

# Wasserstein GAN

Sargur N. Srihari

[srihari@cedar.buffalo.edu](mailto:srihari@cedar.buffalo.edu)

# Topics in GAN

1. GAN Motivation
2. GAN Theory
3. GAN Mode Collapse

## 4. Wasserstein GAN (WGAN)

1. GAN as minimizing JS-Divergence
2. Earth Mover Distance and WGAN
3. Loss in terms of Lipschitz function
4. GAN vs WGAN performance

→ Discussed here

## 5. GAN Variants

# Why study WGAN?

- In ML, unsupervised learning is the cake, supervised learning is the icing and RL the cherry
  - Generative models are the most important class of unsupervised learning
    - GANs are generative models
- WGAN training is backed up by theory
- In WGAN the discriminator is trained to convergence (while it is an art in GAN)
- There is a correlation between discriminator loss and perceptual quality



# Generative Modeling

- Task is one of unsupervised learning
  - To learn a probability density
  - Often done by defining a parametric family of densities  $(P_\theta) \theta \in \mathbb{R}^d$  and finding the one that maximizes likelihood of data
  - if we have real data examples  $\{\mathbf{x}^{(i)}\}, i = 1..m$ , we would solve the problem

$$\max_{\theta \in \mathbb{R}^d} \frac{1}{m} \sum_{i=1}^m \log P_\theta(\mathbf{x}^{(i)})$$

- We can show it is equivalent to minimizing  $KL$  divergence

# Equivalence to minimum $KL$ Divergence

- If real and parameterized densities are  $P_r$  and  $P_\theta$  then it is same as minimizing  $KL(P_r || P_\theta)$ .

- **Proof:**

$$\begin{aligned}
 \lim_{m \rightarrow \infty} \max_{\theta \in \mathbb{R}^d} \frac{1}{m} \sum_{i=1}^m \log P_\theta(x^{(i)}) &= \max_{\theta \in \mathbb{R}^d} \int_x P_r(x) \log P_\theta(x) dx \\
 &= \min_{\theta \in \mathbb{R}^d} - \int_x P_r(x) \log P_\theta(x) dx \\
 &= \min_{\theta \in \mathbb{R}^d} \int_x P_r(x) \log P_r(x) dx - \int_x P_r(x) \log P_\theta(x) dx \\
 &= \min_{\theta \in \mathbb{R}^d} KL(P_r || P_\theta)
 \end{aligned}$$

$$KL(\mathbb{P}_r || \mathbb{P}_g) = \int_{\mathcal{X}} P_r(x) \log \frac{P_r(x)}{P_g(x)} dx$$

- Max. likelihood is same as minimizing  $KL(P_r || P_\theta)$
- For this to make sense we need  $P_\theta$  to exist
  - This is not the case in the rather common situation where we are dealing with distributions supported by low dimensional manifolds
    - See next

# Properties of KL Divergence

- $KL$  has a unique minimum at  $P_g = P_r$ , and
  - it doesn't require knowledge of the unknown  $P_r(\mathbf{x})$  to optimize it (only samples)
- $KL$  is unsymmetrical between  $P_r$  and  $P_g$ :

# Compare $KL(P_r, P_g)$ and $KL(P_g, P_r)$

- If  $P_r(\mathbf{x}) > P_g(\mathbf{x})$ : higher probability  $\mathbf{x}$  is data
  - This is the core of ‘mode dropping’
 

$$KL(\mathbb{P}_r || \mathbb{P}_g) = \int_{\mathcal{X}} P_r(x) \log \frac{P_r(x)}{P_g(x)} dx$$

    - With large regions of high  $P_r$  but small  $P_g$
    - When  $P_r(\mathbf{x}) > 0$  but  $P_g(\mathbf{x}) \rightarrow 0$ , we have  $KL \rightarrow \infty$
    - i.e., high cost to  $G$  not covering parts of data
- If  $P_r(\mathbf{x}) < P_g(\mathbf{x})$ , lower probability  $\mathbf{x}$  is data
  - Case when  $G$ ’s output doesn’t look real
    - When  $P_r(\mathbf{x}) \rightarrow 0$  and  $P_g(\mathbf{x}) > 0$ ,  $KL \rightarrow 0$ ,
    - i.e., low cost for generating fake looking samples
- If we use  $KL(P_g || P_r)$  error weighting reversed
  - high cost for generating implausible pictures

# Jensen-Shannon (JS) Divergence

- GANs, in original formulation, optimize JS divergence (not KL divergence!)

$$JSD(\mathbb{P}_r \parallel \mathbb{P}_g) = \frac{1}{2}KL(\mathbb{P}_r \parallel \mathbb{P}_A) + \frac{1}{2}KL(\mathbb{P}_g \parallel \mathbb{P}_A)$$

- Where  $P_A$  is the average distribution  $(P_r + P_g)/2$

- Proof: GANs are formulated in two steps
  - First train  $D$  to maximize

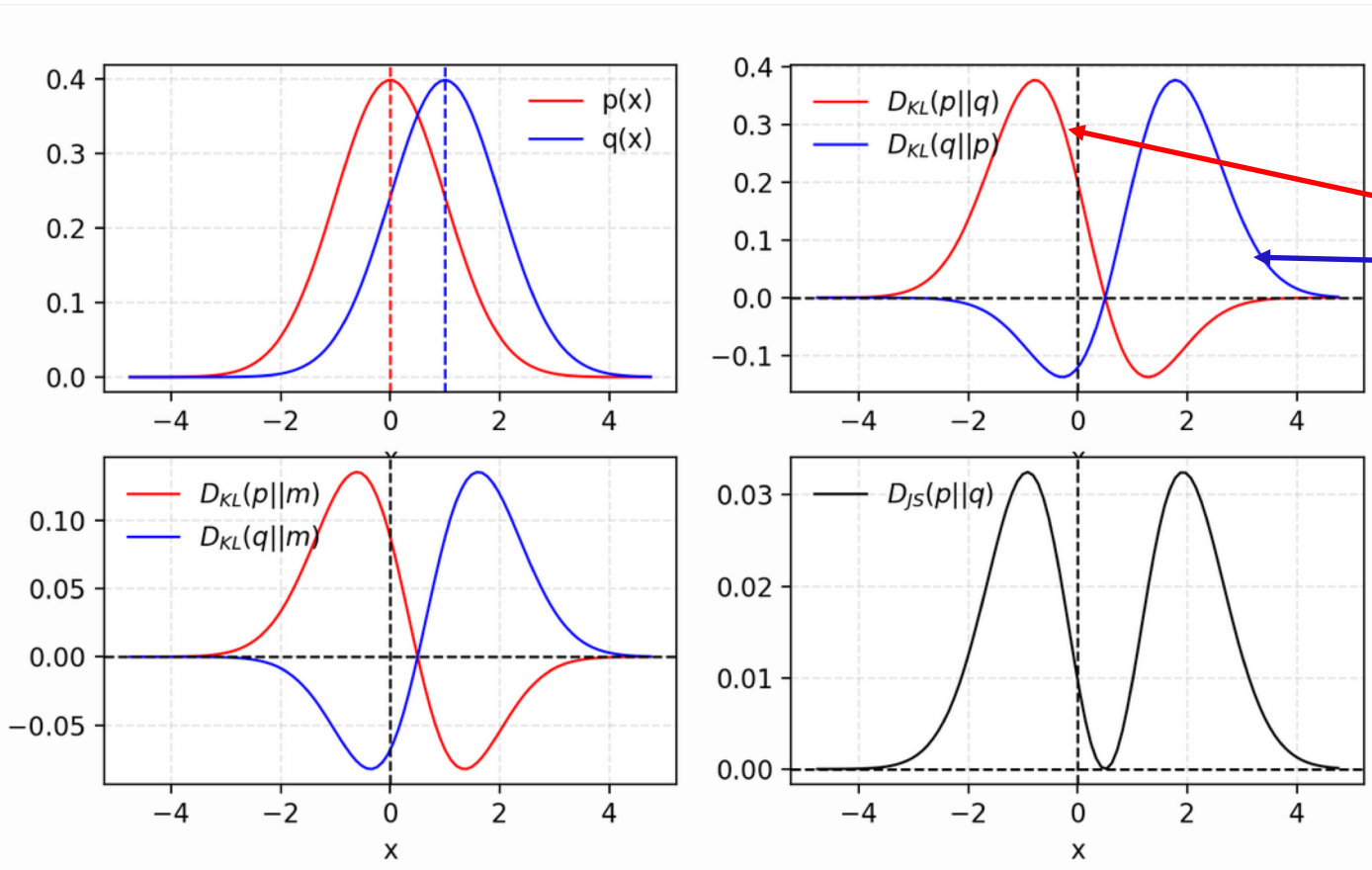
$$L(D, g_\theta) = \mathbb{E}_{x \sim \mathbb{P}_r}[\log D(x)] + \mathbb{E}_{x \sim \mathbb{P}_g}[\log(1 - D(x))]$$

- Optimal discriminator has the shape  $D^*(x) = \frac{P_r(x)}{P_r(x) + P_g(x)}$
- Minimizing  $L(D, g_\theta)$  yields the JS divergence

$$L(D^*, g_\theta) = 2JSD(\mathbb{P}_r \parallel \mathbb{P}_g) - 2 \log 2, \text{ where } JS(P_r, P_g) = \frac{1}{2}KL(P_r \parallel P_m) + \frac{1}{2}KL(P_g \parallel P_m)$$



# Integrands of KL and JS divergence



Areas under curve to be integrated

$\int p(x) \log p(x)/q(x)$

$\int q(x) \log q(x)/p(x)$

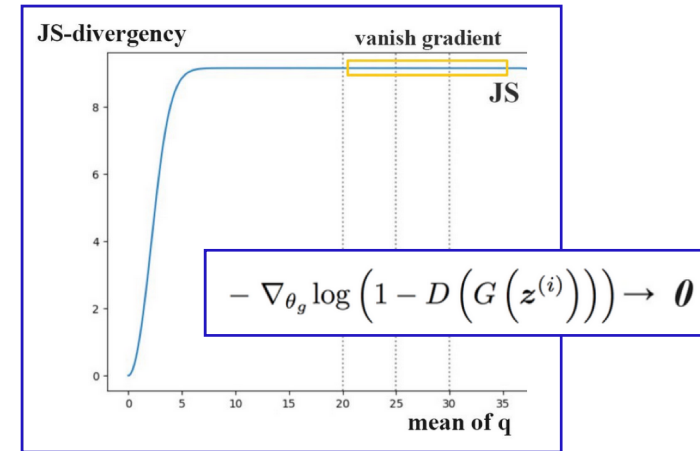
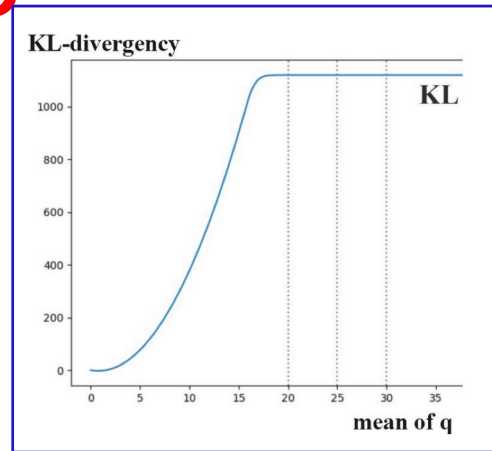
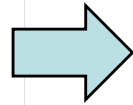
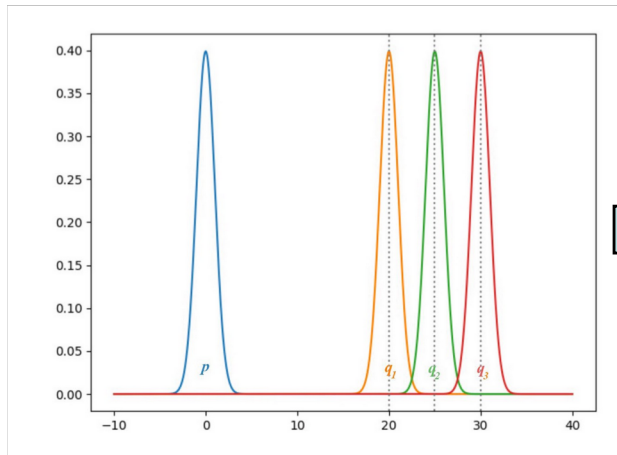
In this case, the areas seem identical, but that wouldn't be true in general

$p \sim N(0,1)$  and  $q \sim N(1,1)$   
 $m = (p+q)/2$

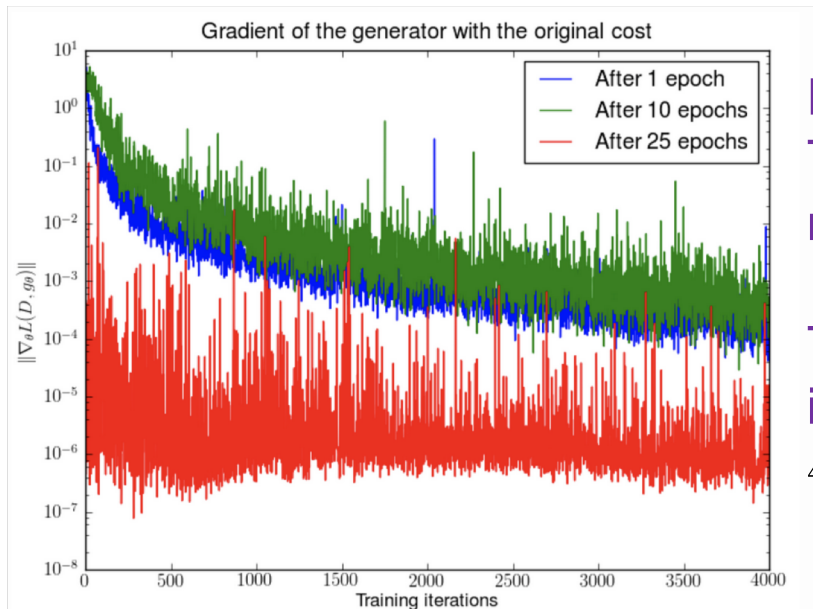
$$D_{JS}(p||q) = \frac{1}{2} D_{KL}(p||\frac{p+q}{2}) + \frac{1}{2} D_{KL}(q||\frac{p+q}{2})$$

$D_{KL}$  integrand is asymmetric in each of the two cases  
 $D_{JS}$  integrand is symmetric

# Vanishing Gradient Problem



When  $D$  is perfect, we are guaranteed with  $D(x)=1, \forall x \in p_r$  and  $D(x)=0, \forall x \in p_g$   
 Therefore loss  $L$  falls to zero and there is no gradient to update  $L$  during learning  
 Figure below: when the discriminator gets better, the gradient vanishes fast



First, a DCGAN is trained for 1, 10 and 25 epochs.  
 Then, with the  $G$  fixed,  $D$  is trained from scratch and measure the gradients with the original cost function

The gradient norms **decay quickly** (in log scale),  
 in the best case 5 orders of magnitude after  
 4000 discriminator iterations

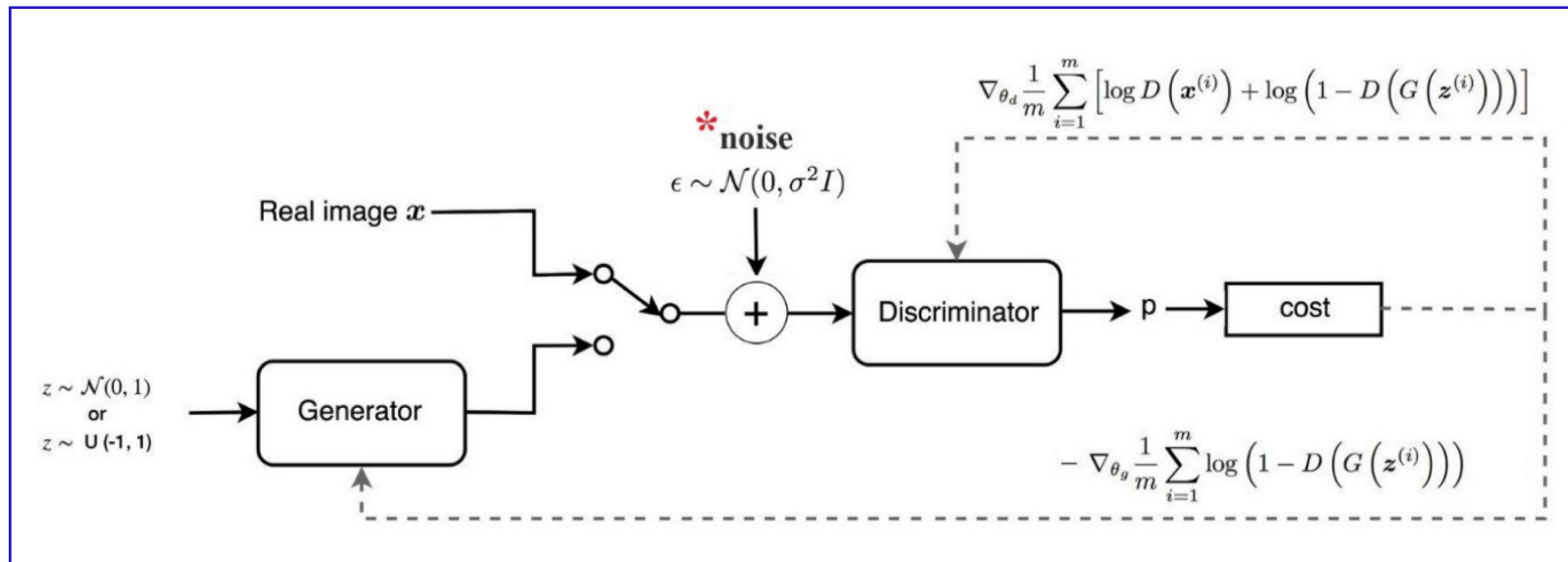
# Vanishing Gradient Proposals

Proposals to address vanishing gradient problem:

- $\nabla_{\theta_g} \log D(G(z^{(i)}))$

- But it is unstable with large variance of gradients

2. Adding noise to generated images to stabilize the model: ↩



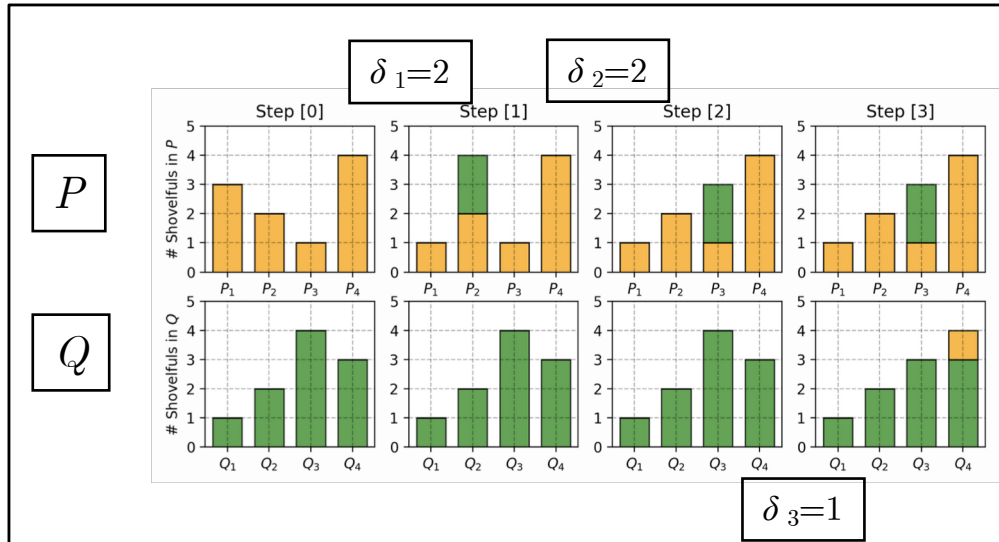
3. New cost function based on Wasserstein distance

# WGAN cost function

- Instead of adding noise, Wasserstein GAN (WGAN) proposes a new cost function using Wasserstein distance that has a smoother gradient everywhere
- It is based on the Earth Mover Distance between probability distributions

# Earth Mover Distance

In order to change  $P$  to look like  $Q$



Move 2 units from  $P_1$  to  $P_2 \rightarrow (P_1, Q_1)$  match  
 Move 2 units from  $P_2$  to  $P_3 \rightarrow (P_2, Q_2)$  match  
 Move 1 unit  $Q_3$  to  $Q_4 \rightarrow (P_3, Q_3), (P_4, Q_4)$  match

Here the EMD is

$$W = \sum |\delta_i| = 5$$

## Wasserstein-1 Definition

$$W(P_r, P_g) = \inf_{\gamma \in \Pi(P_r, P_g)} \mathbb{E}_{(x,y) \sim \gamma} [ \|x - y\| ]$$

Infimum is greatest lower bound



where

$\Pi(P_r, P_g)$  denotes set of all joint distributions  $\gamma(x, y)$  with marginals  $P_r, P_g$

$\gamma(x, y)$  is “mass” to be transported from  $x$  to  $y$  in order to transform  $P_r$  into  $P_g$

EMD then is the “cost” of the optimal transport plan

# Distance measures $\rho(P_r, P_g)$

- Total Variation distance

- Largest possible difference between probability distributions

$$\delta(P_r, P_g) = \sup_A |P_r(A) - P_g(A)|$$

- Kullback-Leibler (KL) Divergence (KL=Max Likelihood)

- Measure of distribution difference from reference distribution

$$KL(P_r \| P_g) = \int_x \log \left( \frac{P_r(x)}{P_g(x)} \right) P_r(x) dx$$

- Jensen-Shannon Divergence (JSD=GAN)

- Symmetrized and smoothed KL divergence

$$JS(P_r, P_g) = \frac{1}{2} KL(P_r \| P_m) + \frac{1}{2} KL(P_g \| P_m)$$

$P_m$  is the mixture  $(P_r + P_g)/2$

- Earth-Mover Distance (EMD=WGAN)

- Minimum cost of moving one pile to another

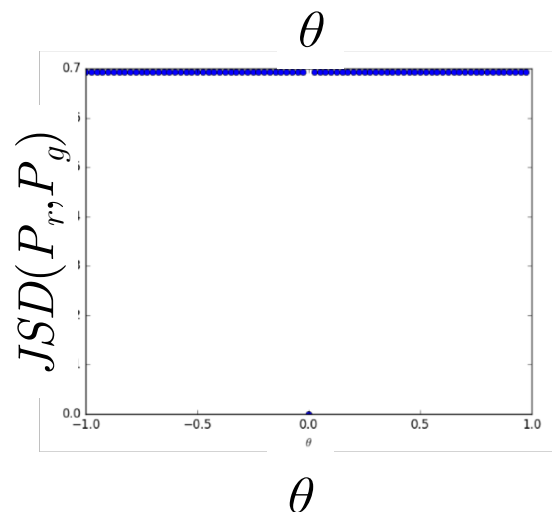
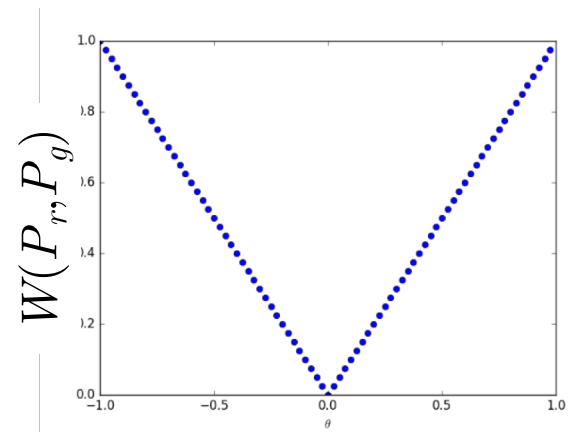
$$W(P_r, P_g) = \inf_{\gamma \in \Pi(P_r, P_g)} \mathbb{E}_{(x,y) \sim \gamma} [\|x - y\|]$$

Infimum is  
greatest lower bound



# Gradient for learning parallel lines

- $Z \sim U[0,1]$  is uniform on unit interval
- Let  $P_0$  be distribution of  $(0, Z) \in \mathbb{R}^2$ 
  - 0 on  $x$ -axis and  $Z$   $y$ -axis,
  - uniform on a straight vertical line passing through origin
- Now let  $g_\theta(z) = (\theta, z)$  with  $\theta$  a single real parameter
  - Then  $W(P_0, P_\theta) = |\theta|$
  - $JS(P_0, P_\theta) = \log 2$  if  $\theta \neq 0$ ,
  - $= 0$  if  $\theta = 0$

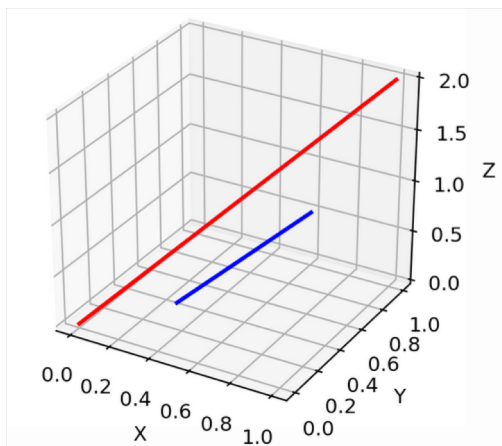


EMD is continuous and provides a usable gradient  
 JSD is discontinuous and there is no usable gradient

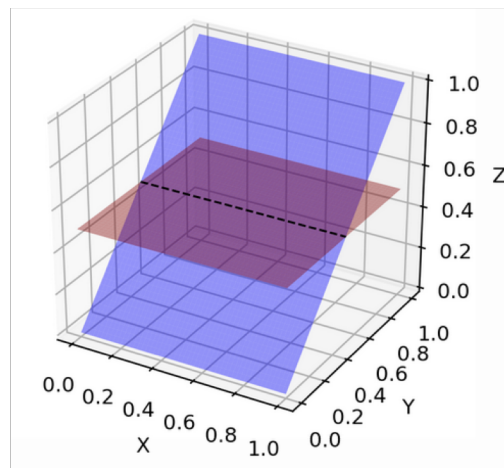
# Intersection of $P_\theta$ and $P_r$

Low dimensional manifolds in high-dimensional space  
Can hardly have overlaps

Two lines in 3-D space



Two surfaces in 3-D space



It is unlikely that model manifold and the true distribution's support have a non-negligible intersection

This means that the KL distance is not defined (or simply infinite)



# Earth Mover (Wasserstein 1)

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

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

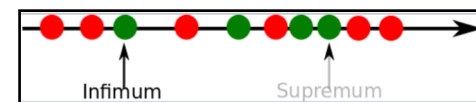
# WGAN loss using a Lipschitz

- EMD  $W(P_r, P_g) = \inf_{\gamma \in \Pi(P_r, P_g)} \mathbb{E}_{(x,y) \sim \gamma} [\|x - y\|]$  is intractable

– Using Kantorovich-Rubinstein duality, simplify to

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

– where sup is the least upper bound



–  $f$  is a 1-Lipschitz function following the constraint:

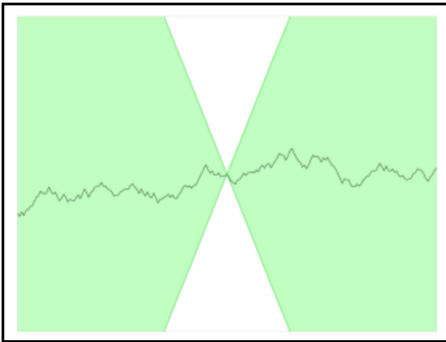
$$|f(x_1) - f(x_2)| \leq |x_1 - x_2|$$

- i.e, a function with a bounded first derivative

- So to calculate the Wasserstein distance, we just need to find a 1-Lipschitz function
- Like other deep learning problems, we can build a deep network to learn it

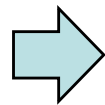
# What is a Lipschitz function?

## A Lipschitz function



A Lipschitz function is such that

$$|f(x) - f(y)| \leq \alpha |x - y|$$



$$\left| \frac{f(x) - f(y)}{x - y} \right| \leq \alpha$$

i.e., slope of secant line joining  $(x, f(x))$  and  $(y, f(y))$  is always bounded above by  $\alpha$

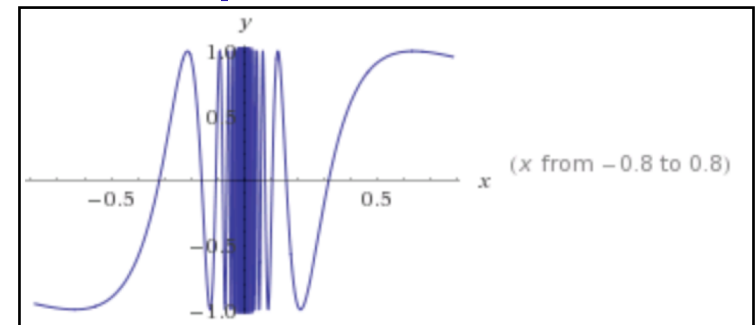
Examples of Lipschitz function are  $\sin x$  and  $x$

A function which is not Lipschitz but is bounded is  $\sin(x^2)$  over  $R$

This is because as we go further towards  $+\infty$ , the oscillation becomes faster and

Thus the slope of the secant lines get nearer to vertical ones

## Not a Lipschitz function



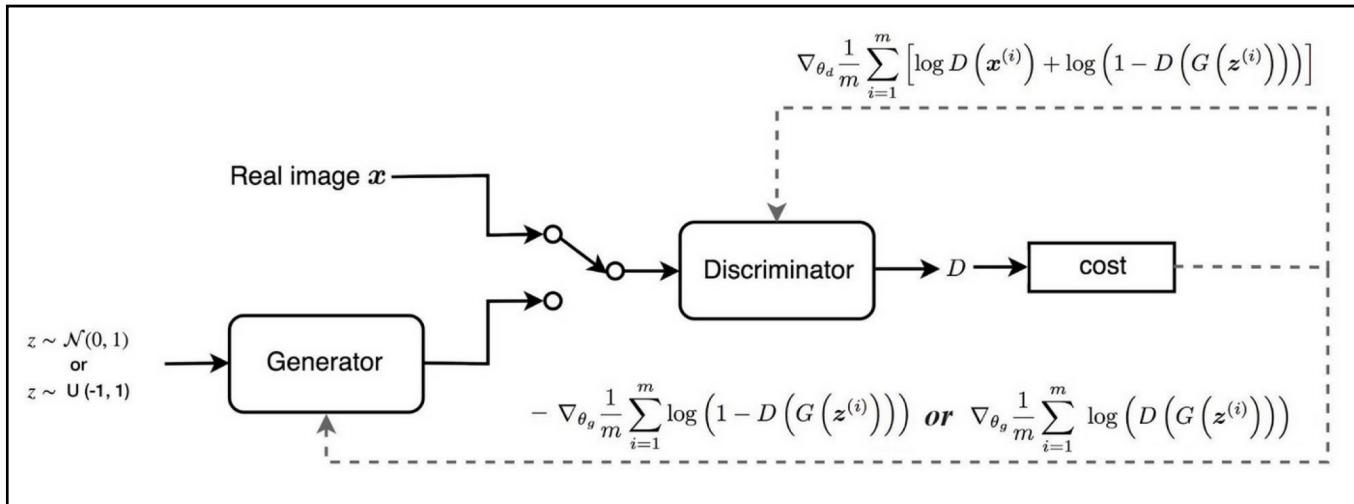
Although this is a bounded function

# Wasserstein Critic Network

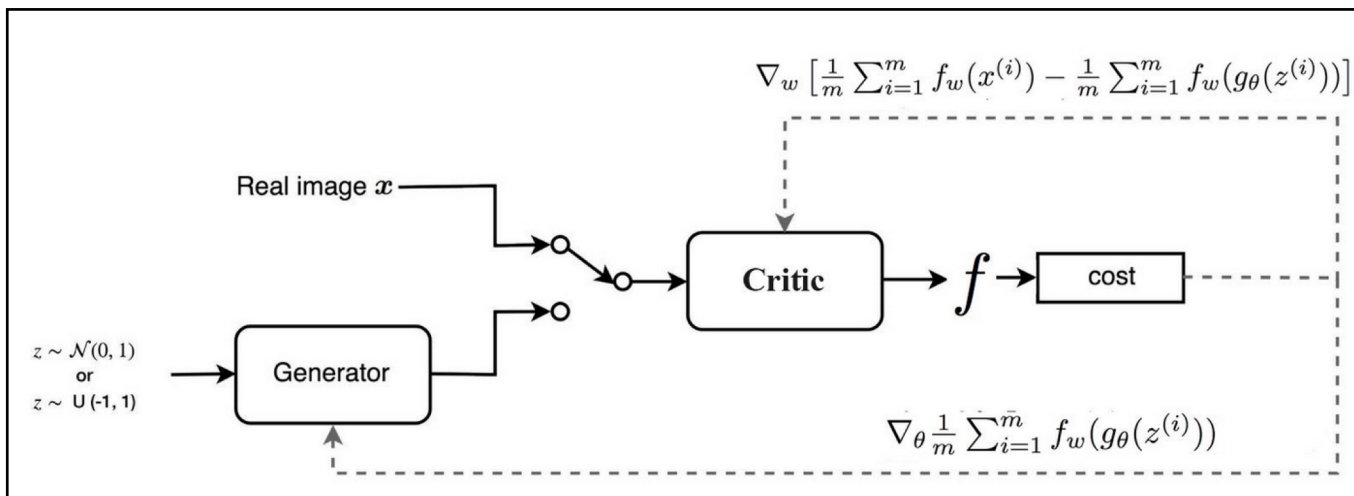
- This network is very similar to  $D$ ,
  - Just without the sigmoid function and outputs a scalar score rather than a probability
  - Score interpreted as how real the input images are
  - In RL, we call it the *value function* which measures how good a state (the input) is
- Rename  $D$  as *critic* to reflect its new role

# GAN and WGAN Architectures

## GAN



## WGAN



# Steps in WGAN Algorithm

- We'd like to train  $P_{\theta} = g_{\theta}(Z)$  to match  $P_r$ , Intuitively,
  - Given  $g_{\theta}$ , can compute optimal  $f_w$  for EMD
  - Then backprop through  $W(P_r, g_{\theta}(Z))$  to get gradient:  

$$\nabla_{\theta} W(P_r, P_{\theta}) = \nabla_{\theta} (E_{x \sim P_r} [f_w(x)] - E_{z \sim Z} [f_w(g_{\theta}(z))]) = -E_{z \sim Z} [\nabla_{\theta} f_w(g_{\theta}(z))]$$
- Training process now has three steps
  1. For fixed  $\theta$ , get  $W(P_r, P_{\theta})$  by training  $f_w$  to convergence
  2. Compute gradient  $-E_{z \sim Z} [\nabla_{\theta} f_w(g_{\theta}(z))]$  by sampling  $z \sim Z$
  3. Update  $\theta$ , and repeat process
- For  $\{f_w\}_{w \in W}$  to be K-Lipschitz use weight clamping:
  - Weights  $w$  are constrained to lie within  $[-c, c]$ , by clipping  $w$  after every update to  $w$

# WGAN Algorithm

---

**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.  $n_{\text{critic}}$ , the number of iterations of the critic per generator iteration.

**Require:** :  $w_0$ , initial critic parameters.  $\theta_0$ , initial generator's parameters.

```

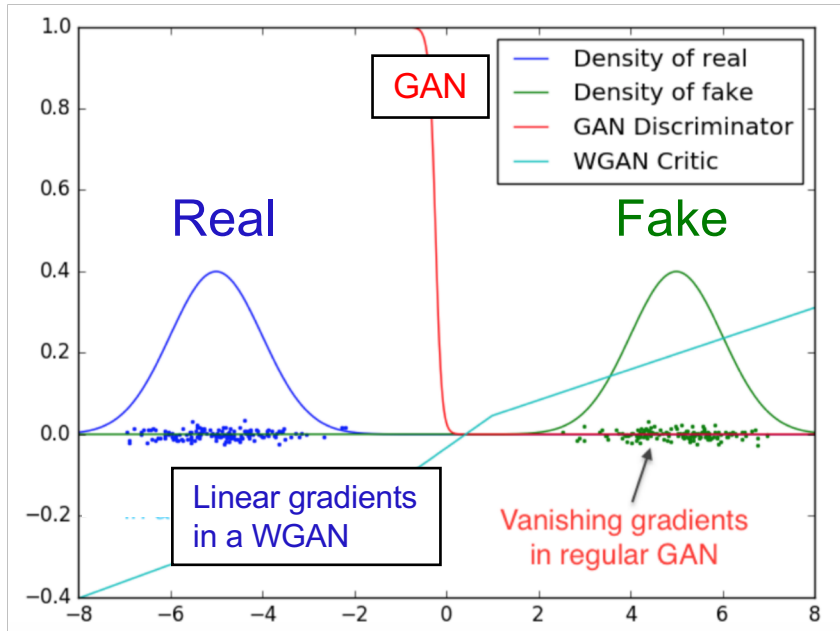
1: while  $\theta$  has not converged do
2:   for  $t = 0, \dots, n_{\text{critic}}$  do
3:     Sample  $\{x^{(i)}\}_{i=1}^m \sim \mathbb{P}_r$  a batch from the real data.
4:     Sample  $\{z^{(i)}\}_{i=1}^m \sim p(z)$  a batch of prior samples.
5:      $g_w \leftarrow \nabla_w \left[ \frac{1}{m} \sum_{i=1}^m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

```

---

# WGAN Critic

- A similar of value of  $D(x)$  for GAN and WGAN



- Optimal discriminator and critic when learning to differentiate two Gaussians.
- The discriminator of a minimax GAN saturates and results in vanishing gradients
- WGAN critic provides very clean gradients on all parts of the space

– GAN (the red line):

- fills with areas with diminishing or exploding gradients

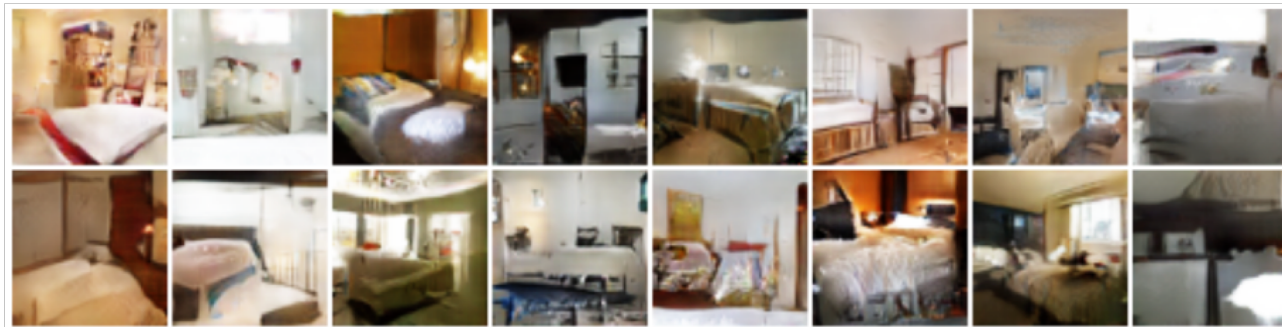
– WGAN (the blue line):

- gradient is smoother everywhere and learns better even the generator is not producing good images



# Performance Comparison 1

- Algorithms trained with a DCGAN generator
  - WGAN



- Standard GAN



- Both algorithms produce high quality samples

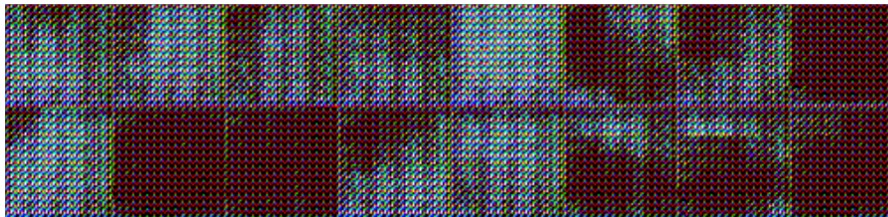
# Performance Comparison 2

- Algorithms trained with a generator
  - without batch normalization and constant no of filters at every layer (as opposed to duplicate every time)

## – WGAN



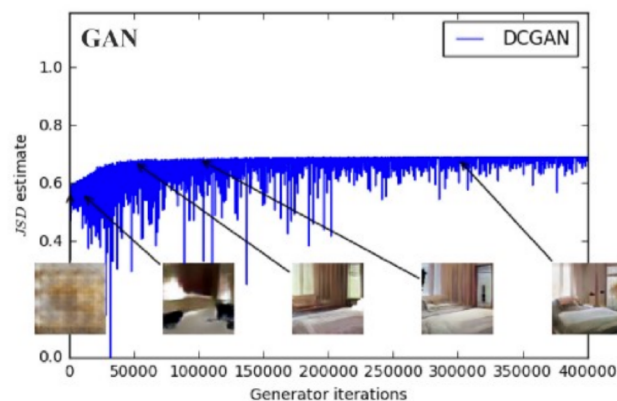
## – Standard GAN



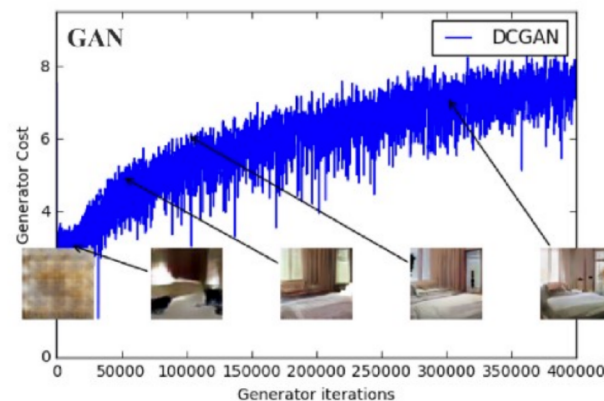
- Standard GAN failed to learn while the WGAN still was able to produce samples

# Quantitative Measure in WGAN

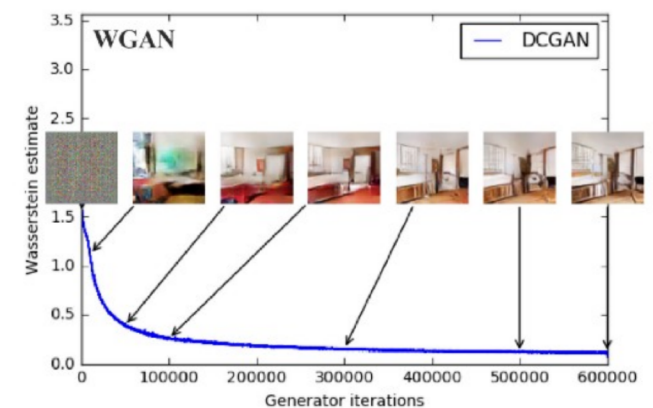
- GAN loss measures how well it fools  $D$  rather than image quality
  - $G$  loss doesn't drop even as image quality improves
  - Hence, we cannot tell progress from its value.
    - We need to save testing images and evaluate it visually.
- WGAN loss function reflects the image quality



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



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



$$\frac{1}{m} \sum_{i=1}^m f_w(x^{(i)}) - \frac{1}{m} \sum_{i=1}^m f_w(g_\theta(z^{(i)}))$$

# Summary

- Maximum likelihood of parameters is same as minimizing KL distance
  - VAEs do this but there is blurring
- WGAN is an alternative to GAN training
  - More stable learning, get rid of problems like mode collapse, and provide meaningful learning curves useful for debugging and hyperparameter searches
- Theoretically
  - The optimization problem is sound
  - There are deep connections to other distances between distributions

# References

- <https://arxiv.org/pdf/1701.07875.pdf>
- <https://arxiv.org/pdf/1701.04862.pdf>
- <https://www.alexirpan.com/2017/02/22/wasserstein-gan.html>
- <https://lilianweng.github.io/lil-log/2017/08/20/from-GAN-to-WGAN.html>