

Generative Adversarial Networks (GAN)

Generative Model

- Generative Adversarial Networks (GANs)
- Variational Autoencoders (VAEs)
 - Previous lectures
- Autoregressive models : PixelRNN and PixelCNN
 - <https://arxiv.org/abs/1601.06759>
 - <https://arxiv.org/abs/1606.05328>
 - <https://openreview.net/pdf?id=BjrFC6ceg>
- Good reference
 - http://cs231n.stanford.edu/slides/2019/cs231n_2019_lecture13.pdf

GAN

I. Goodfellow et al., “Generative Adversarial Nets”
NIPS2014

Generative Adversarial Nets

**Ian J. Goodfellow, Jean Pouget-Abadie,^{*} Mehdi Mirza, Bing Xu, David Warde-Farley,
Sherjil Ozair,[†] Aaron Courville, Yoshua Bengio[‡]**

Département d'informatique et de recherche opérationnelle
Université de Montréal
Montréal, QC H3C 3J7

Abstract

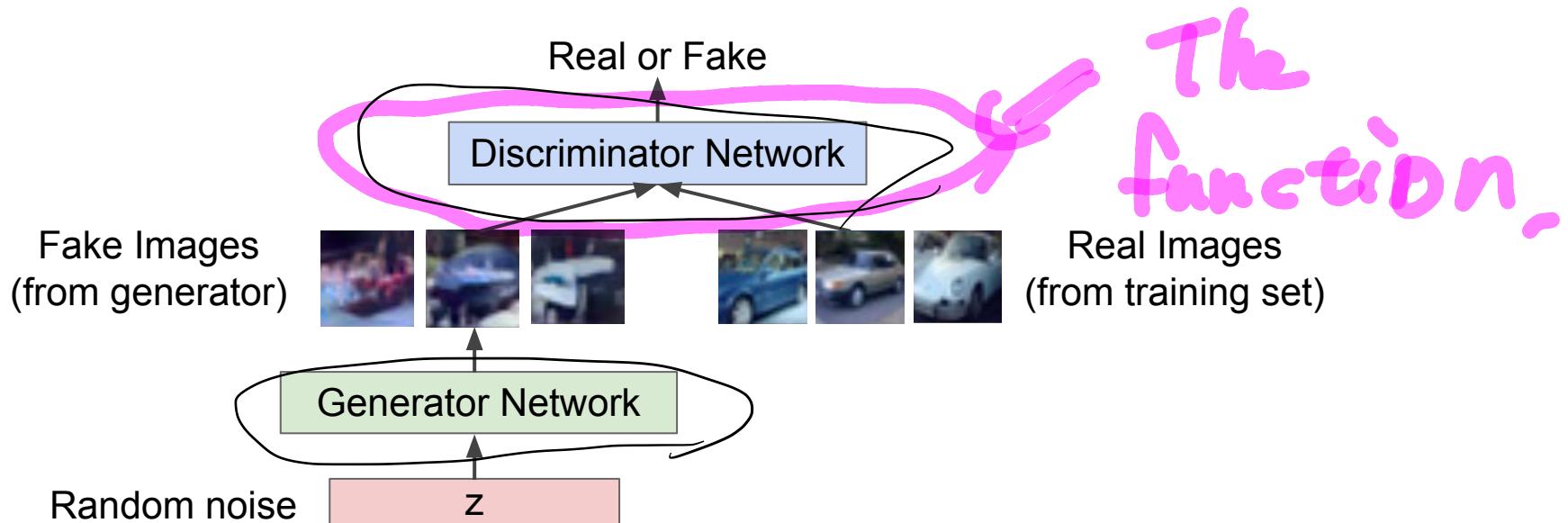
We propose a new framework for estimating generative models via an adversarial process, in which we simultaneously train two models: a generative model G that captures the data distribution, and a discriminative model D that estimates the probability that a sample came from the training data rather than G . The training procedure for G is to maximize the probability of D making a mistake. This framework corresponds to a minimax two-player game. In the space of arbitrary functions G and D , a unique solution exists, with G recovering the training data distribution and D equal to $\frac{1}{2}$ everywhere. In the case where G and D are defined by multilayer perceptrons, the entire system can be trained with backpropagation. There is no need for any Markov chains or unrolled approximate inference networks during either training or generation of samples. Experiments demonstrate the potential of the framework through qualitative and quantitative evaluation of the generated samples.

Training GANs: Two-player game

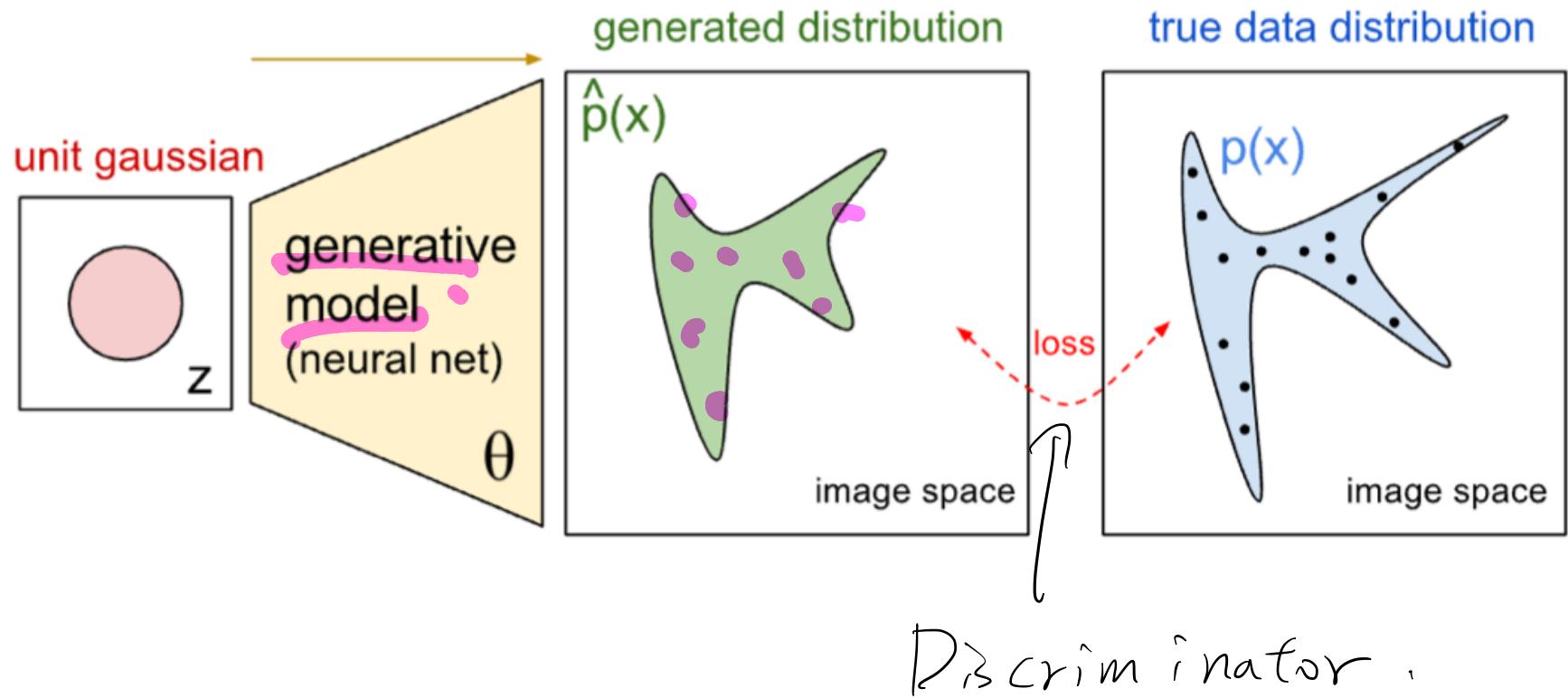
Ian Goodfellow et al., “Generative Adversarial Nets”, NIPS 2014

Generator network: try to fool the discriminator by generating real-looking images

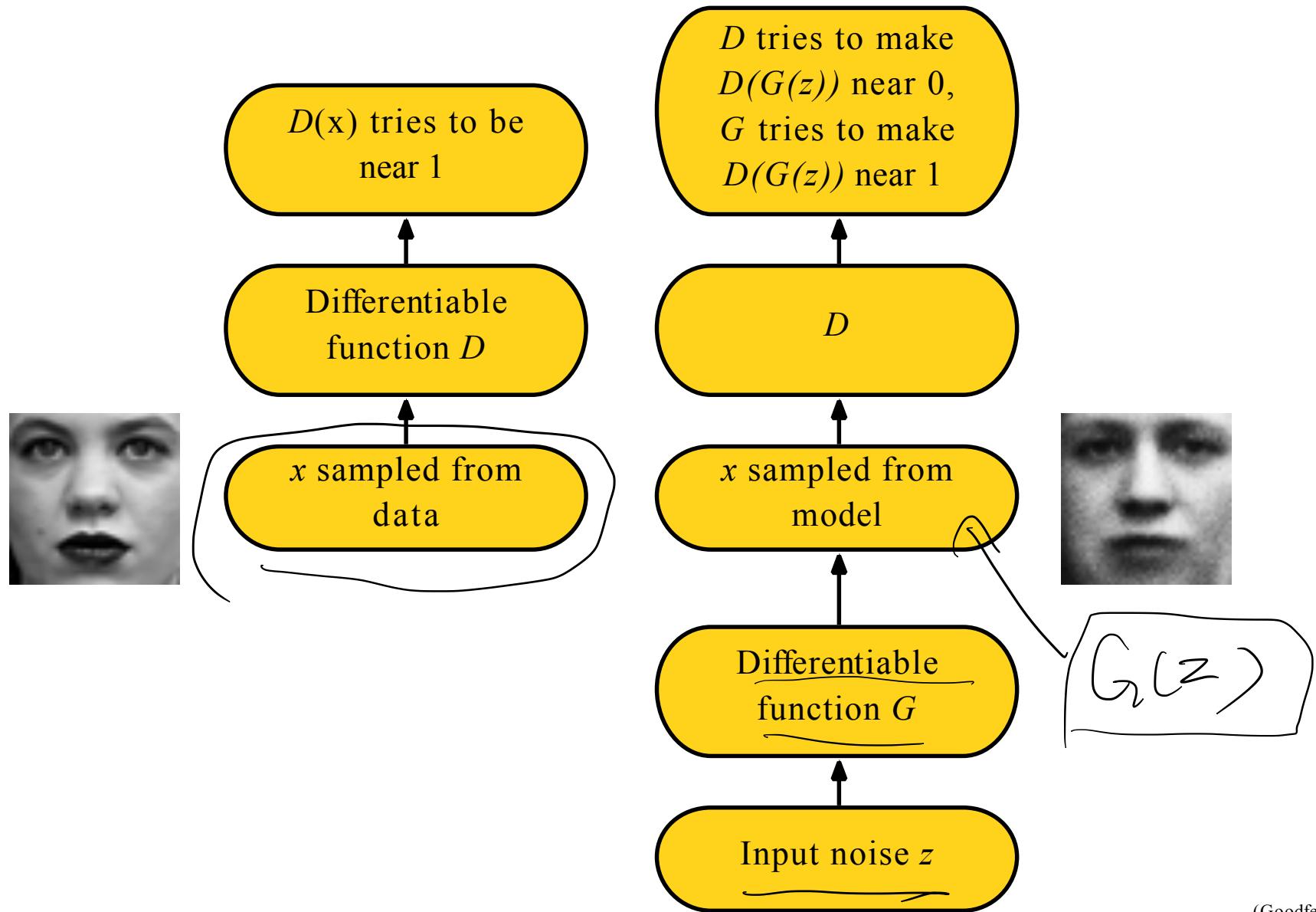
Discriminator network: try to distinguish between real and fake images



Minimax Game



Adversarial Nets Framework



True distribution. $\log(1 - D(G)) \Rightarrow -\log D(G)$

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 $p_{\text{data}}(x)$. \Rightarrow Images in your training set.
- Update the discriminator by ascending its stochastic gradient:

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

maximize
objective for D .

real fake

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:

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

minimize
objective for G .

end for

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

Jensen Shannon (JS)-Divergence

$$JS(P, Q) = \frac{1}{2} KL(P \parallel \frac{P+Q}{2}) + \frac{1}{2} KL(Q \parallel \frac{P+Q}{2})$$

$KL(P, Q)$ vs $JS(P, Q)$ ↗ symmetric
 $P_{\alpha} > 0, Q_{\alpha} = 0 \Rightarrow \exists \alpha$

$$\mathcal{L}(D, G) = \int P_{\text{data}}(x) \log D(x) dx + \int P_{\text{gen}}(x) \log (1 - D(x)) dx$$

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

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

$$\mathcal{L}(D^*, G) = 2 JS(P_{\text{data}}, P_{\text{gen}}) - 2 \log 2.$$

$$\hat{f} = f \circ D$$

Wasserstein GAN

WGAN
f-GAN

Martin Arjovsky¹, Soumith Chintala², and Léon Bottou^{1,2}

¹Courant Institute of Mathematical Sciences

²Facebook AI Research

$$\int P_{\text{data}}(\alpha) \log D(\alpha) d\alpha + \int P_{\text{gen}}(\alpha) \log (1 - D(G(\alpha))). d\alpha$$

$\underbrace{f(D(\alpha))}_{\text{f}(D(\alpha))}$

The *Earth-Mover* (EM) distance or Wasserstein-1

$\rightarrow \mathcal{F}(D(G(\alpha)))$

$$W(\mathbb{P}_r, \mathbb{P}_g) = \inf_{\gamma \in \Pi(\mathbb{P}_r, \mathbb{P}_g)} \mathbb{E}_{(x,y) \sim \gamma} [\|x - y\|], \quad (1)$$

$P_{\text{data}} \quad P_{\text{gen}}$.

where $\Pi(\mathbb{P}_r, \mathbb{P}_g)$ denotes the set of all joint distributions $\gamma(x, y)$ whose marginals are respectively \mathbb{P}_r and \mathbb{P}_g . Intuitively, $\gamma(x, y)$ indicates how much “mass” must be transported from x to y in order to transform the distributions \mathbb{P}_r into the distribution \mathbb{P}_g . The EM distance then is the “cost” of the optimal transport plan.

A Primer on Optimal Transport

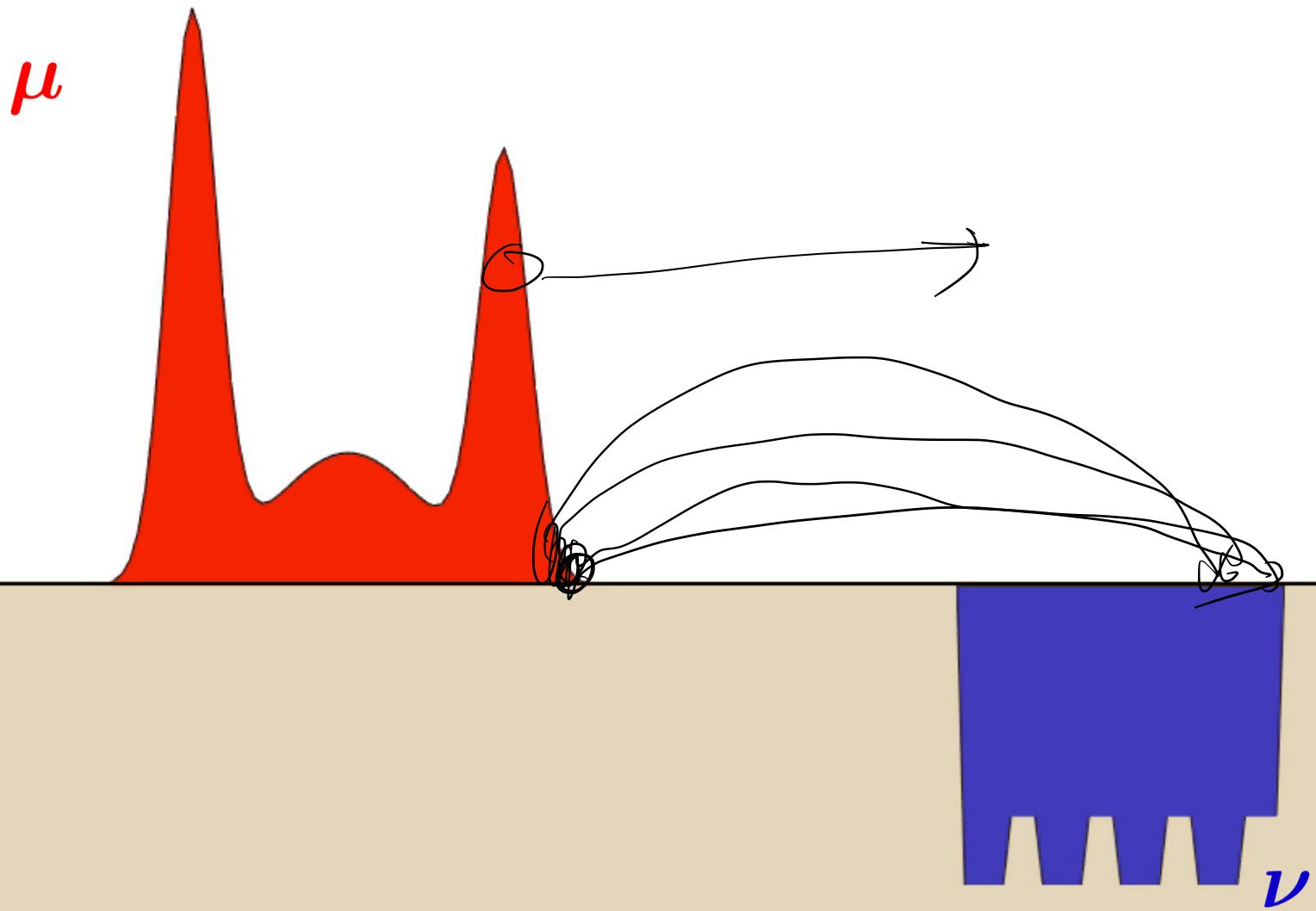
Marco Cuturi



Justin Solomon

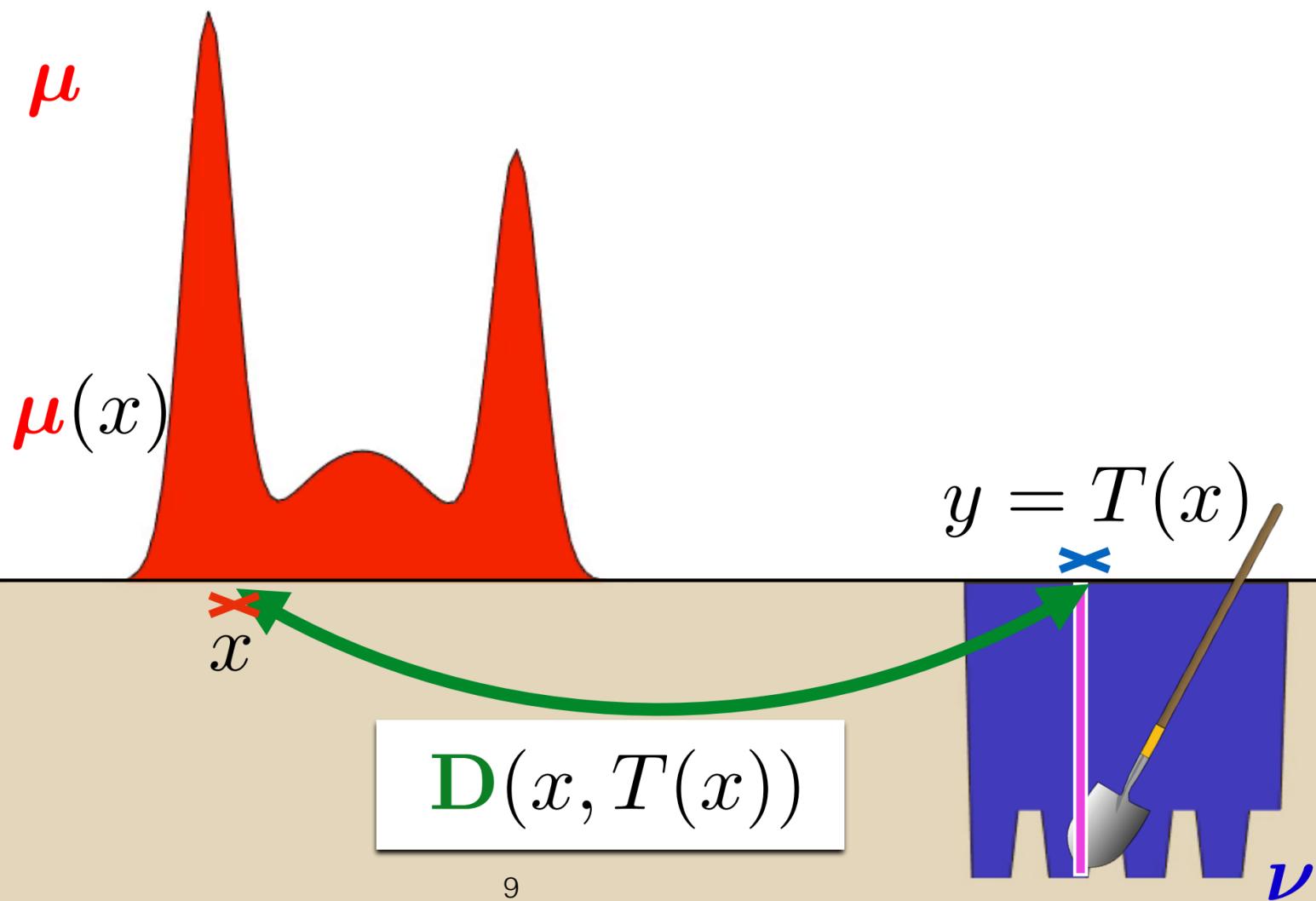


Origins: Monge Problem



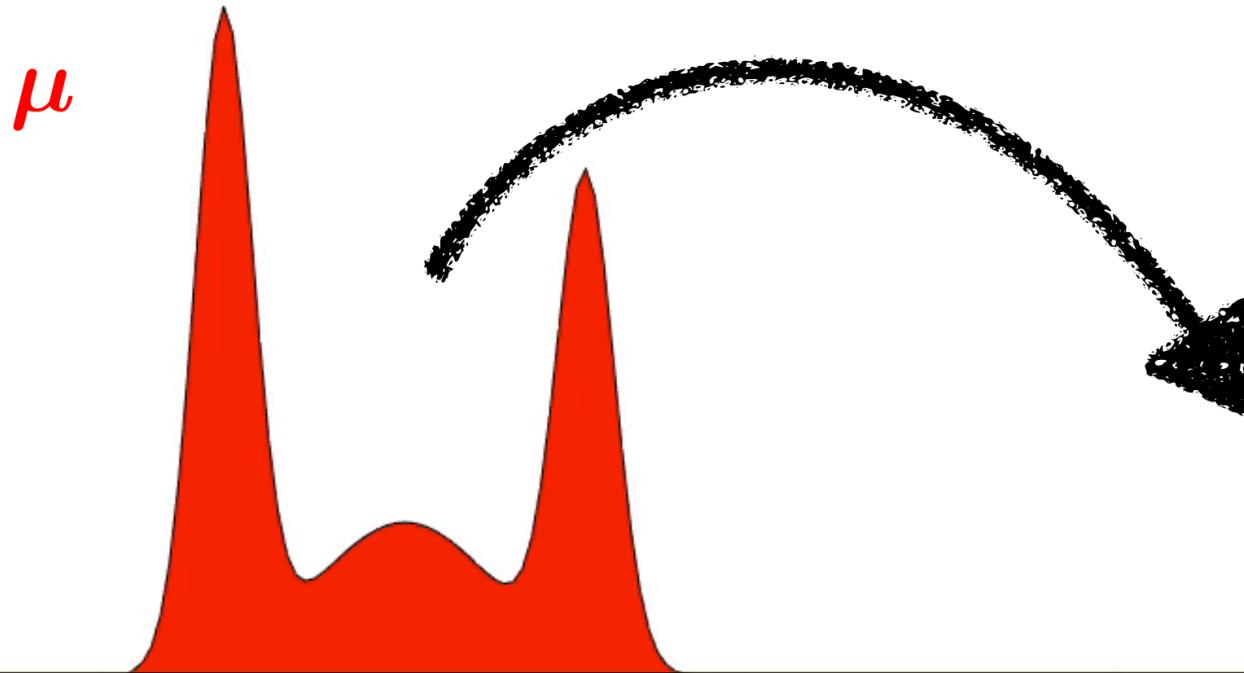
Origins: Monge's Problem

In 1781 however...



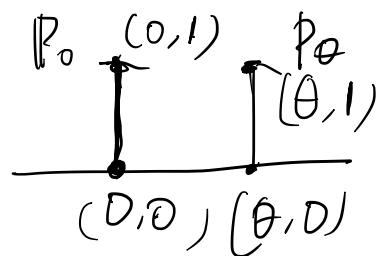
Origins: Monge's Problem

T must **push-forward** the red measure towards the blue



What T s.t. $T_{\sharp}\mu = \nu$
minimizes $\int D(x, T(x)) \mu(dx)$?

JS-Divergence vs. W-distance



Example 1 (Learning parallel lines). Let $Z \sim U[0, 1]$ the uniform distribution on the unit interval. Let \mathbb{P}_0 be the distribution of $(0, Z) \in \mathbb{R}^2$ (a 0 on the x-axis and the random variable Z on the y-axis), uniform on a straight vertical line passing through the origin. Now let $g_\theta(z) = (\theta, z)$ with θ a single real parameter. It is easy to see that in this case,

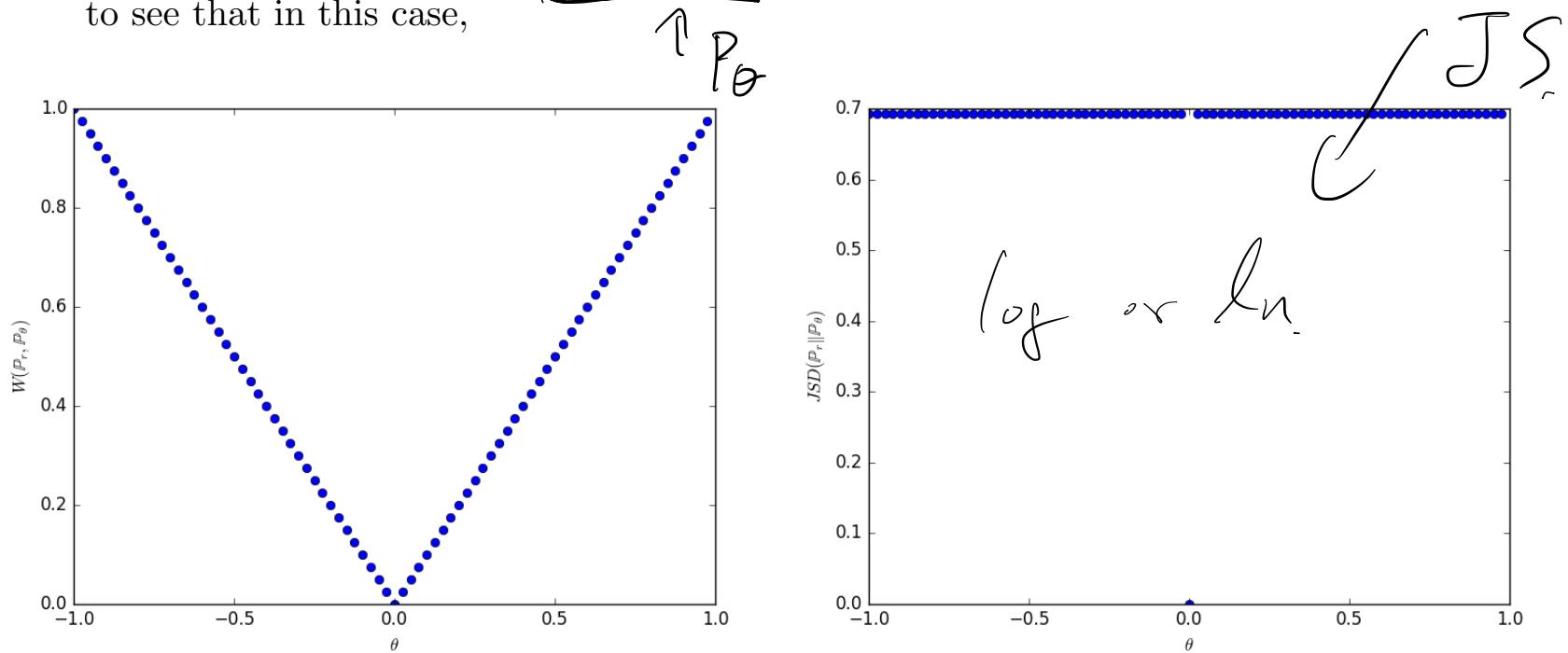


Figure 1: These plots show $\rho(\mathbb{P}_\theta, \mathbb{P}_0)$ as a function of θ when ρ is the EM distance (left plot) or the JS divergence (right plot). The EM plot is continuous and provides a usable gradient everywhere. The JS plot is not continuous and does not provide a usable gradient.

Kantorovich-Rubinstein Duality

$$W(\mathbb{P}_r, \mathbb{P}_\theta) = \sup_{\|f\|_L \leq 1} \mathbb{E}_{x \sim \mathbb{P}_r}[f(x)] - \mathbb{E}_{x \sim \mathbb{P}_\theta}[f(x)]$$

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: : α , the learning rate. c , the clipping parameter. m , the batch size.
 n_{critic} , the number of iterations of the critic per generator iteration.

Require: : w_0 , initial critic parameters. θ_0 , initial generator's parameters.

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

max