

# Generative Modelling with Neural Networks

M2DS Alternants Research Seminar Course 2022

# Outline

Some Applications of Generative Adversarial Networks

Brief Intros on Neural Network

Generative Model: Generative Adversarial Network

Problems with training GANs

Techniques to improve GANs Training

Wasserstein GAN

# Outline

Some Applications of Generative Adversarial Networks

Brief Intros on Neural Network

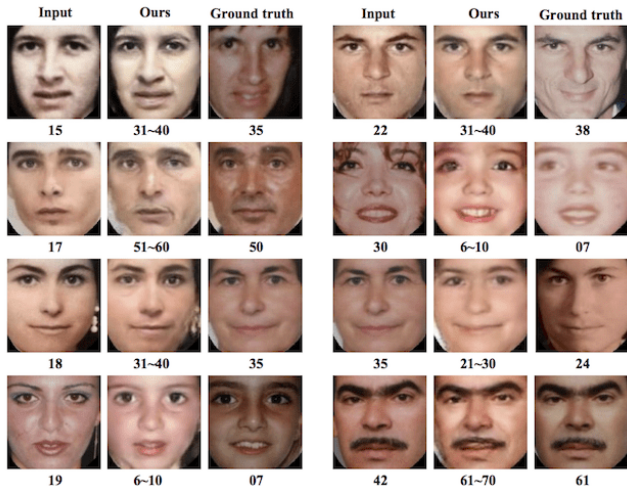
Generative Model: Generative Adversarial Network

Problems with training GANs

Techniques to improve GANs Training

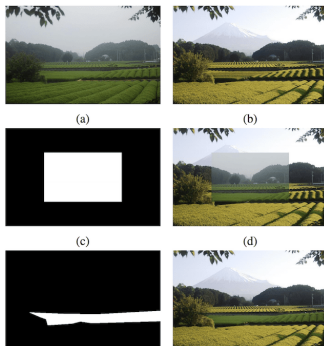
Wasserstein GAN

# Image Generation





# Photo Editing



Real image



Reconstructed images



Blonde ↑

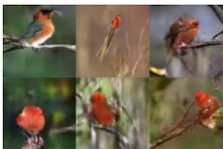
Bangs ↑

Smile ↑

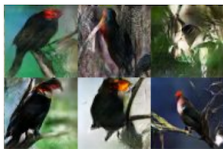
Male ↑

# Text to Image Synthesis

this small bird has a pink breast and crown, and black primaries and secondaries.



this magnificent fellow is almost all black with a red crest, and white cheek patch.



the flower has petals that are bright pinkish purple with white stigma



this white and yellow flower have thin white petals and a round yellow stamen



# Making Music

## GANSynth: Making music with GANs

Feb 25, 2019

Jesse Engel  jesseengel  jesseengel

In this post, we introduce GANSynth, a method for generating high-fidelity audio with Generative Adversarial Networks (GANs).



<https://magenta.tensorflow.org/gansynth>

# Outline

Some Applications of Generative Adversarial Networks

**Brief Intros on Neural Network**

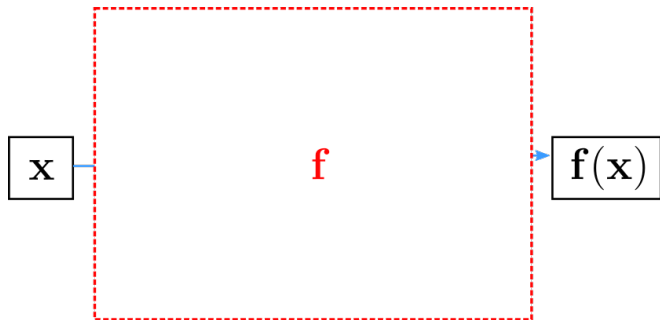
Generative Model: Generative Adversarial Network

Problems with training GANs

Techniques to improve GANs Training

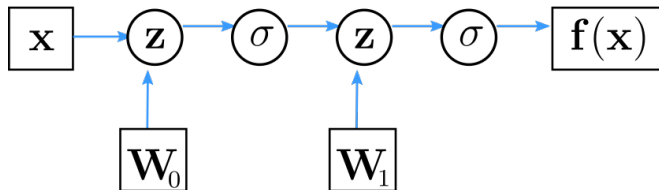
Wasserstein GAN

# Computation Graph



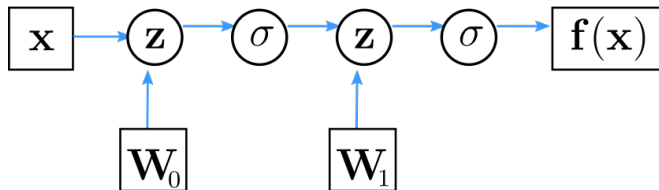
Neural network = parametrized, non-linear function

## Computation Graph



Computation graph: Directed graph of functions, depending on parameters (neuron weights)

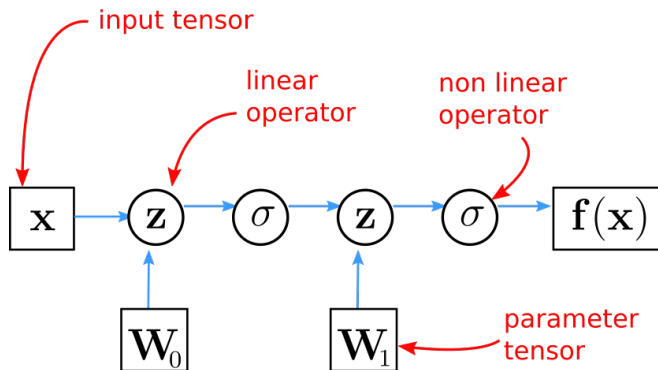
# Computation Graph



Computation graph: Directed graph of functions, depending on parameters (neuron weights)

Question: can you write the full expression of  $f(x)$  for this graph?

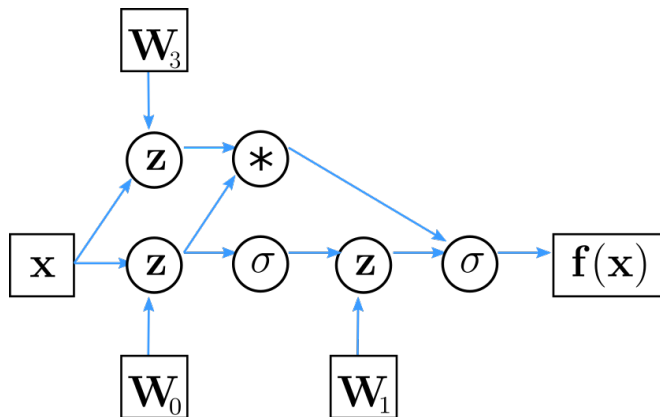
# Computation Graph



Combination of linear (parametrized) and non-linear functions

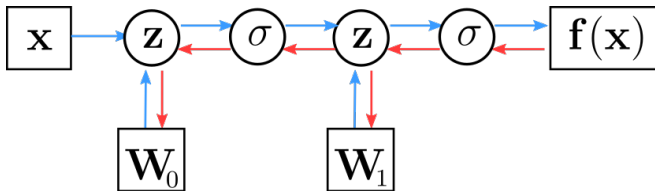


## Computation Graph



Not only sequential application of functions

# Computation Graph



- ▶ Automatic computation of gradients: all modules are *differentiable*!
- ▶ Vector computation on *CPU* and accelerators (*GPU* and *TPU*).



PYTORCH



Microsoft  
CNTK

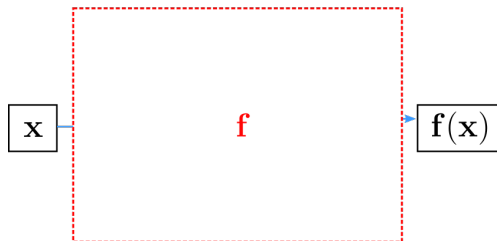
Caffe2

dmlc  
mxnet

gensim spaCy

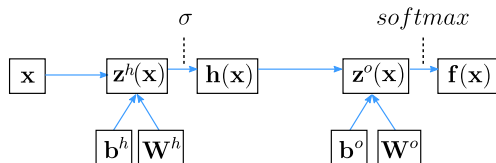
theano

# Formalization for Classification Problem



- ▶ Input:  $\{\mathbf{x}, \mathbf{y}\}_{i=1}^n$  with  $\mathbf{y} \in \{0, 1\}$
- ▶  $f(\cdot, \theta) : \mathbb{R}^p \rightarrow (0, 1)$
- ▶ Output:  $f(\mathbf{x}_i, \theta) = \mathbb{P}(\mathbf{y}_i = 1 \mid \mathbf{x} = \mathbf{x}_i)$

## Simple example: One-hidden layer network



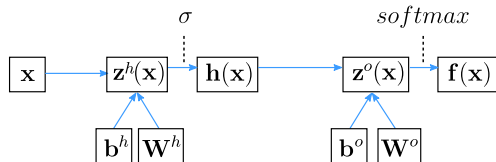
$$f(\mathbf{x}) = \sigma(\mathbf{w}^T \mathbf{x} + b)$$

- ▶  $\mathbf{z}^h(\mathbf{x}) = \mathbf{W}^h \mathbf{x} + \mathbf{b}^h$
- ▶  $\mathbf{h}(\mathbf{x}) = g(\mathbf{z}^h(\mathbf{x})) = g(\mathbf{W}^h \mathbf{x} + \mathbf{b}^h)$
- ▶  $\mathbf{z}^o(\mathbf{x}) = \mathbf{W}^o \mathbf{h}(\mathbf{x}) + \mathbf{b}^o$
- ▶  $\mathbf{f}(\mathbf{x}) = \text{softmax}(\mathbf{z}^o) = \text{softmax}(\mathbf{W}^o \mathbf{h}(\mathbf{x}) + \mathbf{b}^o)$

Where:

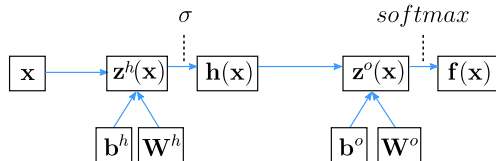
- ▶  $(z(\cdot))$  pre-activation
- ▶  $w, b$  weights and bias
- ▶  $\sigma$  activation function

## Simple example: One-hidden layer network



Question: What is this softmax thing doing here?

## Simple example: One-hidden layer network



Question: What is this softmax thing doing here?

$$\text{softmax}(\mathbf{x}) = \frac{1}{\sum_{i=1}^n e^{x_i}} \cdot \begin{bmatrix} e^{x_1} \\ e^{x_2} \\ \vdots \\ e^{x_n} \end{bmatrix}$$

$$\frac{\partial \text{softmax}(\mathbf{x})_i}{\partial x_j} = \begin{cases} \text{softmax}(\mathbf{x})_i \cdot (1 - \text{softmax}(\mathbf{x})_i) & i = j \\ -\text{softmax}(\mathbf{x})_i \cdot \text{softmax}(\mathbf{x})_j & i \neq j \end{cases}$$

# Training the network

Objective: Finding parameters  $\theta = (W^h; b^h; W^o; b^o)$  that minimize the *negative log likelihood*

The loss function for a given sample  $s \in S$ :

$$\ell(f(\mathbf{x}^s; \theta), \mathbf{y}^s) = -\log f(\mathbf{x}^s; \theta)_{\mathbf{y}^s}$$

# Training the network

Objective: Finding parameters  $\theta = (W^h; b^h; W^o; b^o)$  that minimize the *negative log likelihood*

The loss function for a given sample  $s \in S$ :

$$\ell(f(\mathbf{x}^s; \theta), y^s) = -\log f(\mathbf{x}^s; \theta)_{y^s}$$

→ when computed on full training set:

$$L(\theta) = -\frac{1}{n} \sum_{i=1}^n \log f(\mathbf{x}^s; \theta)_{y^s}$$

or even with  $\ell_2$ -regularization term

$$L(\theta) = -\frac{1}{n} \sum_{i=1}^n \log f(\mathbf{x}^s; \theta)_{y^s} + \lambda(\|W^h\|^2 + \|W^o\|^2)$$



## Training the network

Objective: Finding parameters  $\theta = (W^h; b^h; W^o; b^o)$  that minimize the *negative log likelihood*

The loss function for a given sample  $s \in S$ :

$$\ell(f(\mathbf{x}^s; \theta), y^s) = -\log f(\mathbf{x}^s; \theta)_{y^s}$$

→ when computed on full training set:

$$L(\theta) = -\frac{1}{n} \sum_{i=1}^n \log f(\mathbf{x}^s; \theta)_{y^s}$$

or even with  $\ell_2$ -regularization term

$$L(\theta) = -\frac{1}{n} \sum_{i=1}^n \log f(\mathbf{x}^s; \theta)_{y^s} + \lambda(\|W^h\|^2 + \|W^o\|^2)$$

Question (reminder): how to find  $\theta$ ?

# Gradient Descent

Gradient descent:

$$\theta^{(k+1)} = \theta^{(k)} - \eta \nabla_{\theta} L(\theta^{(k)})$$

Question: what is the potential problem with this?

# Gradient Descent

Gradient descent:

$$\theta^{(k+1)} = \theta^{(k)} - \eta \nabla_{\theta} L(\theta^{(k)})$$

Question: what is the potential problem with this?

→ problematic with large  $n$ : taking long time to train the network

# Stochastic Gradient Descent

1. Initialize  $\theta$  randomly
2. For  $K$  epochs perform:
  - ▶ Important: randomly select a small batch of samples ( $B \subset [n]$ )
  - ▶ Do gradient descent with only this batch of samples:

$$\theta^{(k+1)} = \theta^{(k)} - \eta \nabla_{\theta} L_B(\theta^{(k)})$$

3. Stop when reaching criterion:  
 $L(\theta^{(k)})$  stops decreasing when computed on validation set

# Backpropagation: Computing Gradients Efficiently

Output Weights:

$$\frac{\partial \ell(\mathbf{f}(\mathbf{x}), \mathbf{y})}{\partial W_{i,j}^o}$$

Hidden Weights:

$$\frac{\partial \ell(\mathbf{f}(\mathbf{x}), \mathbf{y})}{\partial W_{i,j}^h}$$

Output bias:

$$\frac{\partial \ell(\mathbf{f}(\mathbf{x}), \mathbf{y})}{\partial b_i^o}$$

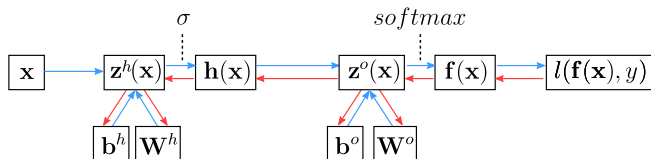
Hidden bias:

$$\frac{\partial \ell(\mathbf{f}(\mathbf{x}), \mathbf{y})}{\partial b_i^h}$$

→ differentiate with *chain-rule* – with  $z(x) = u(v(x))$ :

$$\frac{\partial z}{\partial x} = \left[ \frac{\partial u}{\partial v} \right]^\top \frac{\partial v}{\partial x}$$

# Backpropagation: For one hidden-layer network



Compute activation gradients

$$\nabla_{\mathbf{z}^o(\mathbf{x})} \ell = \mathbf{f}(\mathbf{x}) - \text{one-hot-encode}(y)$$

Compute layer params gradients

$$\nabla_{\mathbf{W}^o} \ell = \nabla_{\mathbf{z}^o(\mathbf{x})} \ell \cdot \mathbf{h}(\mathbf{x})^\top$$

$$\nabla_{\mathbf{b}^o} \ell = \nabla_{\mathbf{z}^o(\mathbf{x})} \ell$$

Compute prev layer activation gradients

$$\nabla_{\mathbf{h}(\mathbf{x})} \ell = \mathbf{W}^{o\top} \nabla_{\mathbf{z}^o(\mathbf{x})} \ell$$

$$\nabla_{\mathbf{z}^h(\mathbf{x})} \ell = \nabla_{\mathbf{h}(\mathbf{x})} \ell \odot \sigma'(\mathbf{z}^h(\mathbf{x}))$$

# Outline

Some Applications of Generative Adversarial Networks

Brief Intros on Neural Network

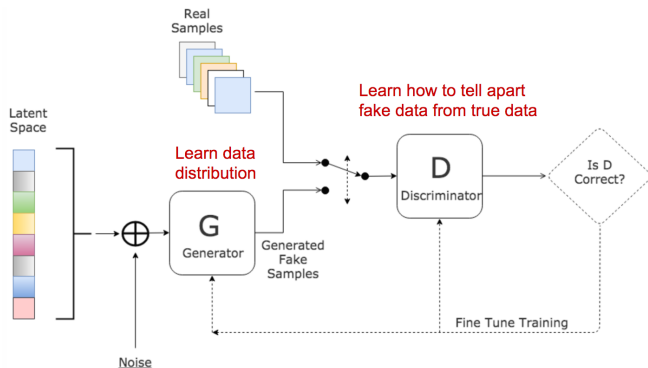
**Generative Model: Generative Adversarial Network**

Problems with training GANs

Techniques to improve GANs Training

Wasserstein GAN

# Generative Adversarial Network (GAN)



Alternate training between:

- ▶ A *generative network*  $G$  (generator), and
- ▶ A *discriminative network*  $D$  (discriminator)

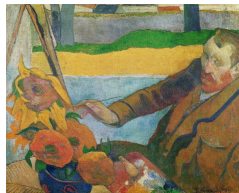
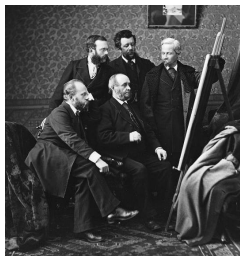
Goodfellow, Ian, et al. Generative Adversarial Nets. NeurIPS 2014



## What does these two networks do?

(intuitively): Discriminator

- ▶ Estimates the probability of a given sample coming from the real dataset
- ▶ Optimized to tell the fake samples from the real ones



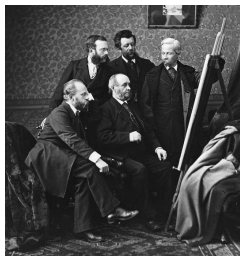
## What do these two networks do?

(intuitively): Discriminator

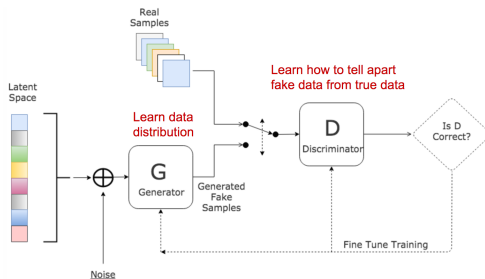
- ▶ Estimates the probability of a given sample coming from the real dataset
- ▶ Optimized to tell the fake samples from the real ones

Generator

- ▶ outputs synthetic samples given a noise variable input ( brings in potential output diversity)
- ▶ trained to capture the real data distribution so to generate samples can be as real as possible



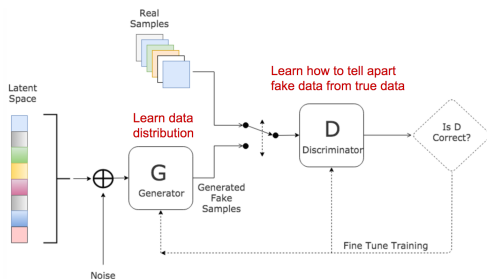
# What do these two networks do?



More formally:

- ▶  $p_z$ : Noise distribution
- ▶  $p_g$ : Generator's distribution over data  $x$
- ▶  $p_r$ : Real data  $x$ 's distribution

# What do these two networks do?



## Discriminator:

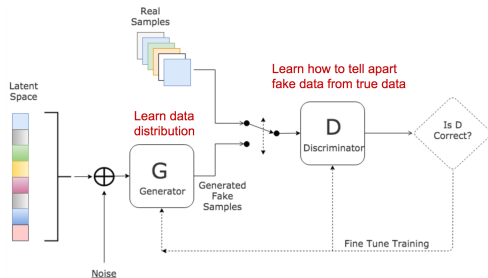
- ▶ Make sure decisions over real data are accurate by

$$\max \mathbb{E}_{x \sim p_r(x)} [\log D(x)]$$

- ▶ At the same time: given a fake sample  $G(z)$  with  $z \sim p_z$

$$\max \mathbb{E}_{z \sim p_z} [\log(1 - D(G(z)))]$$

# What do these two networks do?

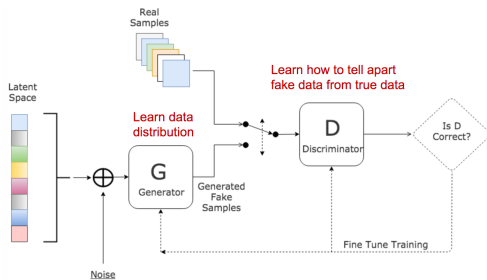


## Generator:

- ▶ Try its best to fool the discriminator:

$$\min \mathbb{E}_{z \sim p_z} [\log(1 - D(G(z)))]$$

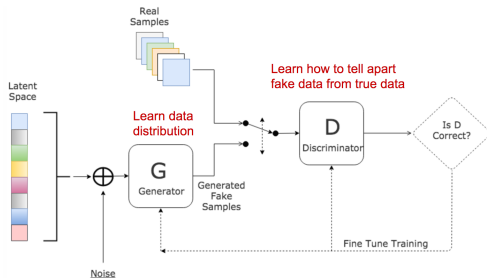
# What do these two networks do?



Putting together:

$$\begin{aligned} L(D, G) &= \mathbb{E}_{x \sim p_r(x)} [\log D(x)] + \mathbb{E}_{z \sim p_z(x)} [\log(1 - D(G(z)))] \\ &= \mathbb{E}_{x \sim p_r(x)} [\log D(x)] + \mathbb{E}_{x \sim p_g(x)} [\log(1 - D(x))] \end{aligned}$$

# What do these two networks do?



Putting together: playing a *minimax* (zero-sum) game

$$\min_G \max_D L(D, G) = \mathbb{E}_{x \sim p_r(x)} [\log D(x)] + \mathbb{E}_{x \sim p_g(x)} [\log(1 - D(x))]$$

Let's examine the loss function

$$L(D, G) = \mathbb{E}_{\mathbf{x} \sim p_r(\mathbf{x})}[\log D(\mathbf{x})] + \mathbb{E}_{\mathbf{x} \sim p_g(\mathbf{x})}[\log(1 - D(\mathbf{x}))]$$



Let's examine the loss function

$$\begin{aligned}L(D, G) &= \mathbb{E}_{x \sim p_r(x)}[\log D(\mathbf{x})] + \mathbb{E}_{x \sim p_g(x)}[\log(1 - D(\mathbf{x}))] \\ &= \int_x (p_r(x) \log D(x) + p_g(x) \log(1 - D(x))) dx\end{aligned}$$

## Let's examine the loss function

$$\begin{aligned}L(D, G) &= \mathbb{E}_{x \sim p_r(x)}[\log D(x)] + \mathbb{E}_{x \sim p_g(x)}[\log(1 - D(x))] \\ &= \int_x (p_r(x) \log D(x) + p_g(x) \log(1 - D(x))) dx\end{aligned}$$

Since we are summing all over sample  $x$ , what we should pay attention to is

$$\min_G \max_D f(x) \stackrel{\text{def.}}{=} p_r(x) \log D(x) + p_g(x) \log(1 - D(x))$$

Question: finding optimal  $D$  is equivalent with?

## Let's examine the loss function

$$\begin{aligned}L(D, G) &= \mathbb{E}_{x \sim p_r(x)}[\log D(x)] + \mathbb{E}_{x \sim p_g(x)}[\log(1 - D(x))] \\ &= \int_x (p_r(x) \log D(x) + p_g(x) \log(1 - D(x))) dx\end{aligned}$$

Since we are summing all over sample  $x$ , what we should pay attention to is

$$\min_G \max_D f(x) \stackrel{\text{def.}}{=} p_r(x) \log D(x) + p_g(x) \log(1 - D(x))$$

Question: finding optimal  $D$  is equivalent with?

$$\frac{\partial f(x)}{\partial D(x)} = 0$$

Optimal value for  $D$

$$\frac{\partial f(\mathbf{x})}{\partial D(\mathbf{x})} =$$

## Optimal value for $D$

$$\frac{\partial f(\mathbf{x})}{\partial D(\mathbf{x})} = p_r(\mathbf{x}) \frac{1}{\ln 10} \frac{1}{D(\mathbf{x})} - p_g(\mathbf{x}) \frac{1}{\ln 10} \frac{1}{1 - D(\mathbf{x})}$$

## Optimal value for $D$

$$\begin{aligned}\frac{\partial f(\mathbf{x})}{\partial D(\mathbf{x})} &= p_r(\mathbf{x}) \frac{1}{\ln 10} \frac{1}{D(\mathbf{x})} - p_g(\mathbf{x}) \frac{1}{\ln 10} \frac{1}{1 - D(\mathbf{x})} \\ &= \frac{1}{\ln 10} \left( \frac{p_r(\mathbf{x})}{D(\mathbf{x})} - \frac{p_g(\mathbf{x})}{1 - D(\mathbf{x})} \right)\end{aligned}$$

## Optimal value for $D$

$$\begin{aligned}\frac{\partial f(\mathbf{x})}{\partial D(\mathbf{x})} &= p_r(\mathbf{x}) \frac{1}{\ln 10} \frac{1}{D(\mathbf{x})} - p_g(\mathbf{x}) \frac{1}{\ln 10} \frac{1}{1 - D(\mathbf{x})} \\ &= \frac{1}{\ln 10} \left( \frac{p_r(\mathbf{x})}{D(\mathbf{x})} - \frac{p_g(\mathbf{x})}{1 - D(\mathbf{x})} \right) \\ &= \frac{1}{\ln 10} \left( \frac{p_r(\mathbf{x}) - (p_r(\mathbf{x}) + p_g(\mathbf{x}))D(\mathbf{x})}{D(\mathbf{x})(1 - D(\mathbf{x}))} \right)\end{aligned}$$

## Optimal value for $D$

$$\begin{aligned}\frac{\partial f(\mathbf{x})}{\partial D(\mathbf{x})} &= p_r(\mathbf{x}) \frac{1}{\ln 10} \frac{1}{D(\mathbf{x})} - p_g(\mathbf{x}) \frac{1}{\ln 10} \frac{1}{1 - D(\mathbf{x})} \\ &= \frac{1}{\ln 10} \left( \frac{p_r(\mathbf{x})}{D(\mathbf{x})} - \frac{p_g(\mathbf{x})}{1 - D(\mathbf{x})} \right) \\ &= \frac{1}{\ln 10} \left( \frac{p_r(\mathbf{x}) - (p_r(\mathbf{x}) + p_g(\mathbf{x}))D(\mathbf{x})}{D(\mathbf{x})(1 - D(\mathbf{x}))} \right)\end{aligned}$$

and so setting this equals to zero gives

$$D^*(\mathbf{x}) = \frac{p_r(\mathbf{x})}{p_r(\mathbf{x}) + p_g(\mathbf{x})}$$



# Global Optimal

$$D^*(\mathbf{x}) = \frac{p_r(\mathbf{x})}{p_r(\mathbf{x}) + p_g(\mathbf{x})}$$

plugging in the loss we have

$$\begin{aligned} L(G, D^*) &= \mathbb{E}_{\mathbf{x} \sim p_r(\mathbf{x})} [\log D^*(\mathbf{x})] + \mathbb{E}_{\mathbf{x} \sim p_g(\mathbf{x})} [\log(1 - D^*(\mathbf{x}))] \\ &= \mathbb{E}_{\mathbf{x} \sim p_r(\mathbf{x})} \left[ \log \frac{p_r(\mathbf{x})}{p_r(\mathbf{x}) + p_g(\mathbf{x})} \right] + \mathbb{E}_{\mathbf{x} \sim p_g(\mathbf{x})} \left[ \frac{p_g(\mathbf{x})}{p_r(\mathbf{x}) + p_g(\mathbf{x})} \right] \end{aligned}$$

# Global Optimal

$$D^*(\mathbf{x}) = \frac{p_r(\mathbf{x})}{p_r(\mathbf{x}) + p_g(\mathbf{x})}$$

plugging in the loss we have

$$\begin{aligned} L(G, D^*) &= \mathbb{E}_{x \sim p_r(x)}[\log D^*(\mathbf{x})] + \mathbb{E}_{x \sim p_g(x)}[\log(1 - D^*(\mathbf{x}))] \\ &= \mathbb{E}_{x \sim p_r(x)} \left[ \log \frac{p_r(\mathbf{x})}{p_r(\mathbf{x}) + p_g(\mathbf{x})} \right] + \mathbb{E}_{x \sim p_g(x)} \left[ \frac{p_g(\mathbf{x})}{p_r(\mathbf{x}) + p_g(\mathbf{x})} \right] \\ &= -2 \log(2) + \mathbb{E}_{x \sim p_r(x)} \left[ \log \frac{2p_r(\mathbf{x})}{p_r(\mathbf{x}) + p_g(\mathbf{x})} \right] + \mathbb{E}_{x \sim p_g(x)} \left[ \frac{2p_g(\mathbf{x})}{p_r(\mathbf{x}) + p_g(\mathbf{x})} \right] \end{aligned}$$

# Global Optimal

$$D^*(\mathbf{x}) = \frac{p_r(\mathbf{x})}{p_r(\mathbf{x}) + p_g(\mathbf{x})}$$

plugging in the loss we have

$$\begin{aligned} L(G, D^*) &= \mathbb{E}_{\mathbf{x} \sim p_r(\mathbf{x})} [\log D^*(\mathbf{x})] + \mathbb{E}_{\mathbf{x} \sim p_g(\mathbf{x})} [\log(1 - D^*(\mathbf{x}))] \\ &= \mathbb{E}_{\mathbf{x} \sim p_r(\mathbf{x})} \left[ \log \frac{p_r(\mathbf{x})}{p_r(\mathbf{x}) + p_g(\mathbf{x})} \right] + \mathbb{E}_{\mathbf{x} \sim p_g(\mathbf{x})} \left[ \frac{p_g(\mathbf{x})}{p_r(\mathbf{x}) + p_g(\mathbf{x})} \right] \\ &= -2 \log(2) + \mathbb{E}_{\mathbf{x} \sim p_r(\mathbf{x})} \left[ \log \frac{2p_r(\mathbf{x})}{p_r(\mathbf{x}) + p_g(\mathbf{x})} \right] + \mathbb{E}_{\mathbf{x} \sim p_g(\mathbf{x})} \left[ \frac{2p_g(\mathbf{x})}{p_r(\mathbf{x}) + p_g(\mathbf{x})} \right] \\ &= -\log(4) + \text{KL}(p_r \parallel (p_r(\mathbf{x}) + p_g(\mathbf{x}))/2) + \text{KL}(p_g \parallel (p_r(\mathbf{x}) + p_g(\mathbf{x}))/2) \end{aligned}$$

## Digression: Statistical Divergence

- ▶ A statistical divergence measures how a probability distribution  $P$  differs from another probability distribution  $Q$
- ▶ Kullback-Leibler divergence:

$$\text{KL}(P \parallel Q) \stackrel{\text{def.}}{=} E_P \left[ \log \left( \frac{P}{Q} \right) \right]$$

- ▶ Jensen-Shannon divergence:

$$\text{JS}(P \parallel Q) \stackrel{\text{def.}}{=} \frac{1}{2} \text{KL}(P \parallel (P + Q)/2) + \frac{1}{2} \text{KL}(Q \parallel (P + Q)/2)$$

→ goes to 0 when  $P$  and  $Q$  are the same

## Back to Global Optimal...

$$D^*(x) = \frac{p_r(x)}{p_r(x) + p_g(x)}$$

plugging in the loss we have

$$\begin{aligned} L(G, D^*) &= -\log(4) + \text{KL}(p_r \parallel (p_r(x) + p_g(x))/2) + \text{KL}(p_g \parallel (p_r(x) + p_g(x))/2) \\ &= -\log(4) + \text{JS}(p_r \parallel p_g) \end{aligned}$$

→ when  $p_g$  is trained to be very close to  $p_r$ , we achieve the theoretical global optimum  $L(G^*, D^*) = -\log(4)$

# Training GAN

---

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

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

**end for**

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

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

**end for**

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

---

# Training GAN

**Proposition 2.** *If  $G$  and  $D$  have enough capacity, and at each step of Algorithm 1, the discriminator is allowed to reach its optimum given  $G$ , and  $p_g$  is updated so as to improve the criterion*

$$\mathbb{E}_{\mathbf{x} \sim p_{data}} [\log D_G^*(\mathbf{x})] + \mathbb{E}_{\mathbf{x} \sim p_g} [\log(1 - D_G^*(\mathbf{x}))]$$

*then  $p_g$  converges to  $p_{data}$*

# Training GAN

**Proposition 2.** *If  $G$  and  $D$  have enough capacity, and at each step of Algorithm 1, the discriminator is allowed to reach its optimum given  $G$ , and  $p_g$  is updated so as to improve the criterion*

$$\mathbb{E}_{\mathbf{x} \sim p_{data}} [\log D_G^*(\mathbf{x})] + \mathbb{E}_{\mathbf{x} \sim p_g} [\log(1 - D_G^*(\mathbf{x}))]$$

*then  $p_g$  converges to  $p_{data}$*

→ but reality is far from theory...



# Outline

Some Applications of Generative Adversarial Networks

Brief Intros on Neural Network

Generative Model: Generative Adversarial Network

**Problems with training GANs**

Techniques to improve GANs Training

Wasserstein GAN

# Hard to achieve a global optimum (Nash Equilibrium)

**Proposition 2.** *If  $G$  and  $D$  have enough capacity, and at each step of Algorithm 1, the discriminator is allowed to reach its optimum given  $G$ , and  $p_g$  is updated so as to improve the criterion*

$$\mathbb{E}_{\mathbf{x} \sim p_{data}} [\log D_G^*(\mathbf{x})] + \mathbb{E}_{\mathbf{x} \sim p_g} [\log(1 - D_G^*(\mathbf{x}))]$$

*then  $p_g$  converges to  $p_{data}$*

Simple example:

$$\min_x \max_y f(x, y) = xy$$

# Hard to achieve a global optimum (Nash Equilibrium)

**Proposition 2.** *If  $G$  and  $D$  have enough capacity, and at each step of Algorithm 1, the discriminator is allowed to reach its optimum given  $G$ , and  $p_g$  is updated so as to improve the criterion*

$$\mathbb{E}_{\mathbf{x} \sim p_{data}} [\log D_G^*(\mathbf{x})] + \mathbb{E}_{\mathbf{x} \sim p_g} [\log(1 - D_G^*(\mathbf{x}))]$$

*then  $p_g$  converges to  $p_{data}$*

Simple example:

$$\min_x \max_y f(x, y) = xy$$

$$\text{equivalent with } \min_x xy \quad \text{and} \quad \min_y -xy$$

# Hard to achieve a global optimum (Nash Equilibrium)

Simple example:

$$\min_x \max_y f(x, y) = xy$$

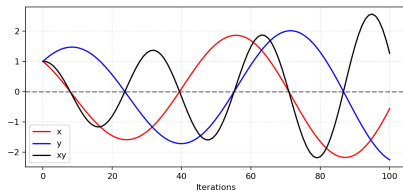
equivalent with  $\min_x xy$  and  $\min_y -xy$

Differentiating gives:

$$\frac{\partial f}{\partial x} = y \quad \frac{\partial f}{\partial y} = -x$$

So when doing GD:

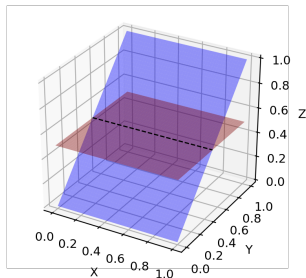
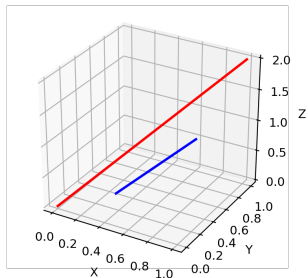
$$x^{(k+1)} = x^{(k)} - \eta y \quad y^{(k+1)} = y^{(k)} + \eta x$$



---

Salimans, Tim, et al. "Improved techniques for training gans." NeurIPS 2016

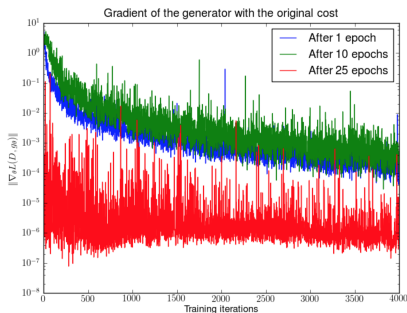
# Supports of $p_r$ and $p_g$ concentrates in lower dimension



Meaning: for high-dimensional data,  $p_r$  and  $p_g$  can be very hard to train to match each others

→ our discriminator has easy time to distinguish real and fake samples

# Vanishing Gradient



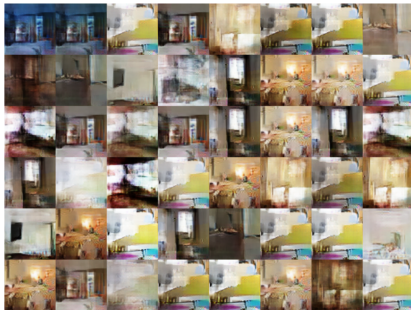
Recall that

$$\min_G \max_D L(D, G) = \mathbb{E}_{x \sim p_r(x)} [\log D(x)] + \mathbb{E}_{x \sim p_g(x)} [\log(1 - D(x))]$$

→ if the discriminator does a great job ( $D(x) = 1$  for all  $x \sim p_r(x)$  and  $D(x) = 0$  for all  $x \sim p_g(x)$ ), the gradient of the loss function drops down to close to zero

→ SGD update no longer works

# Mode Collapse



Mode Collapse: the generator always produces same outputs in training

---

Arjovsky et. al. "Wasserstein generative adversarial networks." ICML 2017

# Lack of Proper Evaluation Metric

When training GANs:

- ▶ We can check whether if the loss decreases – but it does not tell the whole story (recall vanishing gradient slide)
- ▶ How to evaluate the generated images' quality? (asking human being?)
- ▶ *i.e.* lack an objective function to compare performance of different architectures



# Outline

Some Applications of Generative Adversarial Networks

Brief Intros on Neural Network

Generative Model: Generative Adversarial Network

Problems with training GANs

**Techniques to improve GANs Training**

Wasserstein GAN

## To read more

- ▶ This part is based on Section 3 of Salimans, Tim, et al. "Improved techniques for training gans." NeurIPS 2016
- ▶ Related: Arjovsky and Bottou. "Towards Principled Methods for Training Generative Adversarial Networks." ICLR 2017.
- ▶ Also related: Huszár "How (not) to train your generative model: Scheduled sampling, likelihood, adversary?." ICLR 2015.
- ▶ Arguably the original GAN proposal of Goodfellow et al. (2014) worked because they switched the loss from using asymmetric Kullback-Leibler divergence to symmetric Jensen-Shannon divergence

# Feature Matching

- ▶ Objective: lessen the stability of GANs training
- ▶ Idea: discriminator to inspect whether the generator's output matches expected statistics of the real samples
- ▶ Implementation: training the generator to match the expected value of the features on an intermediate layer of the discriminator
- ▶ Denoting  $f(x)$  the summary on an intermediate layer of the discriminator, now the objective is

$$\|\mathbb{E}_{x \sim p_r} f(x) - \mathbb{E}_{x \sim p_z} f(G(z))\|_2^2$$

- ▶ In practice:  $f(x)$  can be mean or median (summary statistics) of the features

# Minibatch Discrimination

- ▶ Objective: avoid mode-collapse in the generator
- ▶ Idea: the discriminator is able to account for the relationship between training data points in one batch, instead of each point independently
- ▶ Implementation: In one minibatch, we approximate the closeness between every pair of samples, denoting  $c_b(\mathbf{x}_i, \mathbf{x}_j)_{i,j \in [n]}$ , and get the overall summary of one data point by summing up how close it is to other samples in the same batch

$$o_b(\mathbf{x}_i) = \sum_j c_b(\mathbf{x}_i, \mathbf{x}_j)$$

- ▶ Then we add  $o_b(\mathbf{x}_i)$  to the input of the model

## Historical Averaging

- ▶ Objective: bypass the low-dimensional manifold support problem and saddle points for the non-convex minimax optimization
- ▶ Idea: adding to the generator and the discriminator's cost a term

$$\|\theta - \frac{1}{t} \sum_{i=1}^t \theta(i)\|_2^2$$

where  $\theta(i)$  is the parameters at the previous time  $i$

- ▶ Implementation: scale better (*i.e.* faster training for large dataset) by adding only  $\theta(i - 1)$  to the  $i$ -th training iteration (online learning)

# One-sided label smoothing

- ▶ Objective: soften the discriminator value to avoid network's vulnerability
- ▶ Idea: smoothing label of discriminator network with  $\alpha$  for positive label (or 1 if the input is real data), and  $\beta$  for negative label
- ▶ The optimal  $D$  is then

$$D(\mathbf{x}) = \frac{\alpha p_r(\mathbf{x}) + \beta p_g(\mathbf{x})}{p_r(\mathbf{x}) + p_g(\mathbf{x})}$$

- ▶ Example:  $\alpha = 0.9$  and  $\beta = 0.1$

# Virtual Batch Normalization

- ▶ Batch-norm<sup>1</sup> is one of the most popular tricks for fastening and stabilizing the training of deep neural networks
- ▶ The idea of batch-norm is very simple: normalized each layer's input (activation) with their mean and covariance
- ▶ However: batch-norm can make output of a neural network highly dependent on some other inputs in the same minibatch
- ▶ Idea: in virtual batch-norm, we normalize the input based on a collection of fixed inputs (called reference batch)
- ▶ Implementation: virtual batch-norm is more expensive (forward pass on two minibatches of data) – so only used in the generator network

---

<sup>1</sup>Ioffe and Szegedy . "Batch Normalization: Accelerating Deep Network Training by Reducing Internal Covariate Shift" (2015)

# Better metric accounting for Distribution Similarity

- ▶ We already argue KL is not the best divergence for GAN loss, but also the JS divergence: when the two distributions lies on low dimensional manifold, it is hard to match them with gradient descent
- ▶ This is why the concept of Wasserstein GAN was introduced <sup>2</sup>

---

<sup>2</sup>Arjovsky et. al. "Wasserstein generative adversarial networks." ICML 2017



# Outline

Some Applications of Generative Adversarial Networks

Brief Intros on Neural Network

Generative Model: Generative Adversarial Network

Problems with training GANs

Techniques to improve GANs Training

Wasserstein GAN

## Previously

$$\min_G \max_D L(D, G) = \mathbb{E}_{x \sim p_r(x)} [\log D(x)] + \mathbb{E}_{x \sim p_g(x)} [\log(1 - D(x))]$$

GANs training suffers from many problems:

- ▶ Stability of the optimization scheme
  - ▶ Vanishing gradient
  - ▶ Mode collapse
- with many training fixes

# Better metric accounting for Distribution Similarity

- ▶ We already argue KL is not the best divergence for GAN loss, but also the JS divergence: when the two distributions lies on low dimensional manifold, it is hard to match them with gradient descent
- ▶ This is why the concept of Wasserstein GAN was introduced <sup>3</sup>

---

## Wasserstein Generative Adversarial Networks

---

Martin Arjovsky<sup>1</sup> Soumith Chintala<sup>2</sup> Léon Bottou<sup>1,2</sup>

### Abstract

We introduce a new algorithm named WGAN, an alternative to traditional GAN training. In this new model, we show that we can improve the stability of learning, get rid of problems like mode collapse, and provide meaningful learning curves useful for debugging and hyperparameter searches. Furthermore, we show that the corresponding optimization problem is sound, and provide extensive theoretical work highlighting the deep connections to different distances between distributions.

The typical remedy is to add a noise term to the model distribution. This is why virtually all generative models described in the classical machine learning literature include a noise component. In the simplest case, one assumes a Gaussian noise with relatively high bandwidth in order to cover all the examples. It is well known, for instance, that in the case of image generation models, this noise degrades the quality of the samples and makes them blurry. For example, we can see in the recent paper (Wu et al., 2016) that the optimal standard deviation of the noise added to the model when maximizing likelihood is around 0.1 to each pixel in a generated image, when the pixels were already normalized to be in the range [0, 1]. This is a very

---

<sup>3</sup>Arjovsky et. al. "Wasserstein generative adversarial networks." ICML 2017

# Statistical Distance

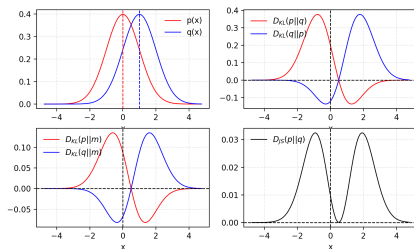
- ▶ Distance in  $\mathbb{R}^d$ ? What is the distance between two vectors? Two matrices?

# Statistical Distance

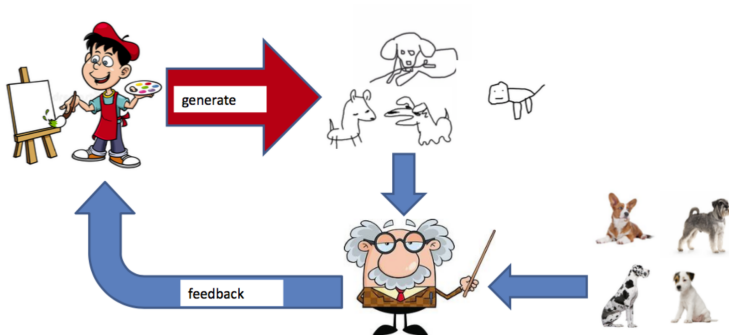
- ▶ Distance in  $\mathbb{R}^d$ ? What is the distance between two vectors? Two matrices?
- ▶ What if we want to measure the distance between two probability distributions  $P$  and  $Q$ ?

# Statistical Distance

- ▶ Distance in  $\mathbb{R}^d$ ? What is the distance between two vectors? Two matrices?
- ▶ What if we want to measure the distance between two probability distributions  $P$  and  $Q$ ?
- ▶ Last week: Kullback-Leibler and Jensen-Shannon divergences



# Reminder: Why do we care about the prob. distance?



In GAN architecture, we have  $p_r(x)$  and  $p_g(x)$   
where  $p_r(x)$  the distribution of real dataset  
 $p_g(x)$  the generator learned distribution over data  $x$

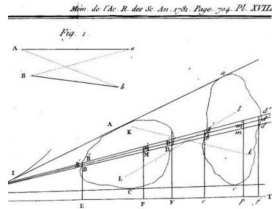
# Background on Optimal Transport – next topic

(1784)

*M É M O I R E*  
*S U R L A*  
*T H É O R I E D E S D É B L A I S*  
*E T D E S R E M B L A I S.*  
Par M. M O N G E.

Lorsqu'on doit transporter des terres d'un lieu dans un autre, on a coutume de donner le nom de *Déblai* au volume des terres que l'on doit transporter, & le nom de *Remblai* à l'espace qu'elles doivent occuper après le transport.

Le prix du transport d'une molécule étant, toutes choses d'ailleurs égales, proportionnel à son poids & à l'espace qu'on lui fait parcourir, & par conséquent le prix du transport total devant être proportionnel à la somme des produits des molécules multipliées chacune par l'espace parcouru, il s'en suit que le déblai & le remblai étant donnés de figure & de position, il n'est pas indifférent que telle molécule du déblai soit transportée dans tel ou tel autre endroit du remblai, mais qu'il y a une certaine distribution à faire des molécules du premier dans le second, d'après laquelle la somme de ces produits sera la moindre possible, & le prix du transport total fera un *minimum*.





# Wasserstein-1 distance as a GAN metric

$$\mathcal{W}_1(p_r, p_g) \stackrel{\text{def.}}{=} \min_{\gamma \in \Pi(p_r, p_g)} \mathbb{E}_{(x, x') \sim \gamma} \|x - x'\|$$

- ▶ Also called Earth-Mover distance (EMD)
- ▶ Imagine the two distributions are a certain amount of piles of dirt over a region
- ▶  $\gamma$  is one way to transport one pile of dirt to the other
- ▶ Therefore we have the constraint

$$\sum_x \gamma(x, x') = p_g(x')$$
$$\sum_{x'} \gamma(x, x') = p_r(x)$$

- ▶ EMD is the minimum cost of moving them

# Wasserstein-1 distance as a GAN metric

$$W_1(p_r, p_g) \stackrel{\text{def.}}{=} \min_{\gamma \in \Pi(p_r, p_g)} \mathbb{E}_{(x, x') \sim \gamma} \|x - x'\|$$

→ empirical version

$$W_1(p_r, p_g) \stackrel{\text{def.}}{=} \min_{\gamma \in \Pi(p_r, p_g)} \sum_{x, x'} \gamma(x, x') \|x - x'\|$$

When using  $W_1$  as a metric, we have the WGAN

---

Arjovsky et. al. "Wasserstein generative adversarial networks." ICML 2017

## Wasserstein-1 vs JS/KL as a GAN loss?

It's all about the smoothness (*i.e.* differentiable). Example 1 in the paper

- ▶ Suppose  $z \sim \mathcal{U}[0, 1]$ ,
- ▶  $\mathbb{P}_0$  is the distribution of  $(0, z) \in \mathbb{R}^2$
- ▶  $\mathbb{P}_\theta$  this distribution of  $g_\theta(z) = (\theta, z)$  with  $\theta \in (0, 1)$  (can view this as our generator distribution)

---

Arjovsky et. al. "Wasserstein generative adversarial networks." ICML 2017

## Wasserstein-1 vs JS/KL as a GAN loss?

What happens when one calculate these divergence? When  $\theta \neq 0$

- ▶ Kullback-Leibler:

$$KL(\mathbb{P}_0, \mathbb{P}_\theta) = KL(\mathbb{P}_\theta, \mathbb{P}_0) = \sum_{x=0, y \sim \mathcal{U}[0,1]} p_x \log \left( \frac{1}{0} \right) = +\infty$$

- ▶ Jensen-Shannon:

$$JS(\mathbb{P}_0, \mathbb{P}_\theta) = 1/2 \left( \sum_{x=0, y \sim \mathcal{U}[0,1]} 1 \log \frac{1}{1/2} + \sum_{x=\theta, y \sim \mathcal{U}[0,1]} 1 \log \frac{1}{1/2} \right) = \log 2$$

- ▶ EMD:

$$W(\mathbb{P}_0, \mathbb{P}_\theta) = |\theta|$$

When  $\theta = 0 \rightarrow KL(\mathbb{P}_0, \mathbb{P}_\theta) = JS(\mathbb{P}_0, \mathbb{P}_\theta) = 0$  but

$$W(\mathbb{P}_0, \mathbb{P}_\theta) = 0 = |\theta|$$

## Wasserstein-1 vs JS/KL as a GAN loss?

Conclusion from the simple example:

- ▶ KL and JS divergences are not smooth (differentiable) everywhere
- ▶ Using Wasserstein-1 distance fix this issue

## Wasserstein-1 as a GAN loss?

$$\mathcal{W}_1(p_r, p_g) \stackrel{\text{def.}}{=} \min_{\gamma \in \Pi(p_r, p_g)} \mathbb{E}_{(x, x') \sim \gamma} \|x - x'\|$$

Problem: this formulation is highly intractable

→ Kantorovich-Rubinstein duality for W1 distance

$$\mathcal{W}_1(p_r, p_g) = \sup_{\|f\|_L \leq 1} \mathbb{E}_{x \sim p_r}[f(x)] - \mathbb{E}_{x \sim p_g}[f(x)]$$

---

Remark 6.5, Villani, Cédric. Optimal transport: old and new. Vol. 338. Berlin: springer, 2009.

## Wasserstein-1 as a GAN loss?

In other words, with the W1 distance, we move from the loss

$$\min_G \max_D L(D, G) = \mathbb{E}_{x \sim p_r(x)}[\log D(x)] + \mathbb{E}_{x \sim p_g(x)}[\log(1 - D(x))]$$

to

$$\max_w \mathbb{E}_{x \sim p_r} [f_w(x)] - \mathbb{E}_{z \sim p_z} [f_w(g_\theta(z))]$$

Important improvements:

- ▶ From minimax loss (hard to solve) to single maximization
- ▶  $f_w$  is a parameterization: now we can use any architectures of neural networks to train

# Wasserstein-1 as a GAN loss?

Theorem 3: with the loss

$$\mathcal{W}_1(p_r, p_g) = \max_w \mathbb{E}_{x \sim p_r} [f_w(x)] - \mathbb{E}_{z \sim p_z} [f_w(g_\theta(z))]$$

Under additional technical assumption (we skip it here, but it relates to local Lipschitz continuity), we have

$$\nabla_w \mathcal{W}_1(p_r, p_g) = -\mathbb{E}_{z \sim p_z} [\nabla_w f_w(g_\theta(z))]$$

---

Arjovsky et. al. "Wasserstein generative adversarial networks." ICML 2017



## Wasserstein-1 as a GAN loss?

$$\nabla_w \mathcal{W}_1(p_r, p_g) = -\mathbb{E}_{z \sim p_z} [\nabla_w f_w(g_\theta(z))]$$

Intuition:

- ▶ Now the “discriminator”  $D$  is not a direct critic of telling the fake samples apart from the real ones anymore.
- ▶ Instead, it is trained to learn  $f_w$  to help compute Wasserstein distance.
- ▶ Loss function decreases  $\longrightarrow$  Wasserstein-1 distance gets smaller
- ▶ Generator  $g_\theta$  output becomes closer to the real data distribution (or  $p_g$  becomes very similar to  $p_r$ )

# Training WGAN

---

**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
```

---

Arjovsky et. al. "Wasserstein generative adversarial networks." ICML 2017

# Fixing vanishing gradient in some settings

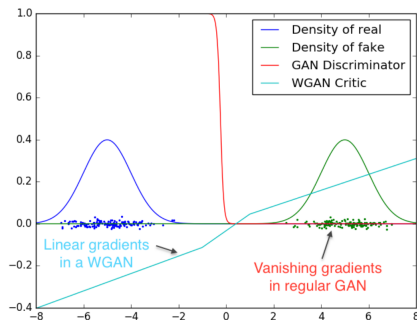


Figure 2: Optimal discriminator and critic when learning to differentiate two Gaussians. As we can see, the discriminator of a minimax GAN saturates and results in vanishing gradients. Our WGAN critic provides very clean gradients on all parts of the space.

# Better at image generation task



Figure 6: Algorithms trained with a DCGAN generator. Left: WGAN algorithm. Right: standard GAN formulation. Both algorithms produce high quality samples.



Figure 7: Algorithms trained with a generator without batch normalization and constant number of filters at every layer (as opposed to duplicating them every time as in (Radford et al., 2015)). Aside from taking out batch normalization, the number of parameters is therefore reduced by a bit more than an order of magnitude. Left: WGAN algorithm. Right: standard GAN formulation. As we can see the standard GAN failed to learn while the WGAN still was able to produce samples.



Figure 8: Algorithms trained with an MLP generator with 4 layers and 512 units with ReLU nonlinearities. The number of parameters is similar to that of a DCGAN, but it lacks a strong inductive bias for image generation. Left: WGAN algorithm. Right: standard GAN formulation. The WGAN method still was able to produce samples, lower quality than the DCGAN, and of higher quality than the MLP of the standard GAN. Note the significant degree of mode collapse in the GAN MLP.

## However...

---

**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:      $w \leftarrow w + \alpha \cdot \text{RMSProp}(w, g_w)$ 
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
```

---

From the paper: ‘Weight clipping is a clearly terrible way to enforce a Lipschitz constrain [...] we leave the topic of enforcing Lipschitz constraints in a NN setting for further investigation’

Arjovsky et. al. "Wasserstein generative adversarial networks." ICML 2017

# Difficulties with Weight Constraints

- ▶ Clipping the weights inside a range (*e.g.*  $(-0.01, 0.01)$ ) actually can lead to optimization difficulties

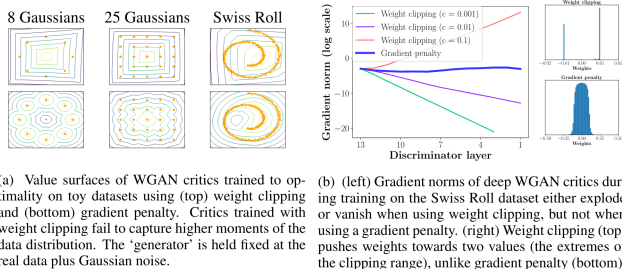


Figure 1: Gradient penalty in WGANs does not exhibit undesired behavior like weight clipping.

---

Gulrajani, Ishaan, et al. "Improved training of wasserstein gans." NeuRIPS 2017.

## WGAN-GP: a variant of GAN variant...

Introducing a regularization term that penalize the gradient to the loss

$$\mathcal{W}_1(p_r, p_g) = \max_w \mathbb{E}_{x \sim p_r} [f_w(x)] - \mathbb{E}_{z \sim p_z} [f_w(g_\theta(z))] \\ + \lambda \mathbb{E}_{z \sim p_z} [(\|\nabla_w f_w(g_\theta(z))\|_2 - 1)^2]$$

- ▶ This is equivalent to constraining the function  $f$  to be 1-Lipschitz
- ▶ Regularization terms is still differentiable

# Training WGAN-GP

---

**Algorithm 1** WGAN with gradient penalty. We use default values of  $\lambda = 10$ ,  $n_{\text{critic}} = 5$ ,  $\alpha = 0.0001$ ,  $\beta_1 = 0$ ,  $\beta_2 = 0.9$ .

---

**Require:** The gradient penalty coefficient  $\lambda$ , the number of critic iterations per generator iteration  $n_{\text{critic}}$ , the batch size  $m$ , Adam hyperparameters  $\alpha, \beta_1, \beta_2$ .

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

```
1: while  $\theta$  has not converged do
2:   for  $t = 1, \dots, n_{\text{critic}}$  do
3:     for  $i = 1, \dots, m$  do
4:       Sample real data  $\mathbf{x} \sim \mathbb{P}_r$ , latent variable  $\mathbf{z} \sim p(\mathbf{z})$ , a random number  $\epsilon \sim U[0, 1]$ .
5:        $\tilde{\mathbf{x}} \leftarrow G_{\theta}(\mathbf{z})$ 
6:        $\hat{\mathbf{x}} \leftarrow \epsilon \mathbf{x} + (1 - \epsilon) \tilde{\mathbf{x}}$ 
7:        $L^{(i)} \leftarrow D_w(\tilde{\mathbf{x}}) - D_w(\mathbf{x}) + \lambda(\|\nabla_{\hat{\mathbf{x}}} D_w(\hat{\mathbf{x}})\|_2 - 1)^2$ 
8:     end for
9:      $w \leftarrow \text{Adam}(\nabla_w \frac{1}{m} \sum_{i=1}^m L^{(i)}, w, \alpha, \beta_1, \beta_2)$ 
10:   end for
11:   Sample a batch of latent variables  $\{\mathbf{z}^{(i)}\}_{i=1}^m \sim p(\mathbf{z})$ .
12:    $\theta \leftarrow \text{Adam}(\nabla_{\theta} \frac{1}{m} \sum_{i=1}^m -D_w(G_{\theta}(\mathbf{z})), \theta, \alpha, \beta_1, \beta_2)$ 
13: end while
```

---

Gulrajani, Ishaan, et al. "Improved training of wasserstein gans." NeuRIPS 2017.



# WGAN-GP vs. WGAN vs. other GANs

DCGAN	LSGAN	WGAN (clipping)	WGAN-GP (ours)
Baseline ( $G$ : DCGAN, $D$ : DCGAN)			
			
$G$ : No BN and a constant number of filters, $D$ : DCGAN			
			
$G$ : 4-layer 512-dim ReLU MLP, $D$ : DCGAN			
			
No normalization in either $G$ or $D$			
			
Gated multiplicative nonlinearities everywhere in $G$ and $D$			
			
tanh nonlinearities everywhere in $G$ and $D$			
			
101-layer ResNet $G$ and $D$			
			

Figure 2: Different GAN architectures trained with different methods. We only succeeded in training every architecture with a shared set of hyperparameters using WGAN-GP.

# Conclusion

- ▶ Using Wasserstein-1 distance is better for GAN training (taken into account we are using the same architecture for generator), which leads to significantly improvement in image-generation task
- ▶ However, clipping weight as a way to enforce Lipschitz continuity is some kind of mysterious engineering technique
- ▶ With adding regularization term to enforce the norm of gradient to be equal to 1, WGAN-GP suffers less instability problem in training