarial**"
  - Generative models w/ discriminative models
  - Min-max process

- "**Networks**"
  - SGD + backprop for problem solving

# Recap: Latent Variable Models
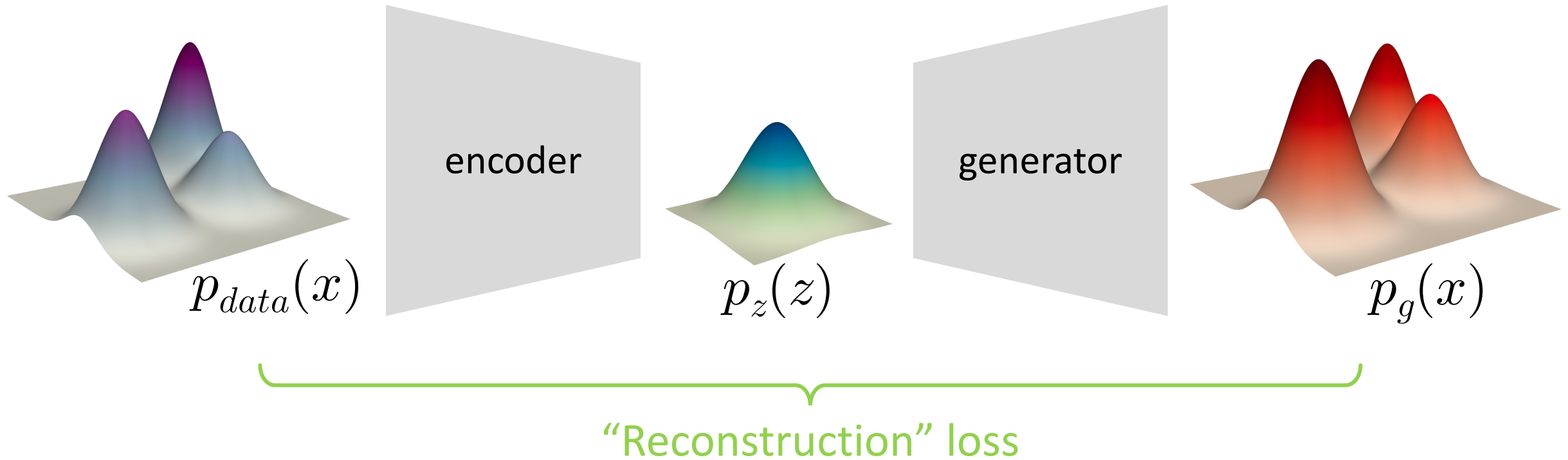
**Represent a distribution by a neural network:**

- $z$ - latent variables

- $x$ - observed variables



$$p_z(z) \qquad \text{generator} \qquad p_g(x)$$

# Recap: Variational Autoencoder (VAE)

**Autoencoding distributions**:

"Encoding" data distribution $p_{data}$ into latent distribution $p_z$



$p_{data}(x)$ — encoder — $p_z(z)$ — generator — $p_g(x)$

"Reconstruction" loss

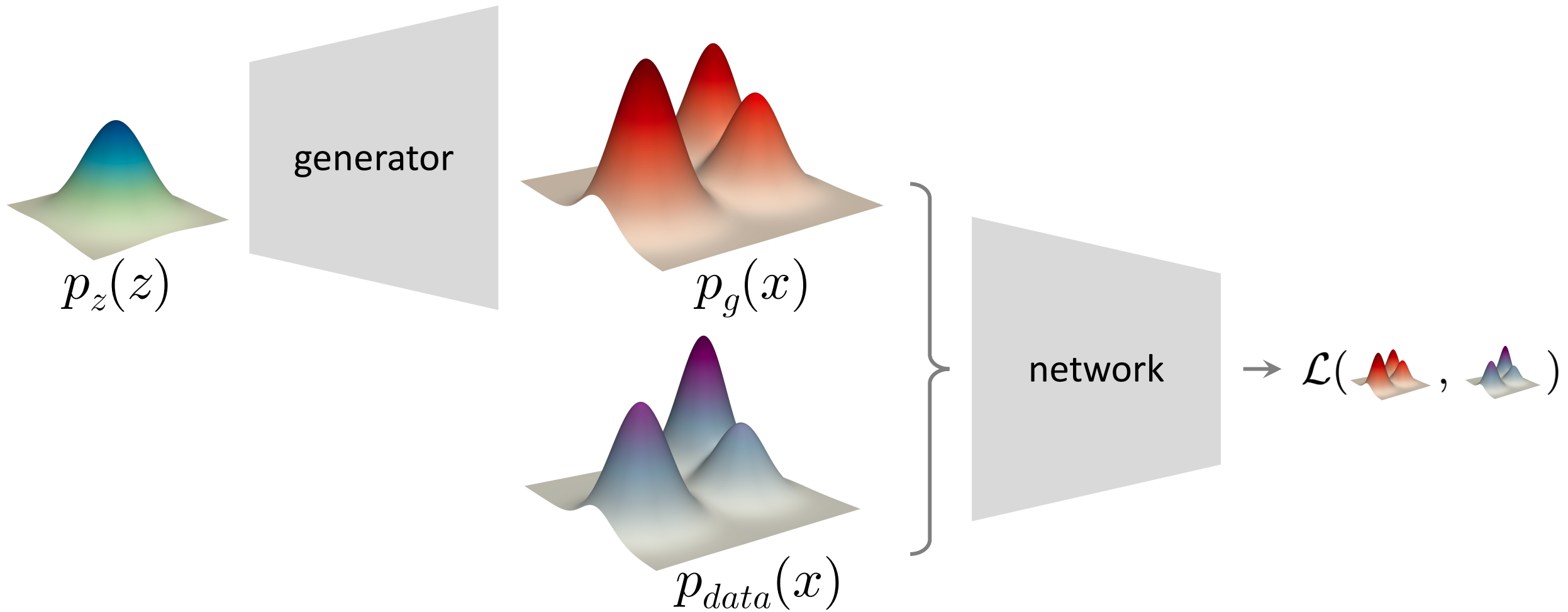# What's the implication of a "*reconstruction*" loss?

- Elements (e.g., pixels) are **independently** distributed
- Each element follows a **simple** distribution (Gaussian/Bernoulli/...)

Assumptions are too strict for **high-dim** variables

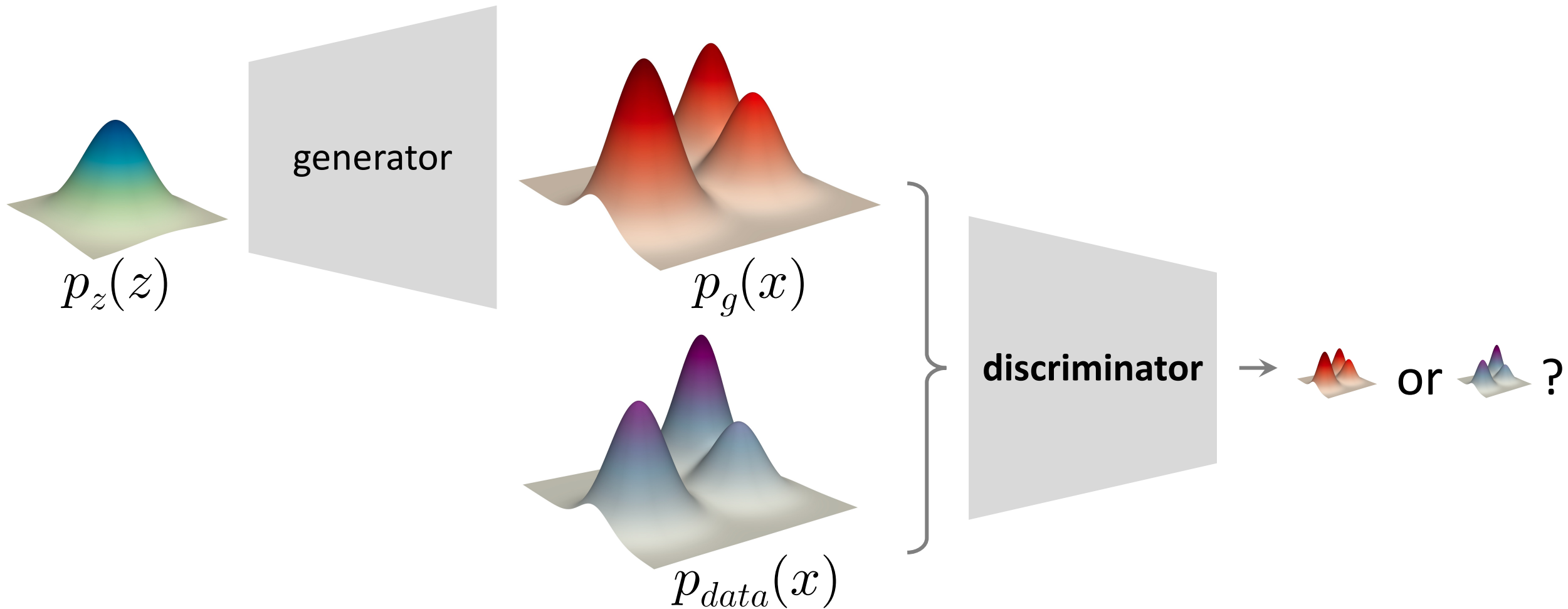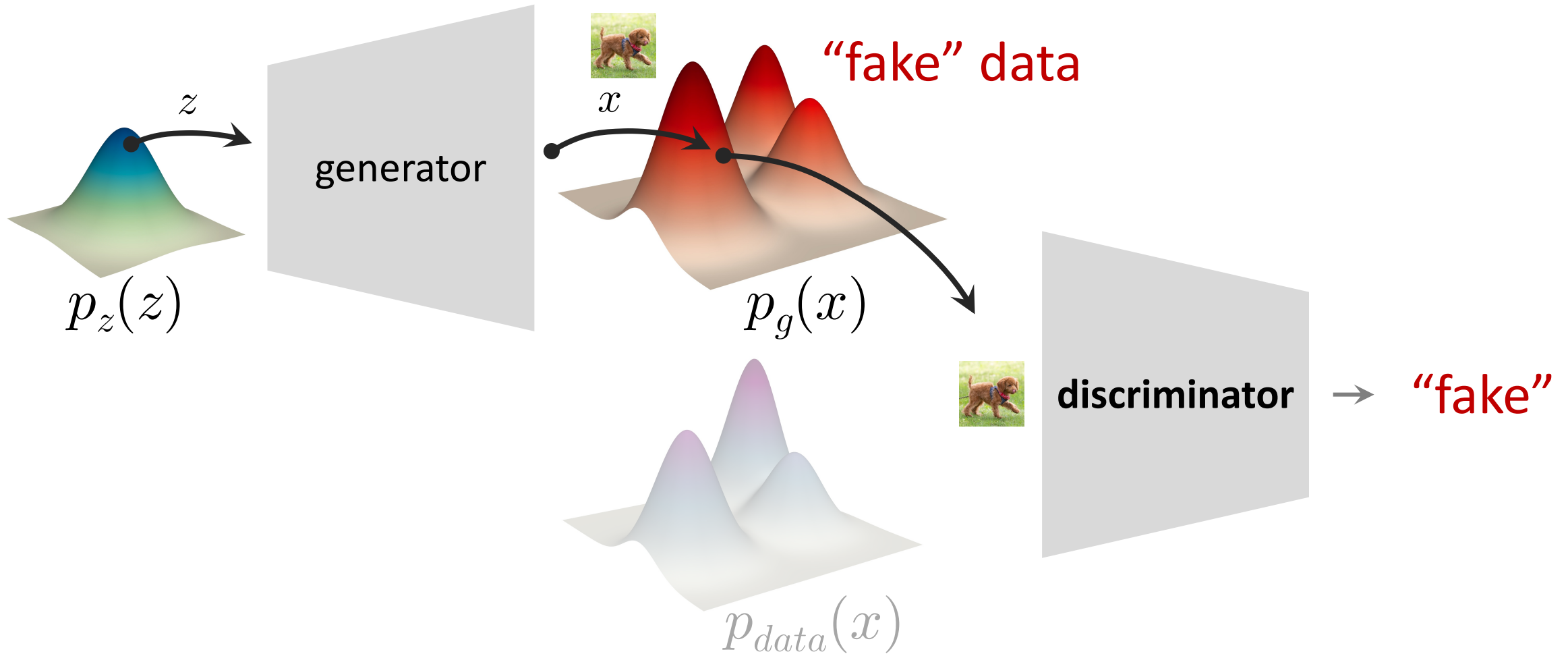*Can we measure the distribution difference in another way?*
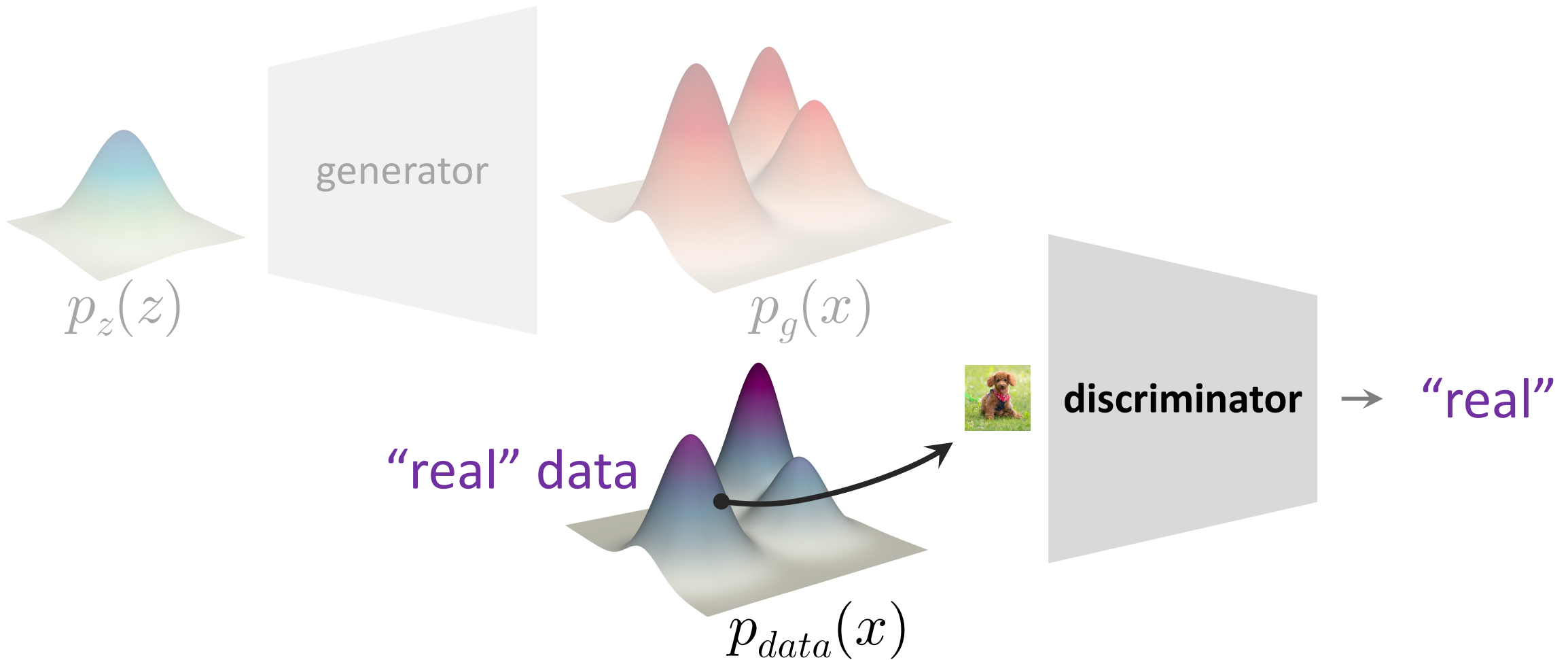
# Generative Adversarial Networks

Representing **distribution difference** by a neural network

# Generative Adversarial Networks

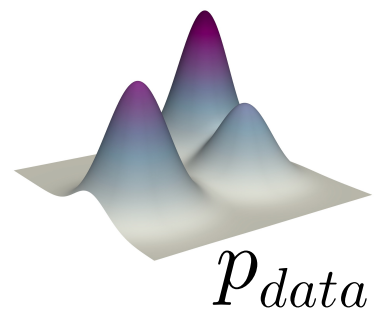Representing **distribution difference** by a neural network
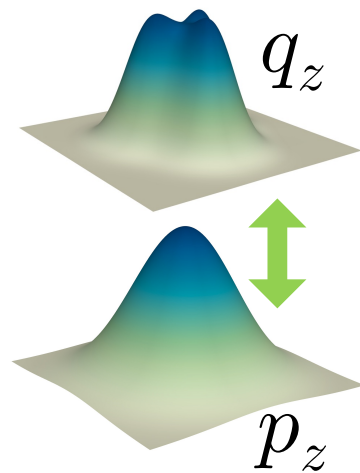
# Generative Adversarial Networks

Representing **distribution difference** by a neural network



$p_z(z)$

generator

$x$

"fake" data

$p_g(x)$

$p_{data}(x)$

**discriminator** → "fake"

# Generative Adversarial Networks

Representing **distribution difference** by a neural network



$p_z(z)$

generator

$p_g(x)$

"real" data

$p_{data}(x)$

**discriminator** $\rightarrow$ "real"

**VAE**

encoder $\quad$ $q_z$ $\quad$ $p_z$ $\quad$ decoder $\quad$ $p_g$

$p_{data}$

**GAN**

generator $\quad$ $p_g$ $\quad$ $p_{data}$ $\quad$ discriminator $\quad$ or $\quad$ ?

$p_z$

**VAE**
generation

$q_z$

encoder

$p_{data}$

$p_z$

decoder

$p_g$

**GAN**
generation

$p_z$

generator

$p_g$

$p_{data}$

discriminator or ?

# Adversarial Objective

$$\min_{G} \max_{D} \mathcal{L}(D, G) = \mathbb{E}_{x \sim p_{\text{data}}}[\log D(x)] + \mathbb{E}_{z \sim p_z}[\log(1 - D(G(z)))]$$

min-max process

(vs. EM's max-max process)

# Adversarial Objective: D-step

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

$D$-step: fix $G$, optimize $D$

# Adversarial Objective: D-step

$$\max_{D} \mathcal{L}(D) = \mathbb{E}_{x \sim p_{\text{data}}}[\log D(x)] + \mathbb{E}_{z \sim p_z}[\log(1 - D(G(z)))]$$

push to 1          push to 0

$D$-step: fix $G$, optimize $D$

- $D$ to classify real or fake

- binary logistic regression (sigmoid + cross-entropy)

# Adversarial Objective: D-step

$$\max_D \mathcal{L}(D) = \mathbb{E}_{x \sim p_{\text{data}}}[\log D(x)] + \mathbb{E}_{x \sim p_g}[\log(1 - D(x))]$$

push to 1                    push to 0

$D$-step: fix $G$, optimize $D$

- $D$ to classify real or fake

- binary logistic regression (sigmoid + cross-entropy)

# Adversarial Objective: G-step

$$\min_{G} \max_{D} \mathcal{L}(D, G) = \mathbb{E}_{x \sim p_{\text{data}}}[\log D(x)] + \mathbb{E}_{z \sim p_z}[\log(1 - D(G(z)))]$$

$G$-step: fix $D$, optimize $G$
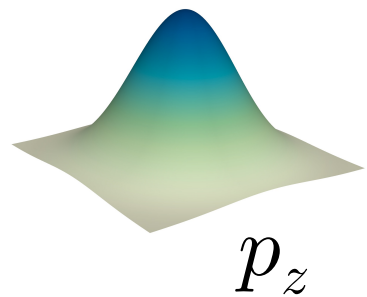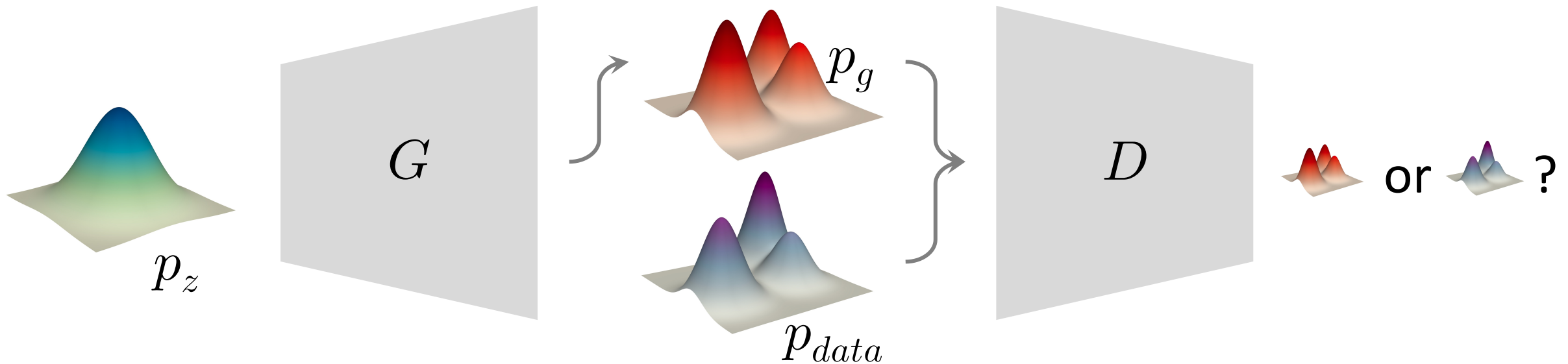
# Adversarial Objective: G-step

$$\min_G \max_D \mathcal{L}(D, G) = \mathbb{E}_{x \sim p_{\text{data}}}[\log D(x)] + \mathbb{E}_{z \sim p_z}[\log(1 - D(G(z)))]$$
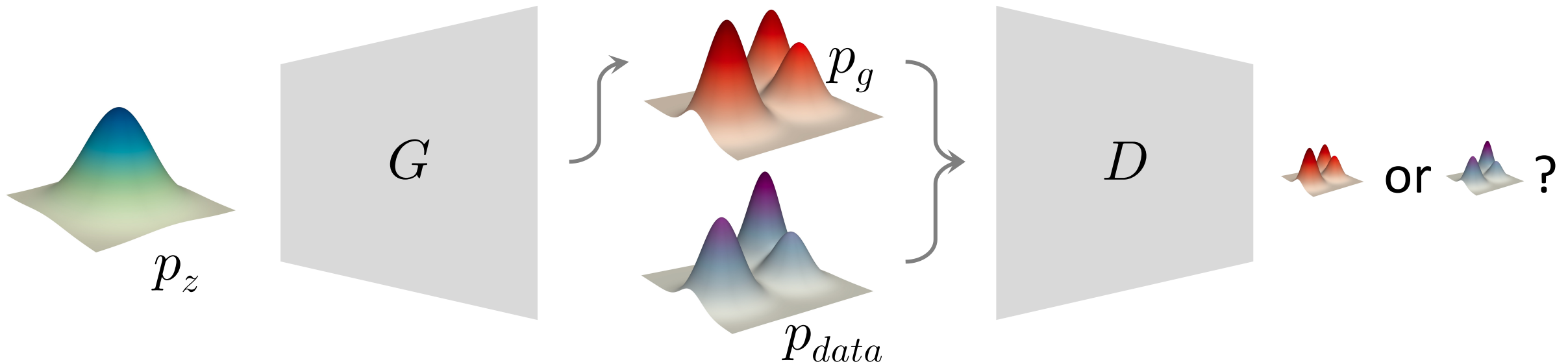
$G$-step: fix $D$, optimize $G$

# Adversarial Objective: G-step

$$\min_G \mathcal{L}(G) = \mathbb{E}_{z \sim p_z}[\log(1 - D(G(z)))]$$

<span style="color:red">push to 1</span>

$G$-step: fix $D$, optimize $G$

- generate fake data such that $D$ classifies it as "real"
- $G$ to "confuse" $D$

# Adversarial Objective: G-step

a "flip" trick:

$$\underset{G}{\overset{\max}{\min}} \mathcal{L}(G) = \mathbb{E}_{z \sim p_z}[\log(1 - D(G(z)))]$$

push to 1

$G$-step: fix $D$, optimize $G$

- generate fake data such that $D$ classifies it as "real"
- $G$ to "confuse" $D$

# Adversarial Objective: G-step

a "flip" trick:

$$\underset{G}{\overset{\max}{\min}} \mathcal{L}(G) = \mathbb{E}_{z \sim p_z}[\log(1 - D(G(z)))]$$

push to 1

Early in training:
- $G$ is poor
- $D(G)$ is near 0



weak grad

strong grad

log(1-D(G(x)))
log(D(G(x)))

D(G(x))

# GAN algorithm annotated

**Algorithm 1** Minibatch stochastic gradient descent training of generative adversarial nets. The number of steps to apply to the discriminator, $k$, is a hyperparameter. We used $k = 1$, the least expensive option, in our experiments.

<span style="color:red">minibatch SGD</span>

**for** number of training iterations **do**
  **for** $k$ steps **do**
    • Sample minibatch of $m$ noise samples $\{z^{(1)}, \ldots, z^{(m)}\}$ from noise prior $p_g(z)$.
    • Sample minibatch of $m$ examples $\{x^{(1)}, \ldots, x^{(m)}\}$ from data generating distribution $p_{\text{data}}(x)$.
    • Update the discriminator by ascending its stochastic gradient:

$$\nabla_{\theta_d} \frac{1}{m} \sum_{i=1}^{m} \left[ \log D\left(x^{(i)}\right) + \log\left(1 - D\left(G\left(z^{(i)}\right)\right)\right) \right].$$

  **end for**
  • Sample minibatch of $m$ noise samples $\{z^{(1)}, \ldots, z^{(m)}\}$ from noise prior $p_g(z)$.
  • Update the generator by descending its stochastic gradient:

$$\nabla_{\theta_g} \frac{1}{m} \sum_{i=1}^{m} \log\left(1 - D\left(G\left(z^{(i)}\right)\right)\right).$$

**end for**
The gradient-based updates can use any standard gradient-based learning rule. We used momentum in our experiments.

# GAN algorithm annotated

**for** number of training iterations **do**

    **for** $k$ steps **do**

- Sample minibatch of $m$ noise samples $\{z^{(1)}, \ldots, z^{(m)}\}$ from noise prior $p_g(z)$.
- Sample minibatch of $m$ examples $\{x^{(1)}, \ldots, x^{(m)}\}$ from data generating distribution $p_{\text{data}}(x)$.
- Update the discriminator by ascending its stochastic gradient:

$$\nabla_{\theta_d} \frac{1}{m} \sum_{i=1}^{m} \left[ \log D\left(x^{(i)}\right) + \log\left(1 - D\left(G\left(z^{(i)}\right)\right)\right) \right].$$
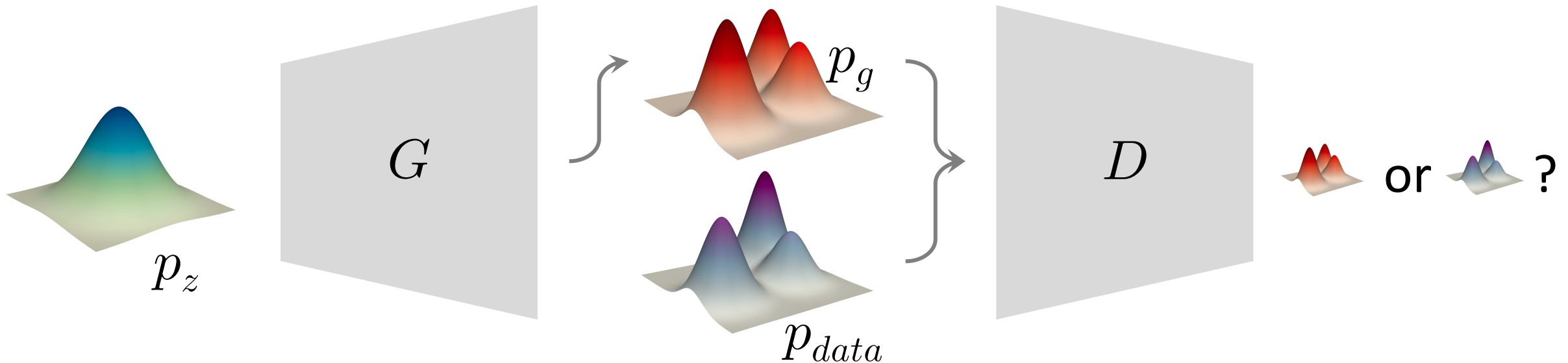
    **end for**

- Sample minibatch of $m$ noise samples $\{z^{(1)}, \ldots, z^{(m)}\}$ from noise prior $p_g(z)$.
- Update the generator by descending its stochastic gradient:

$$\nabla_{\theta_g} \frac{1}{m} \sum_{i=1}^{m} \log\left(1 - D\left(G\left(z^{(i)}\right)\right)\right).$$

**end for**

The gradient-based updates can use any standard gradient-based learning rule. We used momentum in our experiments.



$p_z$    $G$    $p_g$    $p_{data}$    $D$    or    ?

# GAN algorithm annotated

**Algorithm 1** Minibatch stochastic gradient descent training of generative adversarial nets. The number of steps to apply to the discriminator, $k$, is a hyperparameter. We used $k = 1$, the least expensive option, in our experiments.

---

**for** number of training iterations **do**

  **for** $k$ steps **do**

    • Sample minibatch of $m$ noise samples $\{z^{(1)}, \ldots, z^{(m)}\}$ from noise prior $p_g(z)$.

    • Sample minibatch of $m$ examples $\{x^{(1)}, \ldots, x^{(m)}\}$ from data generating distribution $p_{\text{data}}(x)$.

    • Update the discriminator by ascending its stochastic gradient:

$$\nabla_{\theta_d} \frac{1}{m} \sum_{i=1}^{m} \left[ \log D\left(x^{(i)}\right) + \log\left(1 - D\left(G\left(z^{(i)}\right)\right)\right) \right].$$
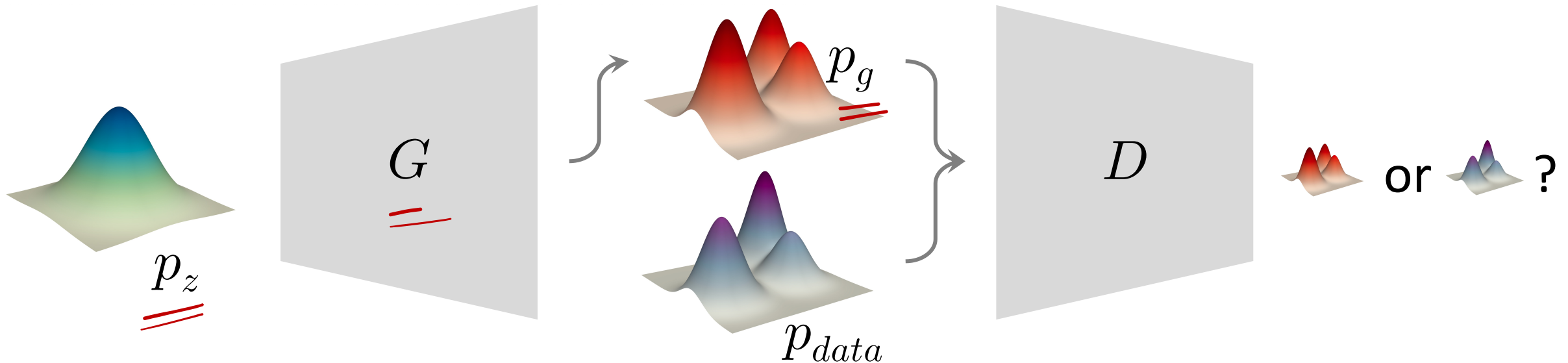
  **end for**

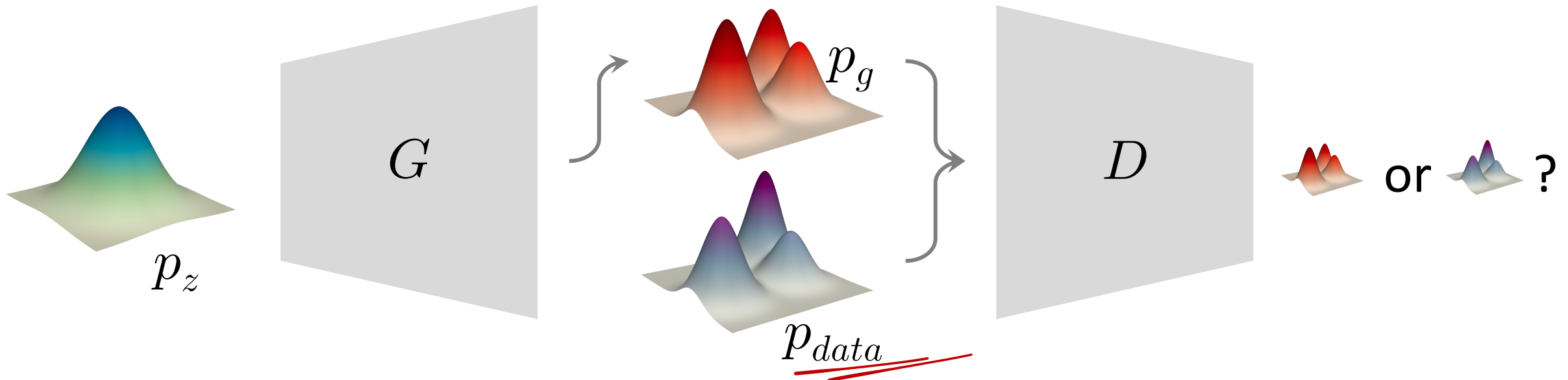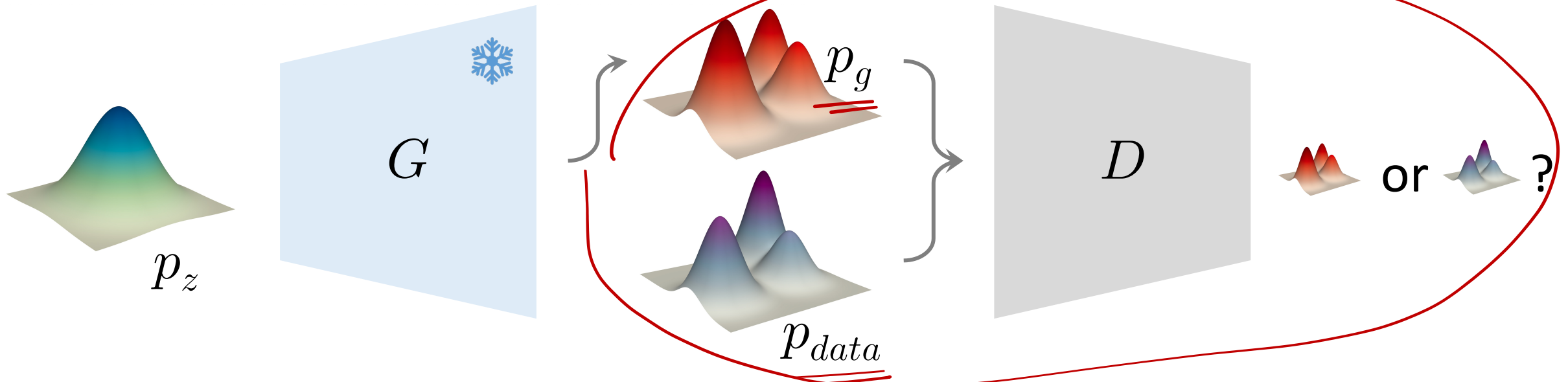  • Sample minibatch of $m$ noise samples $\{z^{(1)}, \ldots, z^{(m)}\}$ from noise prior $p_g(z)$.

  • Update the generator by descending its stochastic gradient:

$$\nabla_{\theta_g} \frac{1}{m} \sum_{i=1}^{m} \log\left(1 - D\left(G\left(z^{(i)}\right)\right)\right).$$

**end for**

The gradient-based updates can use any standard gradient-based learning rule. We used momentum in our experiments.

---

# GAN algorithm annotated

**Algorithm 1** Minibatch stochastic gradient descent training of generative adversarial nets. The number of steps to apply to the discriminator, $k$, is a hyperparameter. We used $k = 1$, the least expensive option, in our experiments.

**for** number of training iterations **do**

   **for** $k$ steps **do**

     • Sample minibatch of $m$ noise samples $\{\boldsymbol{z}^{(1)}, \ldots, \boldsymbol{z}^{(m)}\}$ from noise prior $p_g(\boldsymbol{z})$.

     • Sample minibatch of $m$ examples $\{\boldsymbol{x}^{(1)}, \ldots, \boldsymbol{x}^{(m)}\}$ from data generating distribution $p_{\text{data}}(\boldsymbol{x})$.

     • Update the discriminator by ascending its stochastic gradient:

**D-step**
$$\nabla_{\theta_d} \frac{1}{m} \sum_{i=1}^{m} \left[ \log D\left(\boldsymbol{x}^{(i)}\right) + \log\left(1 - D\left(G\left(\boldsymbol{z}^{(i)}\right)\right)\right) \right].$$

**end for**

   • Sample minibatch of $m$ noise samples $\{\boldsymbol{z}^{(1)}, \ldots, \boldsymbol{z}^{(m)}\}$ from noise prior $p_g(\boldsymbol{z})$.

   • Update the generator by descending its stochastic gradient:

$$\nabla_{\theta_g} \frac{1}{m} \sum_{i=1}^{m} \log\left(1 - D\left(G\left(\boldsymbol{z}^{(i)}\right)\right)\right).$$

**end for**

The gradient-based updates can use any standard gradient-based learning rule. We used momentum in our experiments.

gradient **ascend**
(maximize)

# GAN algorithm annotated

**for** number of training iterations **do**

    **for** $k$ steps **do**

- Sample minibatch of $m$ noise samples $\{z^{(1)}, \dots, z^{(m)}\}$ from noise prior $p_g(z)$.
- Sample minibatch of $m$ examples $\{x^{(1)}, \dots, x^{(m)}\}$ from data generating distribution $p_{\text{data}}(x)$.
- Update the discriminator by ascending its stochastic gradient:

$$\nabla_{\theta_d} \frac{1}{m} \sum_{i=1}^{m} \left[ \log D\left(x^{(i)}\right) + \log\left(1 - D\left(G\left(z^{(i)}\right)\right)\right) \right].$$
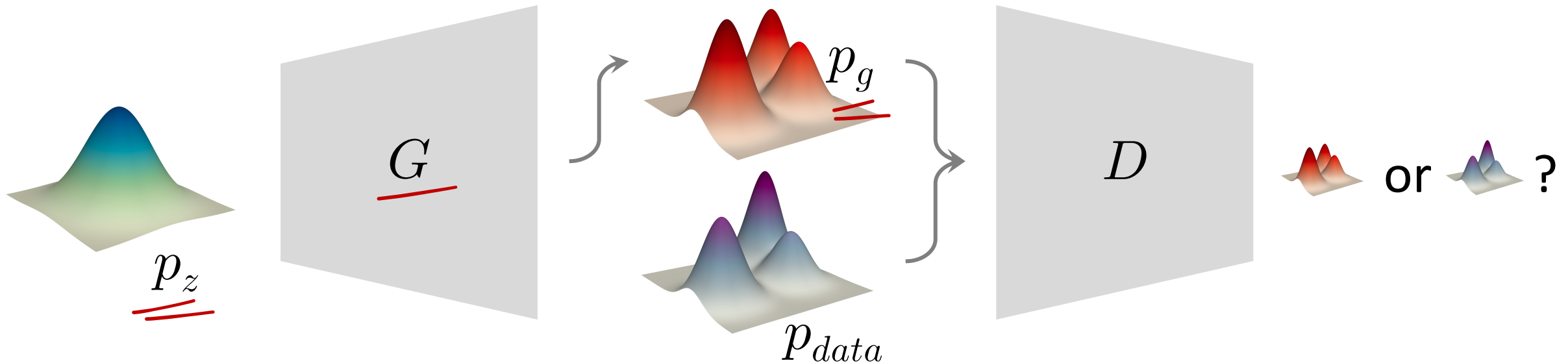
    **end for**

- Sample minibatch of $m$ noise samples $\{z^{(1)}, \dots, z^{(m)}\}$ from noise prior $p_g(z)$.
- Update the generator by descending its stochastic gradient:

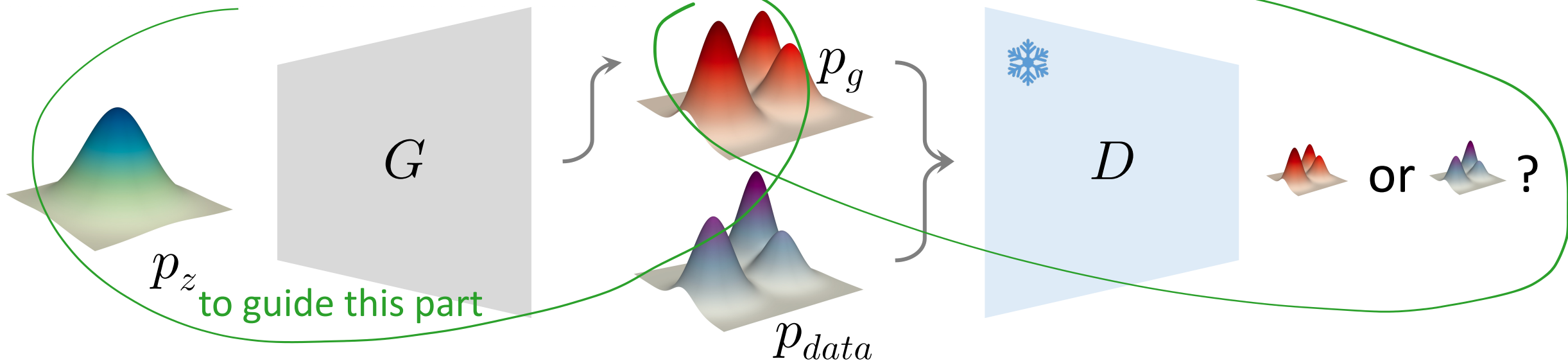$$\nabla_{\theta_g} \frac{1}{m} \sum_{i=1}^{m} \log\left(1 - D\left(G\left(z^{(i)}\right)\right)\right).$$

**end for**

The gradient-based updates can use any standard gradient-based learning rule. We used momentum in our experiments.

# GAN algorithm annotated

**for** number of training iterations **do**

    **for** $k$ steps **do**

- Sample minibatch of $m$ noise samples $\{\boldsymbol{z}^{(1)}, \ldots, \boldsymbol{z}^{(m)}\}$ from noise prior $p_g(\boldsymbol{z})$.
- Sample minibatch of $m$ examples $\{\boldsymbol{x}^{(1)}, \ldots, \boldsymbol{x}^{(m)}\}$ from data generating distribution $p_{\text{data}}(\boldsymbol{x})$.
- Update the discriminator by ascending its stochastic gradient:

$$\nabla_{\theta_d} \frac{1}{m} \sum_{i=1}^{m} \left[ \log D\left(\boldsymbol{x}^{(i)}\right) + \log\left(1 - D\left(G\left(\boldsymbol{z}^{(i)}\right)\right)\right) \right].$$

    **end for**

- Sample minibatch of $m$ noise samples $\{\boldsymbol{z}^{(1)}, \ldots, \boldsymbol{z}^{(m)}\}$ from noise prior $p_g(\boldsymbol{z})$.
- Update the generator by descending its stochastic gradient:

$$\nabla_{\theta_g} \frac{1}{m} \sum_{i=1}^{m} \log\left(1 - D\left(G\left(\boldsymbol{z}^{(i)}\right)\right)\right).$$

**end for**

The gradient-based updates can use any standard gradient-based learning rule. We used momentum in our experiments.

gradient **descend**
(minimize)

G-step

a parameterized
loss function



$p_z$

to guide this part

$G$

$p_g$

$p_{data}$

$D$

or ?

# GAN algorithm annotated

**Algorithm 1** Minibatch stochastic gradient descent training of generative adversarial nets. The number of steps to apply to the discriminator, $k$, is a hyperparameter. We used $k = 1$, the least expensive option, in our experiments.

---

**for** number of training iterations **do**

  **for** $k$ steps **do**

    • Sample minibatch of $m$ noise samples $\{z^{(1)}, \dots, z^{(m)}\}$ from noise prior $p_g(z)$.

    • Sample minibatch of $m$ examples $\{x^{(1)}, \dots, x^{(m)}\}$ from data generating distribution $p_{\text{data}}(x)$.

    • Update the discriminator by ascending its stochastic gradient:

$$\nabla_{\theta_d} \frac{1}{m} \sum_{i=1}^{m} \left[ \log D\left(x^{(i)}\right) + \log\left(1 - D\left(G\left(z^{(i)}\right)\right)\right) \right].$$

  **end for**

  • Sample minibatch of $m$ noise samples $\{z^{(1)}, \dots, z^{(m)}\}$ from noise prior $p_g(z)$.

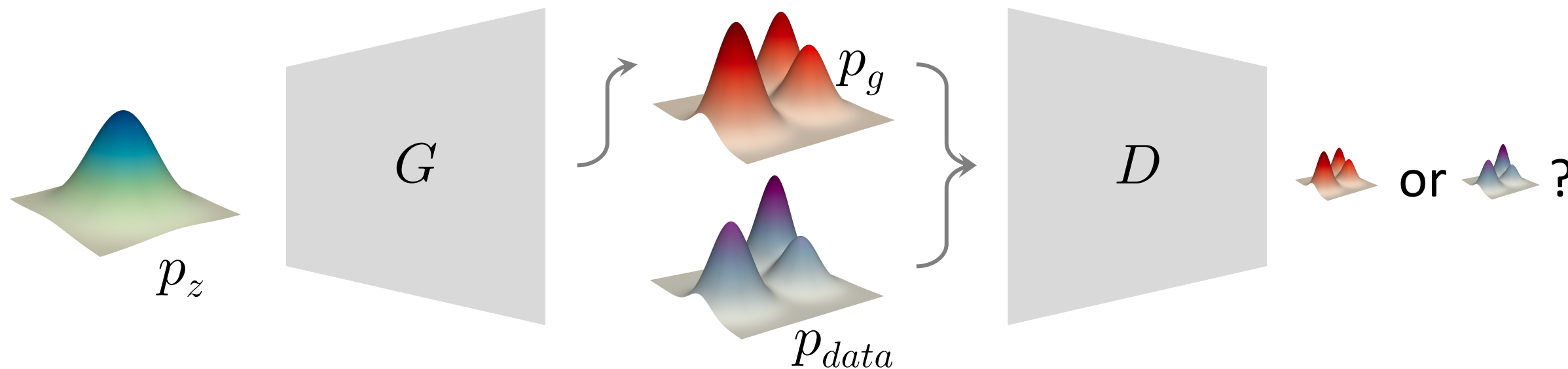  • Update the generator by descending its stochastic gradient:

$$\nabla_{\theta_g} \frac{1}{m} \sum_{i=1}^{m} \log\left(1 - D\left(G\left(z^{(i)}\right)\right)\right).$$

**end for**

The gradient-based updates can use any standard gradient-based learning rule. We used momentum in our experiments.

---

iterating
min-max



$p_z$    $G$    $p_g$    $p_{data}$    $D$    or    ?
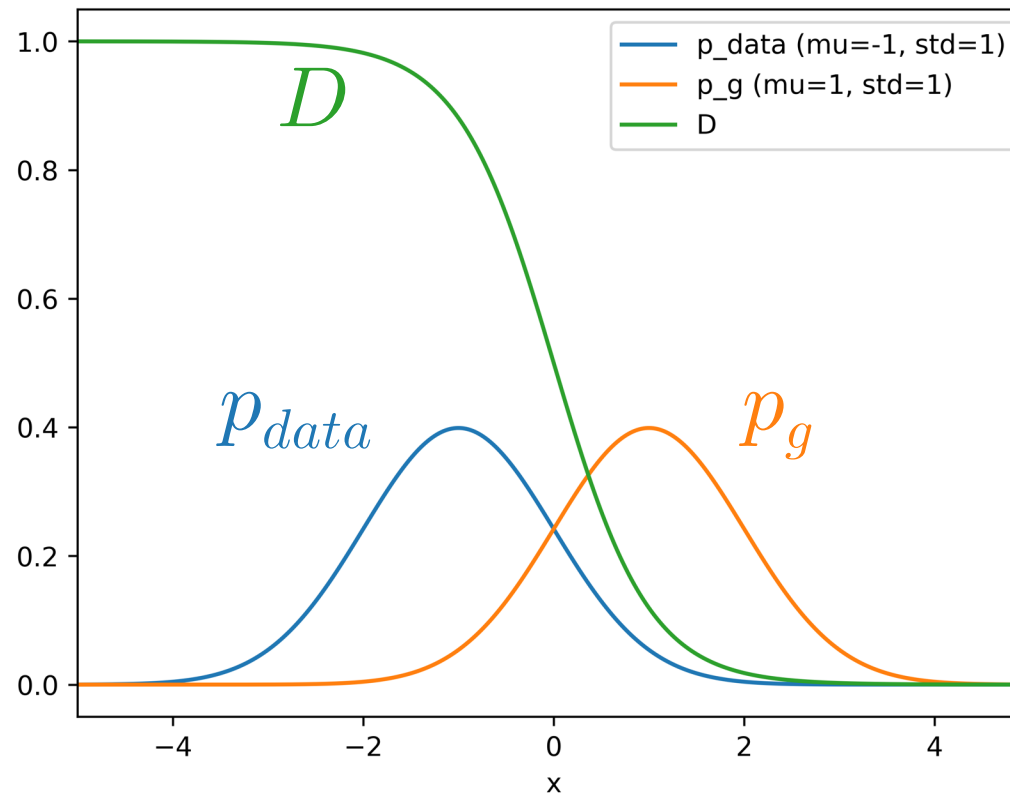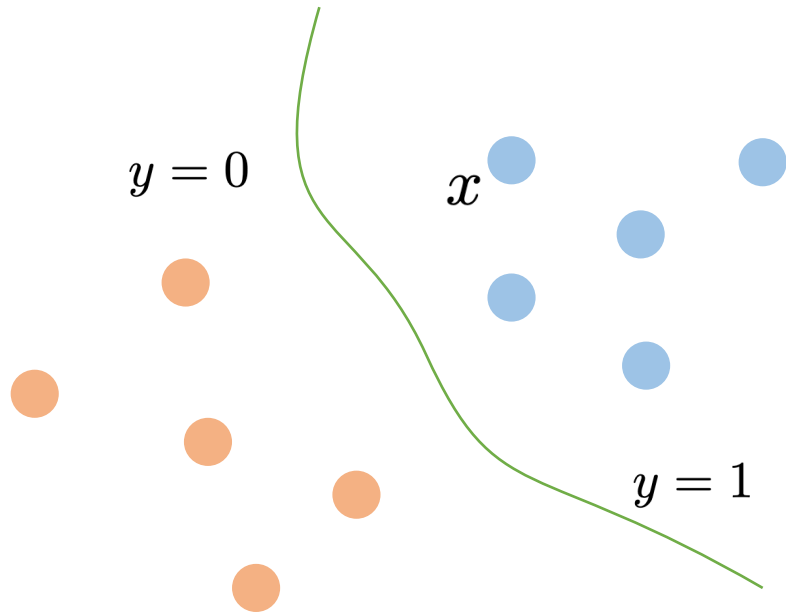
# Theoretical Results

1. For any given $G$, the optimal $D$ is:

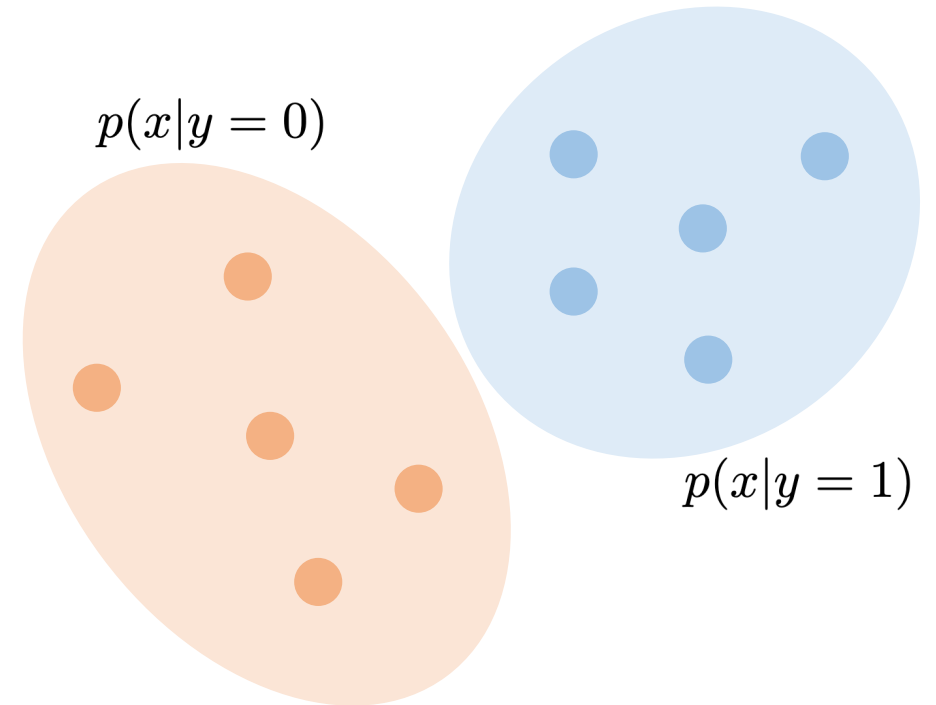$$D_G^*(x) = \frac{p_{\text{data}}(x)}{p_{\text{data}}(x) + p_g(x)}$$

# Recap (Lec. 1): Discriminative vs. Generative

**discriminative**

$y = 0$

$x$

$y = 1$

**generative**

$p(x|y=0)$

$p(x|y=1)$

# Theoretical Results

2. With the optimal $D_G$, the objective function is:

$$\mathcal{L}(D^*, G) = 2D_{JS}(p_{\text{data}} \| p_g) - 2\log 2$$

where $D_{JS}$ is Jensen–Shannon divergence

# Background: Jensen–Shannon divergence

$D_{JS}$: "total divergence to the average"

$$D_{JS}(p\|q) \triangleq \frac{1}{2}D_{KL}(p\|\frac{p+q}{2}) + \frac{1}{2}D_{KL}(q\|\frac{p+q}{2})$$

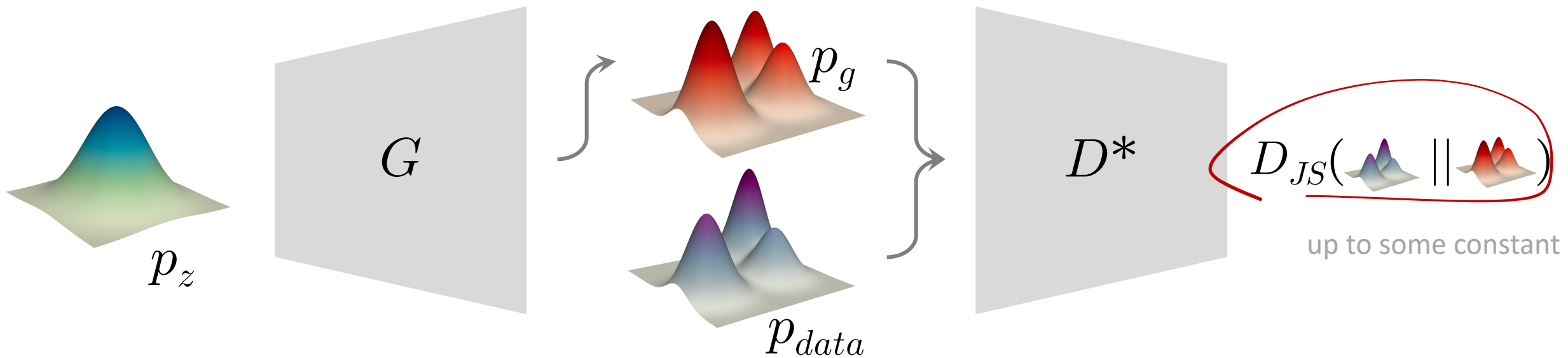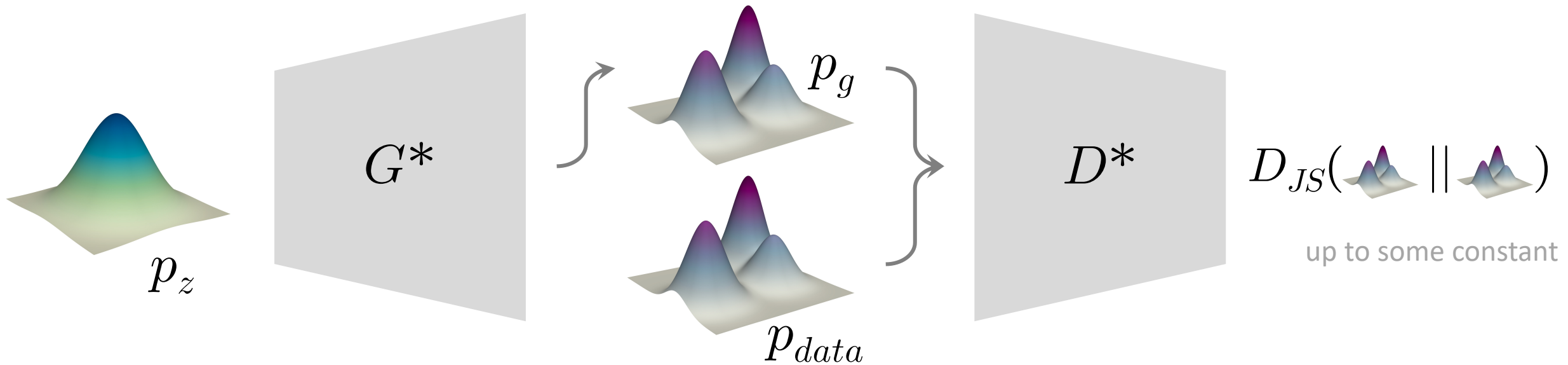# Background: Jensen–Shannon divergence

$D_{JS}$: "total divergence to the average"

$$D_{JS}(p\|q) \triangleq \frac{1}{2}D_{KL}(p\|\frac{p+q}{2}) + \frac{1}{2}D_{KL}(q\|\frac{p+q}{2})$$

Properties:

- $D_{JS}$ is symmetric; $D_{KL}$ is not

- $D_{JS}$ is bounded: $[0, \log2]$; $D_{KL}$ is unbounded: $[0, \inf)$

- $D_{JS}$ is more stable

# Theoretical Results

2. With the optimal $D_G$, the objective function is:

$$\mathcal{L}(D^*, G) = 2D_{JS}(p_{\text{data}} \| p_g) - 2\log 2$$

**GAN optimizes for Jensen–Shannon divergence.**



up to some constant

See proof in "Generative Adversarial Nets", Goodfellow, et al., 2014

# Theoretical Results

3. Global optimality is achieved at $p_g = p_{data}$

$$\mathcal{L}(D^*, \overset{*}{G}) = \cancel{2D_{JS}(p_{\text{data}}\|p_g)} - 2\log 2$$

$$= 0$$



$p_z$     $G*$     $p_g$   $p_{data}$     $D*$     $D_{JS}(\;\|\;)$

up to some constant

See proof in "Generative Adversarial Nets", Goodfellow, et al., 2014

# Theoretical Results: Summary

1. For any given $G$, the optimal $D$ is:

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

2. With optimal $D_G$, GAN optimizes for Jensen–Shannon divergence:

$$\mathcal{L}(D^*, G) = 2D_{JS}(p_{\text{data}}\|p_g) - 2\log 2$$

3. Global optimality is achieved at $p_g = p_{data}$
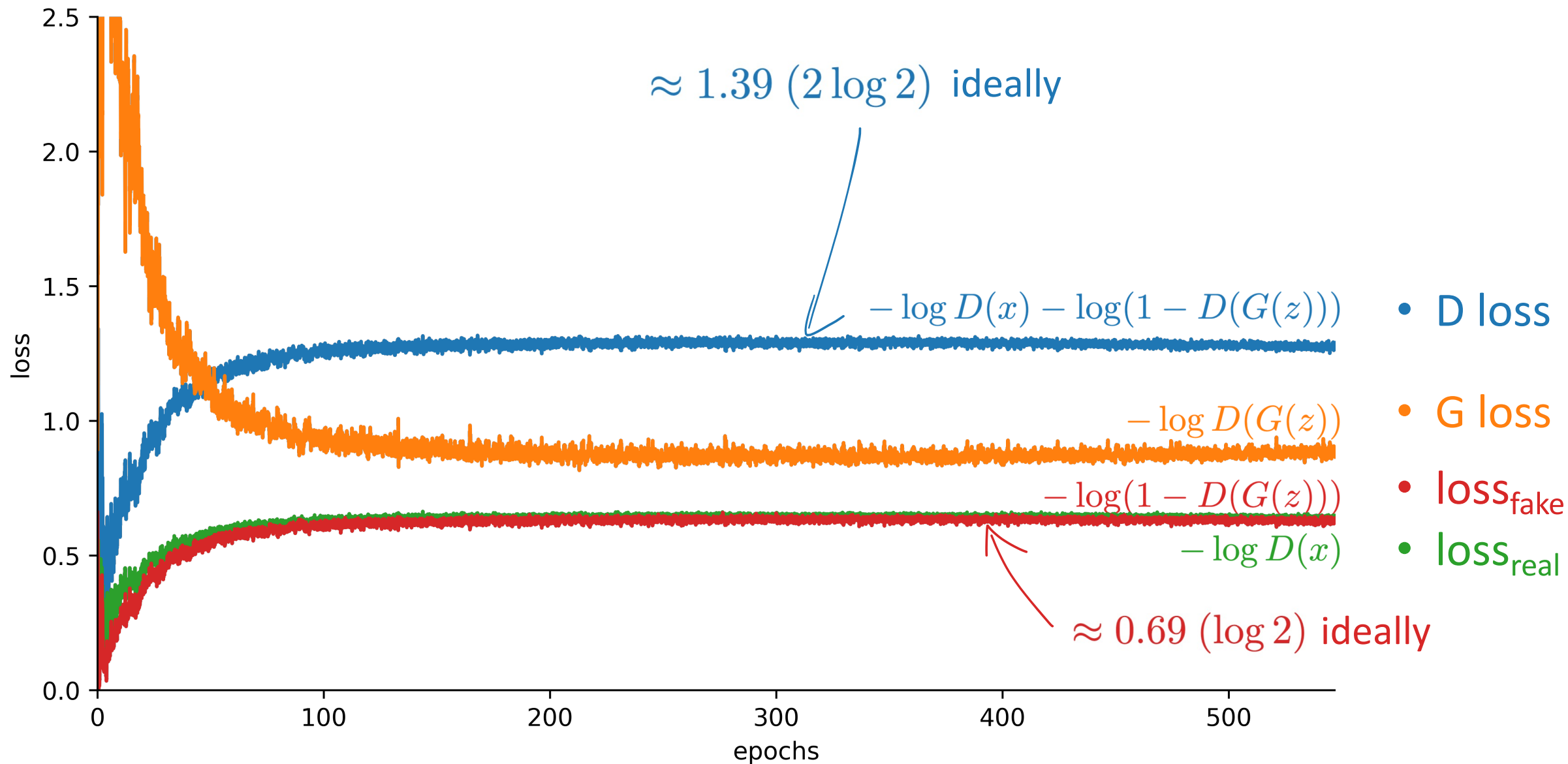
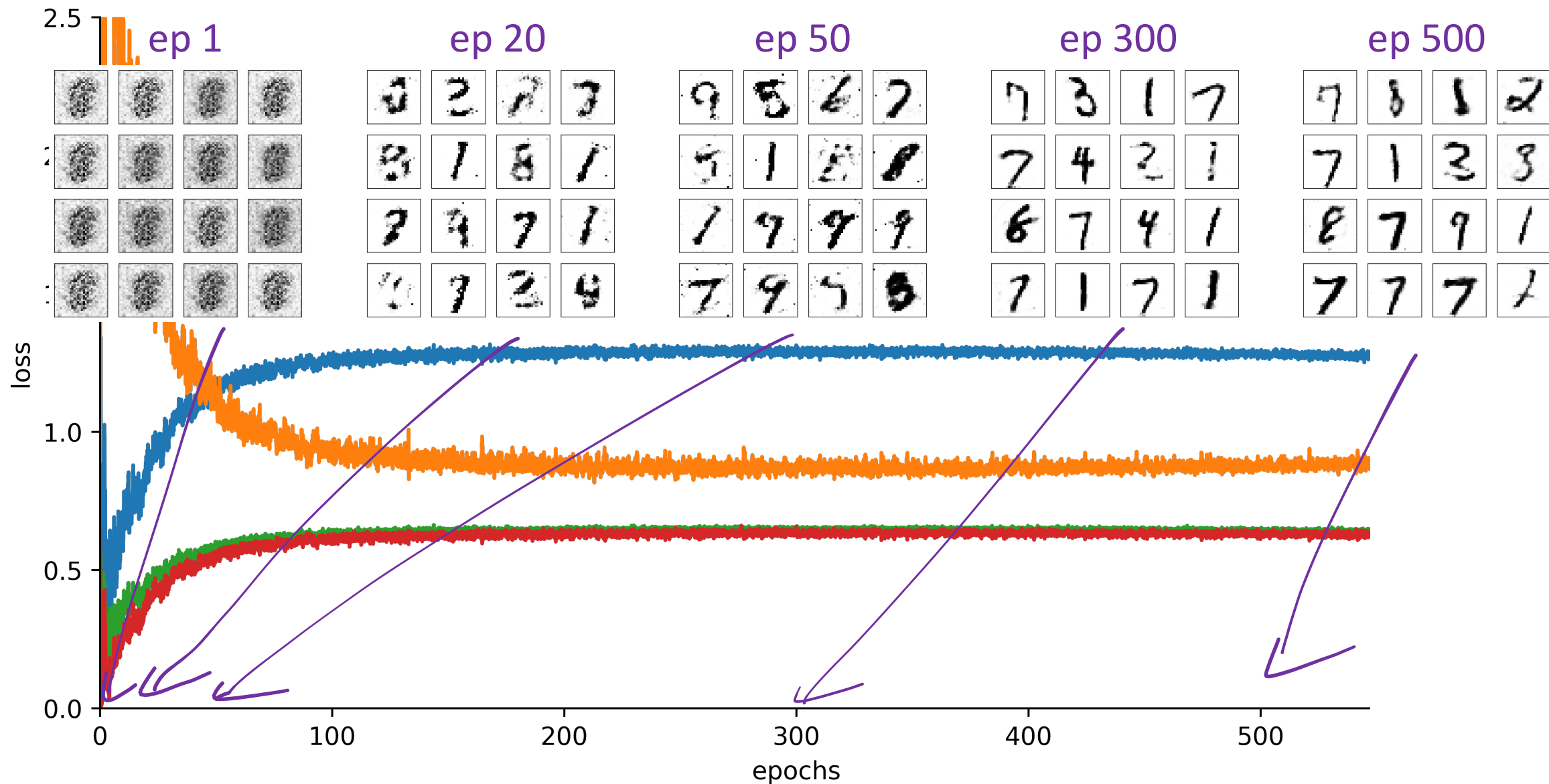$$\mathcal{L}(D^*, G^*) = -2\log 2$$

# Running example: MNIST



Legend:
- $-\log D(x) - \log(1 - D(G(z)))$ — D loss
- $-\log D(G(z))$ — G loss
- $-\log(1 - D(G(z)))$ — loss$_{\text{fake}}$
- $-\log D(x)$ — loss$_{\text{real}}$

Axes: loss vs epochs

Code adapted from: https://github.com/prcastro/pytorch-gan/tree/master

# Running example: MNIST



$-\log D(x) - \log(1 - D(G(z)))$   • D loss

$-\log D(G(z))$   • G loss

$-\log(1 - D(G(z)))$   • loss$_{fake}$

$-\log D(x)$   • loss$_{real}$

max w/ G   min w/ G

min w/ D

min w/ D

*All objectives are negative of their original form

Code adapted from: https://github.com/prcastro/pytorch-gan/tree/master

# Running example: MNIST

# Running example: MNIST



ep 1      ep 20      ep 50      ep 300      ep 500

loss

epochs

# Problems of GAN

Difficult to train/converge

- Hard to achieve equilibrium

- Vanishing gradients

- Mode collapse



oscillating

J. Brownlee, "How to Identify and Diagnose GAN Failure Modes"



vanishing grad

Arjovsky & Bottou, "Towards Principled Methods for Training GANs"



mode collapse

Step 0    Step 5k    Step 10k    Step 15k    Step 20k    Step 25k    Target

L. Metz, "Unrolled Generative Adversarial Networks"

# Running example: GAN Lab

# Wasserstein GAN

# W-GAN in Short

For mathematicians:

- Wasserstein distance, instead of JS divergence

For engineers:

- remove logarithms
- clip weights

For laymen:

- art critic, instead of forgery expert

# Recap: GAN optimizes for $D_{JS}$

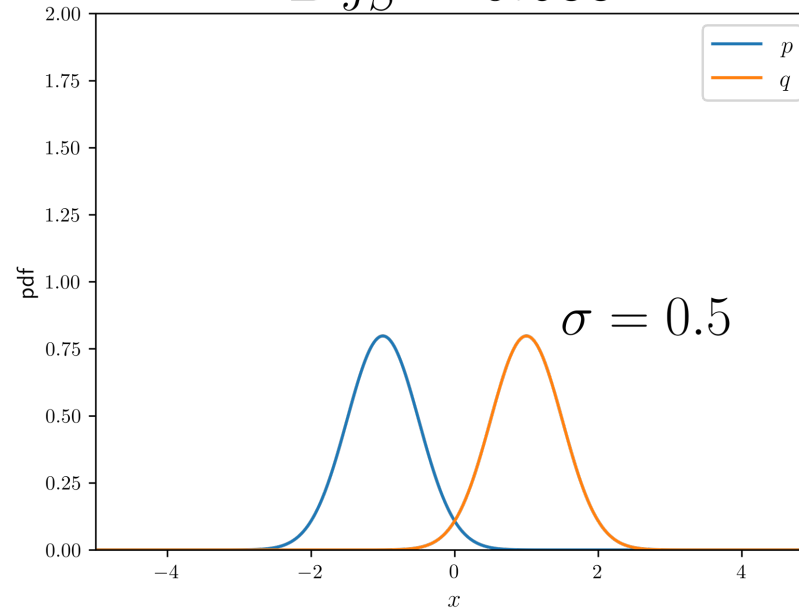$$\mathcal{L}(D^*, G) = 2D_{JS}(p_{\text{data}} \| p_g) - 2\log 2$$

# Problems of $D_{JS}$

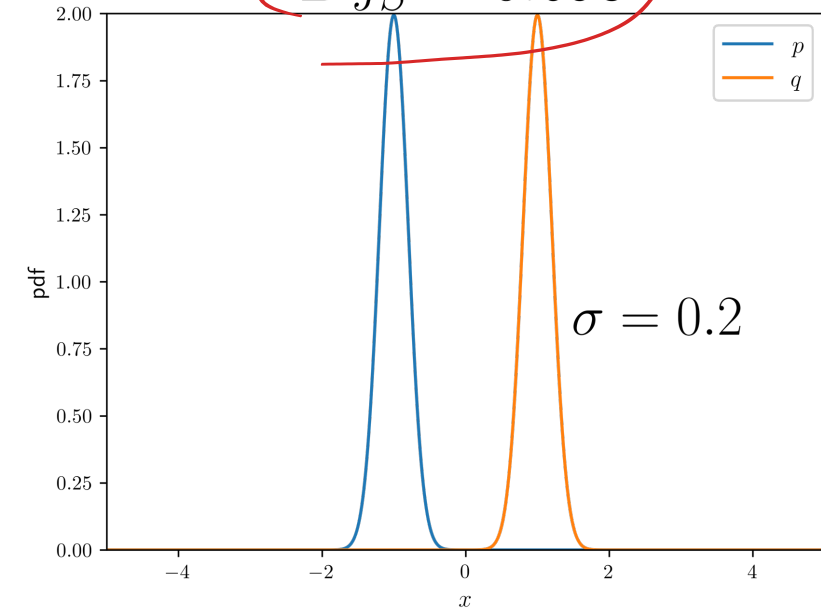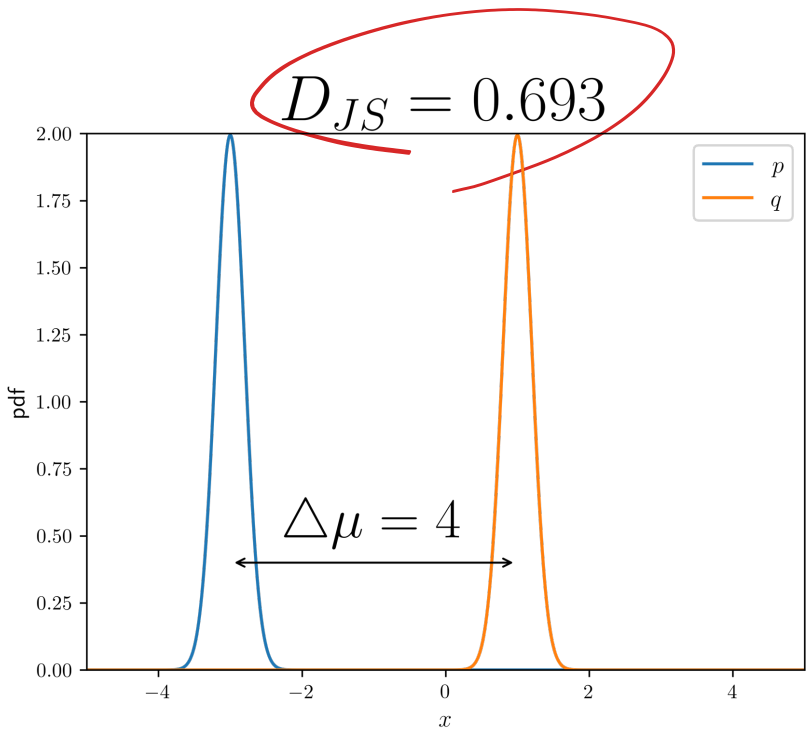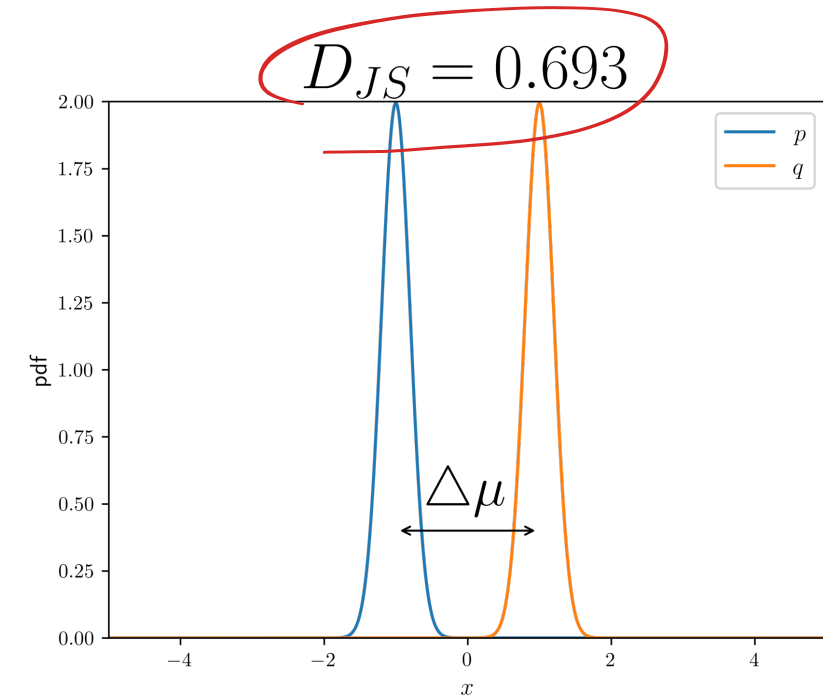If $p$ and $q$ don't overlap, $D_{JS}$ is a constant (log2), i.e., no gradient
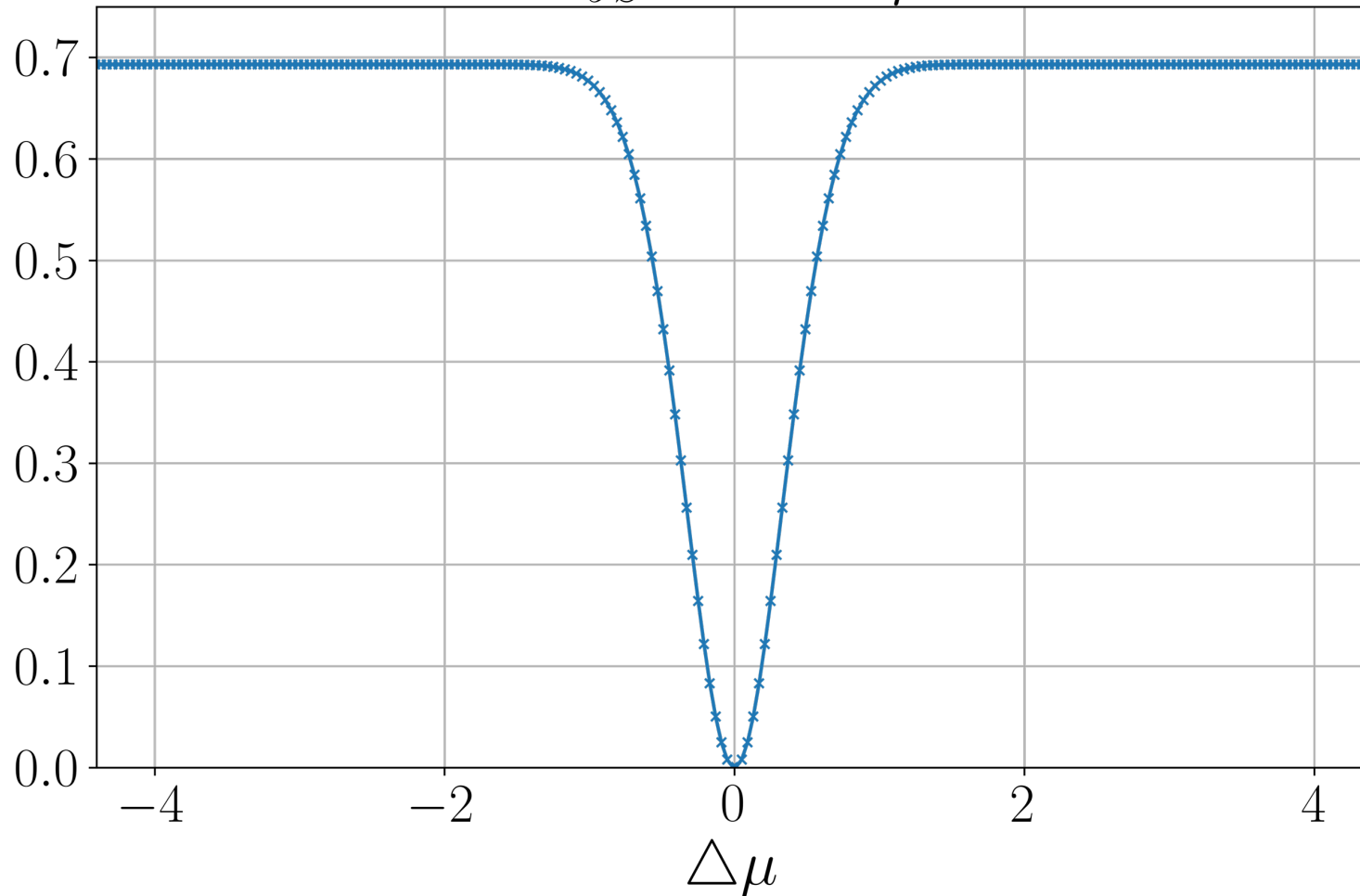
# Problems of $D_{JS}$

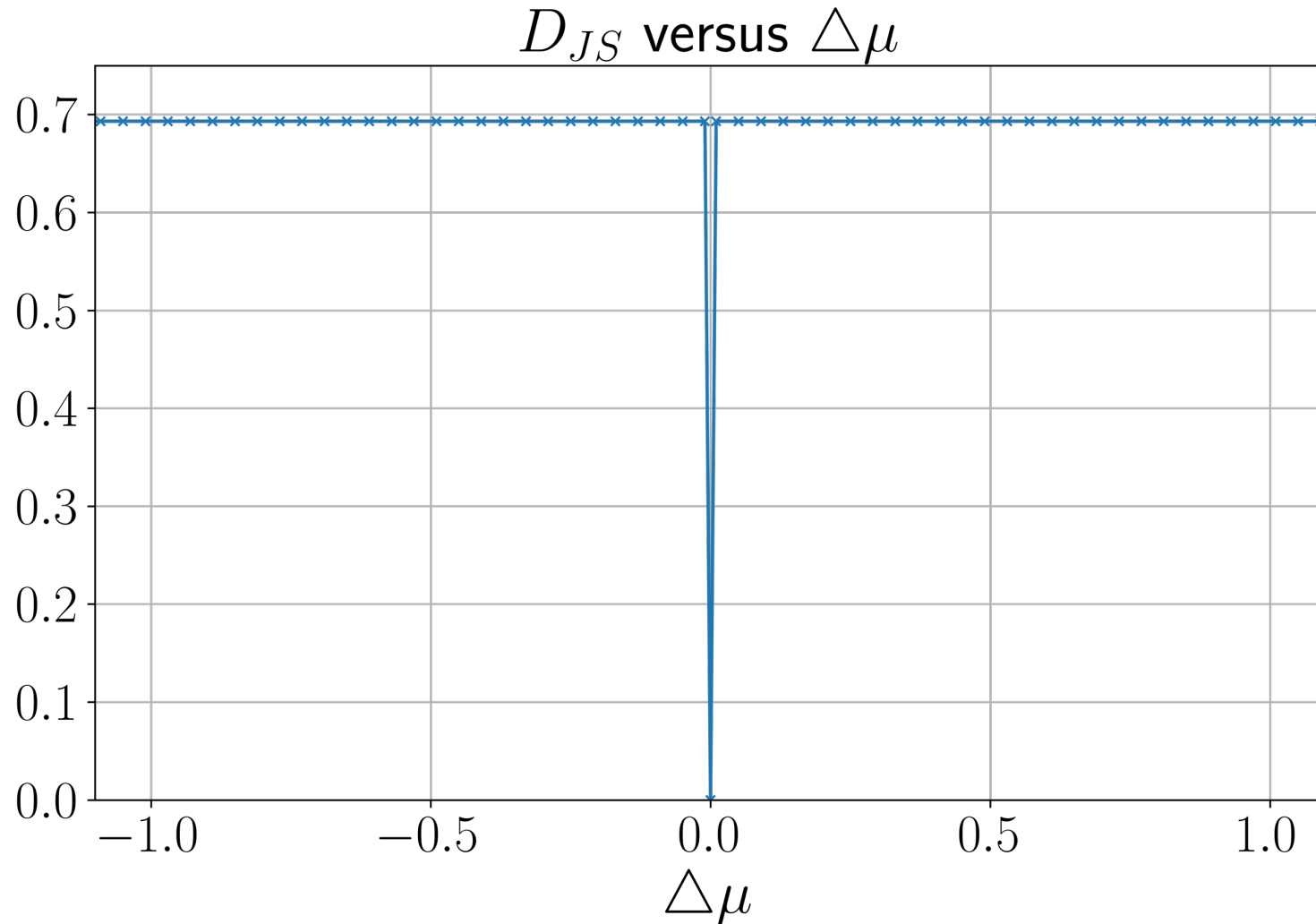If $p$ and $q$ don't overlap, $D_{JS}$ is a constant (log2), i.e., no gradient

# Problems of $D_{JS}$

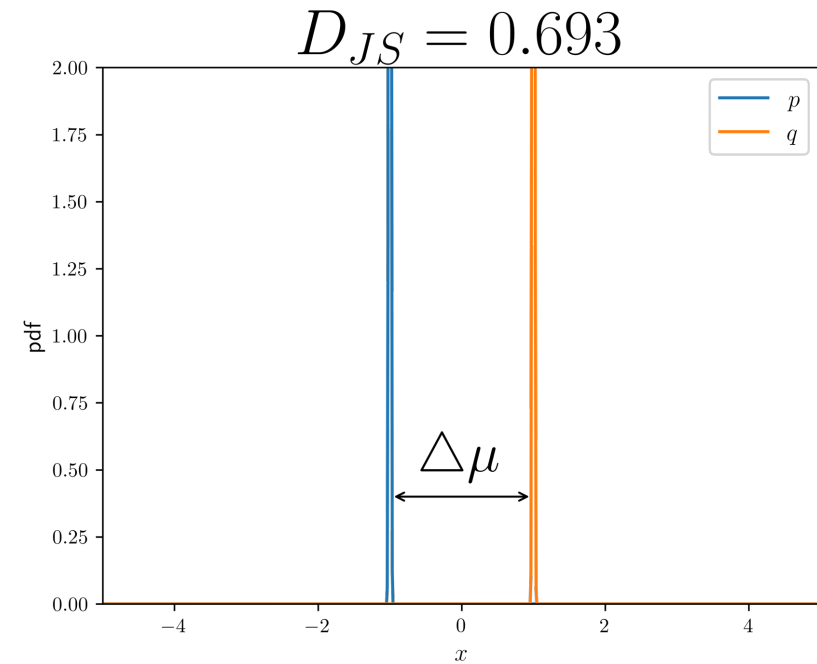- $D_{JS}$ is useful only if $p$ and $q$ are close

# Problems of $D_{JS}$
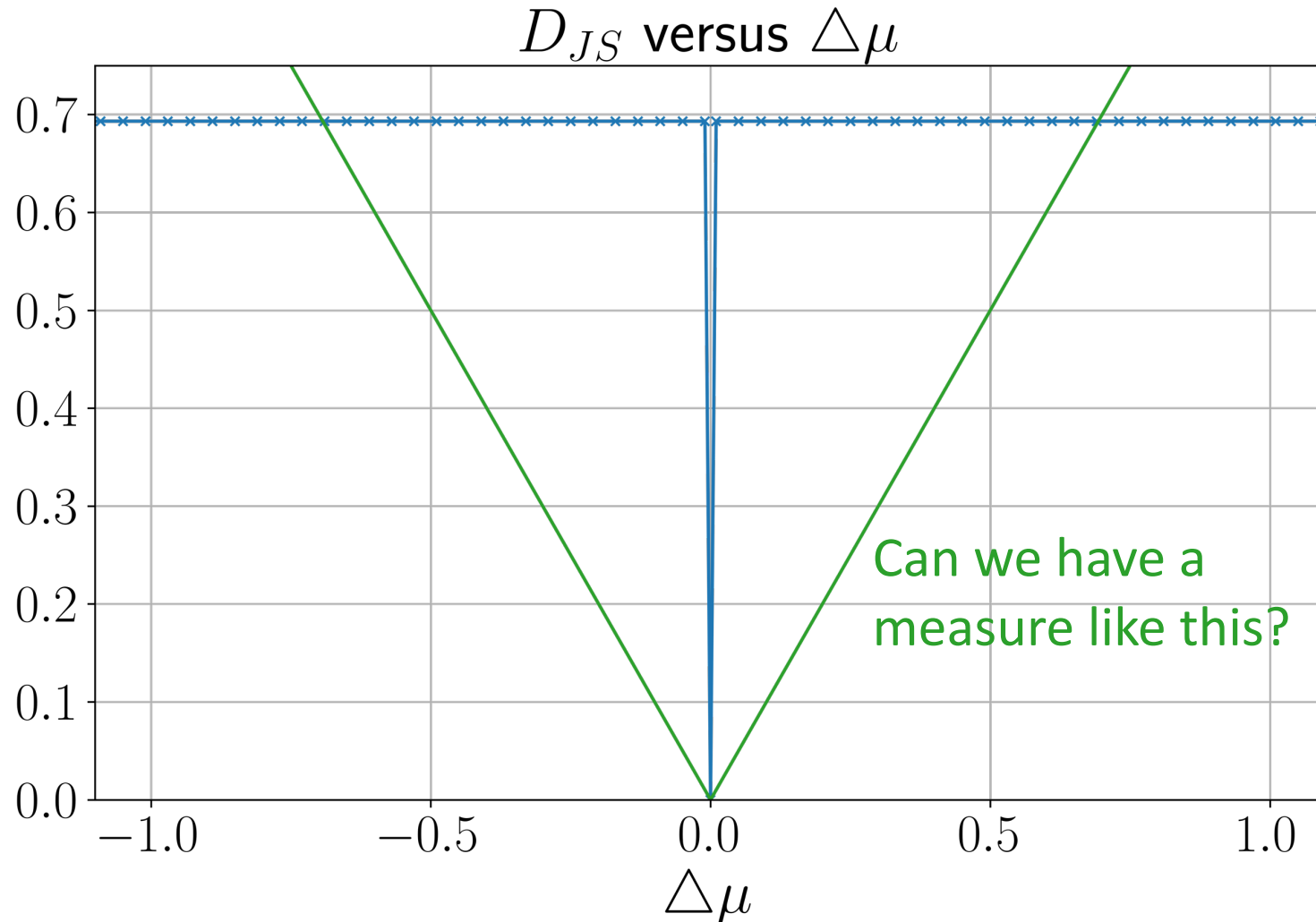
- $D_{JS}$ is a delta function when $p$ and $q$ are delta functions

# Problems of $D_{JS}$

- $D_{JS}$ is a delta function when $p$ and $q$ are delta functions



$D_{JS}$ versus $\triangle\mu$

Can we have a measure like this?

$D_{JS} = 0.693$

$\triangle\mu$

# Wasserstein Distance

"Earth Mover's Distance"

# Running example: Wasserstein Distance

# Running example: Wasserstein Distance

# Running example: Wasserstein Distance



Figure inspired by: Lilian Weng, "From GAN to WGAN", arXiv:1904.08994

# Running example: Wasserstein Distance

2 shovelfuls

# Running example: Wasserstein Distance



1 shovelful

$P_1$  $P_2$  $P_3$  $P_4$

$Q_1$  $Q_2$  $Q_3$  $Q_4$

# Running example: Wasserstein Distance

$$W(P, Q) = 5 \times \blacksquare$$

# Running example: Wasserstein Distance

# Running example: Wasserstein Distance



$P_1 \quad P_2 \quad P_3 \quad P_4$

$\mathbf{cdf}_P$

$\mathbf{cdf}_Q$

$Q_1 \quad Q_2 \quad Q_3 \quad Q_4$

- cdf: cumulative distribution function

# Running example: Wasserstein Distance



$$W(P, Q) = 5 \times \blacksquare$$

# Wasserstein Distance

- 1-Wasserstein Distance (1-d, discrete)

$l_1$-norm

$$W_1(P, Q) = \sum_i |\mathbf{cdf}_P(i) - \mathbf{cdf}_Q(i)|$$

- 1-Wasserstein Distance (1-d, continuous)

$$W_1(p, q) = \int_x |\mathbf{cdf}_p(x) - \mathbf{cdf}_q(x)| dx$$

# Recap: $D_{JS}$

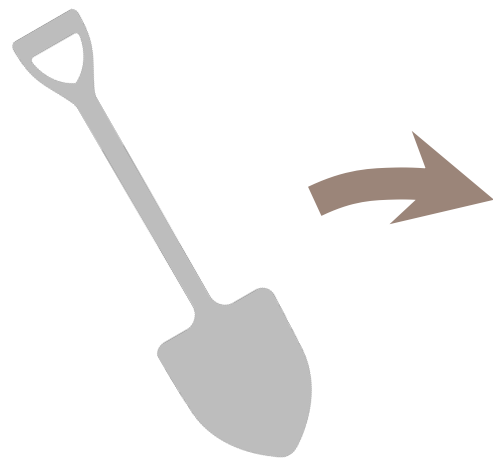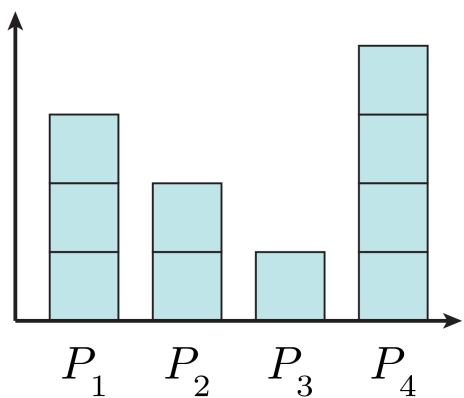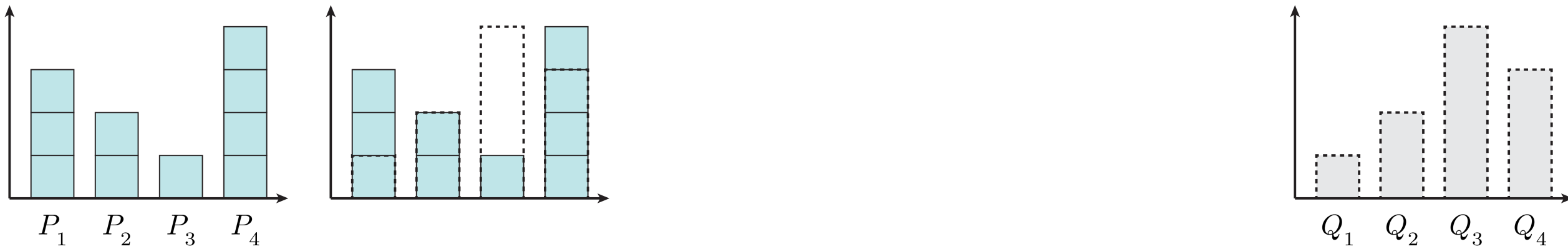- $D_{JS}$ is a delta function when $p$ and $q$ are delta functions



$D_{JS}$ versus $\triangle\mu$

$\triangle\mu$

Can we have a measure like this?

$D_{JS} = 0.693$

# Wasserstein Distance

- when $p$ and $q$ are delta functions:

$$W_1(p, q) = |\mu_p - \mu_q|$$

# Wasserstein Distance

- 1-Wasserstein Distance (high-dim, continuous)

$$W_1(p, q) = \inf_{\gamma \in \Pi(p,q)} \mathbb{E}_{(x,y) \sim \gamma}[\|x - y\|]$$

- all joint distributions $\gamma(x, y)$
  whose marginals are $p$ and $q$

# W-GAN optimizes for Wasserstein Distance

$$W_1(p, q) = \inf_{\gamma \in \Pi(p,q)} \mathbb{E}_{(x,y) \sim \gamma}[\|x - y\|]$$

# W-GAN optimizes for Wasserstein Distance

- Kantorovich-Rubinstein duality:

$$W_1(p, q) = \sup_{\|f\|_L \leq 1} \mathbb{E}_{x \sim p}[f(x)] - \mathbb{E}_{x \sim q}[f(x)]$$

- all 1-Lipschitz functions

# W-GAN optimizes for Wasserstein Distance

- Kantorovich-Rubinstein duality:

$$W_1(p, q) = \sup_{\|f\|_L \leq 1} \mathbb{E}_{x \sim p}[f(x)] - \mathbb{E}_{x \sim q}[f(x)]$$

- all 1-Lipschitz functions

$K$-Lipschitz continuity:

$$|f(x) - f(y)| \leq K|x - y|, \quad \forall x, y$$

gradient is bounded:

$$\frac{|f(x) - f(y)|}{|x - y|} \leq K$$



Figure from: https://en.wikipedia.org/wiki/Lipschitz_continuity

# W-GAN optimizes for Wasserstein Distance

- Kantorovich-Rubinstein duality:

$$W_1(p, q) = \frac{1}{K} \sup_{\|f\|_L \leq K} \mathbb{E}_{x \sim p}[f(x)] - \mathbb{E}_{x \sim q}[f(x)]$$

$K$-Lipschitz continuity:

$$|f(x) - f(y)| \leq K|x - y|, \quad \forall x, y$$

# W-GAN optimizes for Wasserstein Distance

- W-GAN's objective function:

$$\max_{w \in \mathcal{W}} \mathbb{E}_{x \sim p_{\text{data}}}[f_w(x)] - \mathbb{E}_{x \sim p_g}[f_w(x)]$$

# W-GAN optimizes for Wasserstein Distance

- W-GAN's objective function:

$$\max_{w \in \mathcal{W}} \mathbb{E}_{x \sim p_{\text{data}}}[f_w(x)] - \mathbb{E}_{x \sim p_g}[f_w(x)]$$

# W-GAN optimizes for Wasserstein Distance

- W-GAN's objective function:

$$\max_{w \in \mathcal{W}} \mathbb{E}_{x \sim p_{\text{data}}}[f_w(x)] - \mathbb{E}_{x \sim p_g}[f_w(x)]$$

- weights are bounded: in practice, clipped [-0.01, 0.01]



$p_g$

$f_w$

$p_{data}$

$W_1( \; || \; )$

# W-GAN vs. original GAN

- W-GAN's objective function:

$$\max_{w \in \mathcal{W}} \mathbb{E}_{x \sim p_{\text{data}}}[f_w(x)] - \mathbb{E}_{x \sim p_g}[f_w(x)]$$

- clip weights

- remove logarithms

- original GAN's objective function (D-step):

$$\max_{D} \mathbb{E}_{x \sim p_{\text{data}}}[\log D(x)] + \mathbb{E}_{x \sim p_g}[\log(1 - D(x))]$$

# W-GAN vs. original GAN

"art critic"
- value/merit/quality/...
- direction to improve (gradients)

- W-GAN's objective function:

$$\max_{w \in \mathcal{W}} \mathbb{E}_{x \sim p_{\text{data}}}[f_w(x)] - \mathbb{E}_{x \sim p_g}[f_w(x)]$$

- original GAN's objective function (D-step):

$$\max_{D} \mathbb{E}_{x \sim p_{\text{data}}}[\log D(x)] + \mathbb{E}_{x \sim p_g}[\log(1 - D(x))]$$

"forgery expert"
- real/fake

# W-GAN algorithm annotated

**Algorithm 1** WGAN, our proposed algorithm. All experiments in the paper used the default values $\alpha = 0.00005$, $c = 0.01$, $m = 64$, $n_{\text{critic}} = 5$.

**Require:** : $\alpha$, the learning rate. $c$, the clipping parameter. $m$, the batch size.
$\quad$ $n_{\text{critic}}$, the number of iterations of the critic per generator iteration.
**Require:** : $w_0$, initial critic parameters. $\theta_0$, initial generator's parameters.

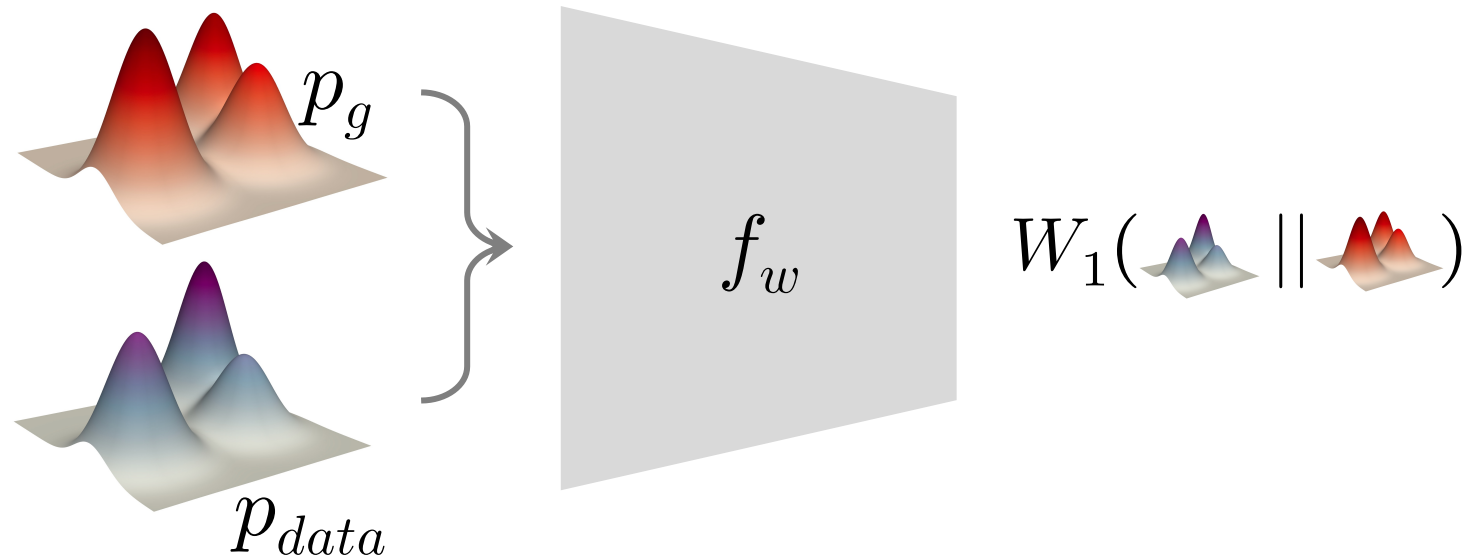1: **while** $\theta$ has not converged **do**
2: $\quad$ **for** $t = 0, ..., n_{\text{critic}}$ **do**
3: $\qquad$ Sample $\{x^{(i)}\}_{i=1}^{m} \sim \mathbb{P}_r$ a batch from the real data.
4: $\qquad$ Sample $\{z^{(i)}\}_{i=1}^{m} \sim p(z)$ a batch of prior samples.
5: $\qquad$ $g_w \leftarrow \nabla_w \left[ \frac{1}{m} \sum_{i=1}^{m} f_w(x^{(i)}) - \frac{1}{m} \sum_{i=1}^{m} f_w(g_\theta(z^{(i)})) \right]$ <span style="color:red">remove logarithms</span>
6: $\qquad$ $w \leftarrow w + \alpha \cdot \text{RMSProp}(w, g_w)$
7: $\qquad$ $w \leftarrow \text{clip}(w, -c, c)$
8: $\quad$ **end for** <span style="color:purple">clip weights</span>
9: $\quad$ Sample $\{z^{(i)}\}_{i=1}^{m} \sim p(z)$ a batch of prior samples.
10: $\quad$ $g_\theta \leftarrow -\nabla_\theta \frac{1}{m} \sum_{i=1}^{m} f_w(g_\theta(z^{(i)}))$
11: $\quad$ $\theta \leftarrow \theta - \alpha \cdot \text{RMSProp}(\theta, g_\theta)$
12: **end while**



$p_z \qquad g \qquad p_g \qquad p_{data} \qquad f_w \qquad W_1(\ ||\ )$

# W-GAN vs. original GAN

# W-GAN vs. original GAN



original GAN

W-GAN

KDE

Samples

KDE

Samples

**Epoch 0**  **Epoch 1**  **Epoch 5**  **Epoch 10**  **Epoch 20**  **Epoch 50**  **Epoch 100**

# W-GAN in Short

For mathematicians:

- Wasserstein distance, instead of JS divergence

For engineers:

- remove logarithms    Wasserstein distance
- clip weights
          Lipschitz continuity

For laymen:

- art critic instead of a forgery expert
          gradients

# Brief: LSGAN, EBGAN

- Least Square (LS) GAN:

$$\mathbb{E}_{x \sim p_{\text{data}}}(D(x) - b)^2 \quad + \quad \mathbb{E}_{x \sim p_g}(D(x) - a)^2$$

- Energy-based (EB) GAN:

$$\mathbb{E}_{x \sim p_{\text{data}}} D(x) \quad + \quad \mathbb{E}_{x \sim p_g}[m - D(x)]^+$$

Mao, et al., "Least Squares Generative Adversarial Networks", ICCV 2017
Zhao, et al., "Energy-based Generative Adversarial Networks", ICLR 2017

# Adversary as a Loss Function

# Adversary as a Loss Function

- GAN essentially defines an **adversarial loss** function

- Input to networks is **not** necessarily random/noise

- **Beyond L2/L1**: adversarial loss encourages output to look "realistic"

- **Combined with L2/L1**: reconstruction loss largely stabilizes training

# Adversary as a Loss Function

- GAN: input is random

# Adversary as a Loss Function

- Input can be from another source

# Adversary as a Loss Function

- Input can be from another source



reconstruction loss

- parameterized loss function
- trained alternately

adversarial loss

# Example: Super-Resolution GAN

- better PSNR

- worse PSNR, but better visual quality

original

bicubic
(21.59dB/0.6423)

SRResNet
(23.44dB/0.7777)

SRGAN
(20.34dB/0.6562)



Ledig, et al., "Photo-Realistic Single Image Super-Resolution Using a Generative Adversarial Network", CVPR 2016

# Example: Super-Resolution GAN



original       reconstruction       reconstruction + adversarial

Ledig, et al., "Photo-Realistic Single Image Super-Resolution Using a Generative Adversarial Network", CVPR 2016

# Example: Context Encoder



| Image | Ours($L2$) | Ours(Adv) | Ours(L2+Adv) |

Pathak, et al., "Context Encoders: Feature Learning by Inpainting", CVPR 2016

# Example: pix2pix



| Input | Ground truth | L1 | cGAN | L1 + cGAN |

Isola, et al., "Image-to-Image Translation with Conditional Adversarial Networks", CVPR 2017

# Example: CycleGAN



zebra $\longrightarrow$ horse

horse $\longrightarrow$ zebra

Zhu, et al., "Unpaired Image-to-Image Translation using Cycle-Consistent Adversarial Networks", ICCV 2017

# From VQ-VAE to VQ-GAN



codebook

$e_0 \; e_1 \; \dots$

$x$

**ConvNet** encoder

$z_e$

VQ encoder

$z$

VQ decoder

$z_q$

**ConvNet** decoder

$x$'

**VQ-VAE**: L2-only

$l_2$

$D$

**VQ-GAN**: L2 + Adv

$h \times w \times c$

$h \times w \times k$

$h \times w \times c$

Esser, et al., "Taming Transformers for High-Resolution Image Synthesis", CVPR 2021

# From VQ-VAE to VQ-GAN

**VQ-VAE**

**VQ-GAN**



Esser, et al., "Taming Transformers for High-Resolution Image Synthesis", CVPR 2021

# Discussion

- To be precise: **VQ-GAN** = **VQ-VAE** + Adv Loss + Perceptual Loss

- w/o VQ, it's **VAE** + Adv Loss + Perceptual Loss

- Both are the *de facto* **tokenizers** in image generation
  - w/ VQ: e.g., Autoregressive Models
  - w/o VQ: e.g., Diffusion Models

- Commercial models (e.g., Stable Diffusion, Sora) use these tokenizers

It involves everything!

# This Lecture

- Generative Adversarial Networks (GAN)

- Wasserstein GAN (W-GAN)

- Adversary as a Loss Function

**Main References**

- Goodfellow et al. "Generative Adversarial Nets", NeurIPS 2014
- Arjovsky et al. "Wasserstein Generative Adversarial Networks", ICML 2017