

# Introducing Generative Adversarial Network (GAN)

Presentation on the module Advance Machine Learning

3048781      Sanchit Bhavsar  
3044093      Sthitadhee Panthadas

Course: Advance Machine Learning  
Prof. Dr. Ulf Brefeld

# Agenda

- ❖ **Understanding Generative Modeling**
- ❖ **Introducing GANs**
  - Model Architecture
  - Theoretical Foundation & Algorithm
  - Limitations of Vanilla GAN
  - Wasserstein GAN
- ❖ **Timeline of GAN based Architecture**
  - Conditional GAN (CGAN)
  - Deep-Convolutional GAN (DCGAN)
  - GAN with different loss function (Least Square GAN)
- ❖ **Training GAN on MNIST Data**
- ❖ **Conclusion and Future Outlook**



# Quick Test: Which face is real?



A



B



C



# Supervised vs Unsupervised Learning

Supervised Learning	Unsupervised Learning
<p><b>Data:</b> (x,y)</p> <p><math>x</math> is data, <math>y</math> is label</p> <p><b>Goal:</b> Learn function to map <math>x \rightarrow y</math></p> <p><b>Examples:</b> Classification, regression, Object detection etc.</p>	<p><b>Data:</b> (x)</p> <p><math>x</math> is data, <i><b>no labels!</b></i></p> <p><b>Goal:</b> Learn the hidden or underlying structure of data</p> <p><b>Examples:</b> Clustering, dimensionality reduction etc.</p>



# Generative Models

**Question:** can we build a model to approximate a data distribution?

- We are given and a finite sample from this distribution ,

$$X = \{x | x \sim p_{\text{data}}(x)\}, |X| = n$$

- **Goal:** Given training data, learn a model that represents new samples from same distribution

$$p_{\text{model}}(x; \theta) \approx p_{\text{data}}(x)$$



Training data  $\sim p_{\text{data}}(x)$



Generated samples  $\sim p_{\text{model}}(x)$

# Why Care About Generative Models?

Often overused quote:

“What I cannot create, I do not understand” -R. Feynman

Capable of uncovering **underlying features** in a dataset,

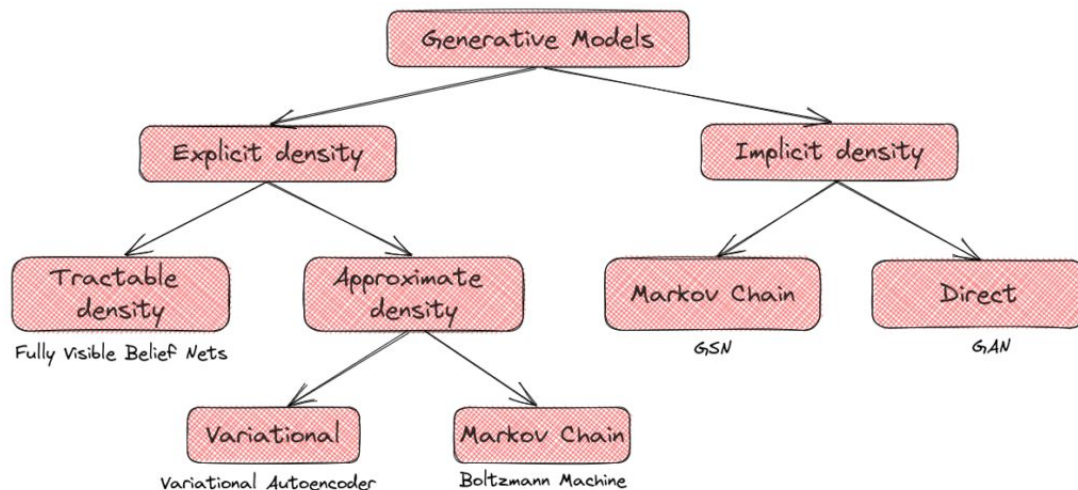
- ❖ Noisy Input
- ❖ Simulated Data
- ❖ Features Representation of Data
- ❖ Prediction of Future State
- ❖ Missing Data



source: StyleGAN

# Taxonomy of Generative Models

- ❖ Addresses density estimation, a core problem in unsupervised learning
  - **Explicit density estimation**: explicitly define and solve for  $p_{\text{model}}(x)$
  - **Implicit density estimation**: learn model that can sample from  $p_{\text{model}}(x)$  w/o explicitly defining it



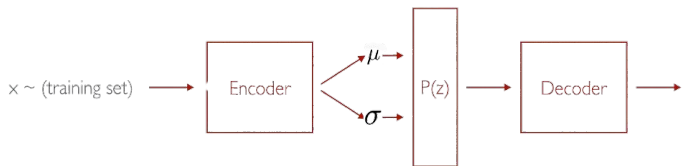
# WHAT ARE GANS?

## Generative Adversarial Networks

### Generative Models

We try to learn the underlying the distribution from which our dataset comes from.

Eg: Variational AutoEncoders (VAE)



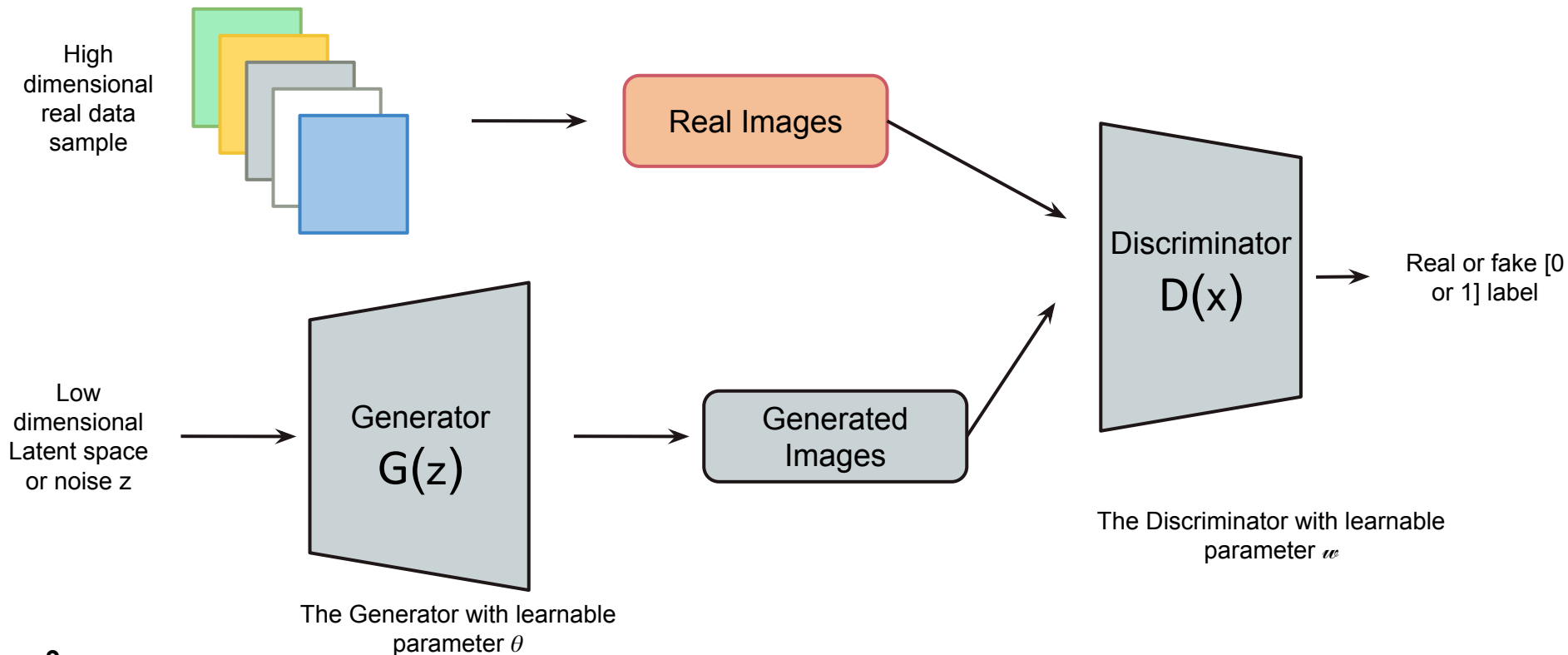
### Neural Networks

### Adversarial Training

GANS are made up of two competing networks (adversaries) that are trying beat each other.



# GAN Architecture: Two Player Game



# Objective Function



$$B.C.E = H(p, q) = -p \log q + (1 - p) \log(1 - q)$$

Notation	Description
$D(x)$	The discriminator network
$G(z)$	The generator network
$z$	STD Gaussian noise
$p$	Probability distribution or labels
$q$	Probability distribution or labels
$H(p, q)$	Binary Cross entropy

i

$$Err(1, D(x)) + Err(0, D(G(z)))$$

$$= -[1 \cdot \log D(x) + 0 \cdot \log(1 - D(x))] - [0 \cdot \log D(G(z)) + 1 \cdot \log(1 - D(G(z)))]$$

$$= -\log D(x) - \log(1 - D(G(z)))$$

$$p = 1 \\ q = D(.)$$

ii

$$\log D(x) + \log(1 - D(G(z)))$$

$$\log(1 - D(G(z)))$$

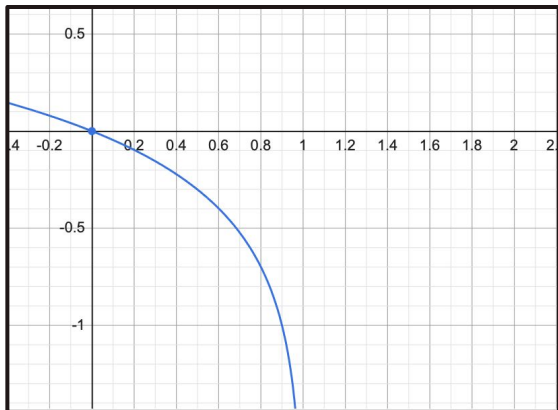
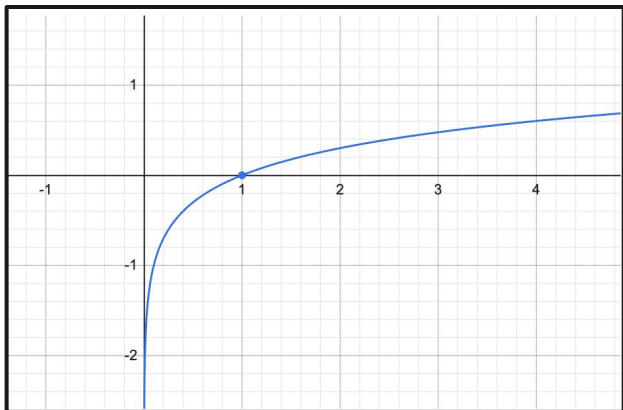
Discriminator loss → must maximise

Generator loss → must minimise

# Objective Function

Notation	Description
$w$	The discriminator parameters
$\theta$	The generator parameters
$E$	Expectation with respect to distribution
$x$	Real data sample

$$\min_{\theta} \max_w \mathbb{E}_{x \sim p_{data}} [\log D_w(x)] + \mathbb{E}_{z \sim p_{noise}} [\log(1 - D_w(G_{\theta}(z)))]$$



# GAN Algorithm

---

**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.

---

# Global Optimality

$$V(G, D) = \int_{\mathbf{x}} p_{\text{data}}(\mathbf{x}) \log(D(\mathbf{x})) d\mathbf{x} + \int_{\mathbf{z}} p_{\mathbf{z}}(\mathbf{z}) \log(1 - D(g(\mathbf{z}))) d\mathbf{z}$$

Training criterion for D  
(discriminator) at a fixed G

$$= \int_{\mathbf{x}} \underbrace{p_{\text{data}}(\mathbf{x})}_{a} \log(\underbrace{D(\mathbf{x})}_{y}) + \underbrace{p_g(\mathbf{x})}_{b} \log(1 - \underbrace{D(\mathbf{x})}_{y}) d\mathbf{x}$$

Can be proven in 2 ways

$$y \rightarrow a \log(y) + b \log(1 - y) \xrightarrow{\text{differentiate}} \frac{a}{a+b}$$

Note  $(a, b) \in \mathbb{R}^2 \setminus \{0, 0\}$  and  $y$  achieves its maximum in the range  $[0, 1]$

$$D_G^*(\mathbf{x}) = \frac{p_{\text{data}}(\mathbf{x})}{p_{\text{data}}(\mathbf{x}) + p_g(\mathbf{x})} = 0.5$$

- Global Optimality is reached at  $p_g = p_{\text{data}}$
- At that point  $D_G^*(\mathbf{x}) = 0.5$

Notation	Description
$p_{\text{data}}$	Prob Distribution of real data space
$p_g$	Prob Distribution of generated data
$D^*(\mathbf{x})$	Optimal discriminator at fixed G
$p_z$	Prob distribution of std normal
$V(G, D)$	General Loss function of GAN

# Connection to Jensen-Shannon Divergence

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

$$= \mathbb{E}_{\mathbf{x} \sim p_{\text{data}}} [\log D_G^*(\mathbf{x})] + \mathbb{E}_{\mathbf{z} \sim p_z} [\log(1 - D_G^*(G(\mathbf{z})))]$$

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

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

Can be proved via  $\rightarrow$

- Change of Variables
- Radon-Nikodym theorem (in Measure Theory)

Notation	Description
$D^*(\mathbf{x})$	Optimal discriminator at fixed $G$
$C(G)$	Training Criterion at fixed $G$
$V(G, D)$	General Loss function of GAN

- Global Optimality is reached at  $p_g = p_{\text{data}}$
- At that point  $D_G^*(\mathbf{x}) = 0.5$

# Connection to Jensen-Shannon Divergence

- Global Optimality is reached at  $p_g = p_{data}$
- At that point  $D_G^*(x) = 0.5$



$$\mathbb{E}_{x \sim p(x)}[f(x)] = \sum p(x)f(x)$$

$$\begin{aligned}
 C(G) &= \mathbb{E}_{x \sim p_{data}} \left[ \log \frac{p_{data}}{p_{data} + p_G} \right] + \mathbb{E}_{x \sim p_G} \left[ \log \frac{p_G}{p_{data} + p_G} \right] \\
 &= \mathbb{E}_{x \sim p_{data}} \left[ \log \frac{\frac{p_{data}}{2}}{\frac{p_{data} + p_G}{2}} \right] + \mathbb{E}_{x \sim p_G} \left[ \log \frac{\frac{p_G}{2}}{\frac{p_{data} + p_G}{2}} \right] \\
 &= -\log 4 + \mathbb{E}_{x \sim p_{data}} \left[ \log \frac{p_{data}}{\frac{p_{data} + p_G}{2}} \right] + \mathbb{E}_{x \sim p_G} \left[ \log \frac{p_G}{\frac{p_{data} + p_G}{2}} \right] \\
 &= -\log 4 + KL \left( p_{data} \parallel \frac{p_{data} + p_G}{2} \right) + KL \left( p_G \parallel \frac{p_{data} + p_G}{2} \right) \\
 &= -\log 4 + 2JS(p_{data} \parallel p_G)
 \end{aligned}$$

$p_g = p_{data}$  then  $D_G^*(x) = 0.5$

→  $-\log 4$  or  $-2 \log 2$


Notation	Description
$C(G)$	Virtual Training Criterion for fixed G.
$p_G$	Prob Distribution of generated data
$KL$	Kullback-Liebler divergence
$JS$	Jensen Shannon Divergence
$D_G^*(x)$	Optimal discriminator
$f(x)$	A general function, here $\log(\cdot)$

# Convergence of Algorithm

The argument follows that  $\rightarrow$

- ★ Due to linearity of expectations
- ★ Due to each term being convex functions in  $p_g$

Intuitive explanation of the terms + the proof and concept of sampling from  $p_g$



Sufficient small Gradient descent updates for  $p_g$  having a unique global optima  $D^*$ , allows  $p_g$  to converge to  $p_{\text{data}}$ .

**In practice,**

- ❖  $p_g \rightarrow$  **neural network** of the generating resulting in generated samples,
- ❖ we optimize not  $p_g$  but the parameters of the generator neural network  $\theta_g$ .



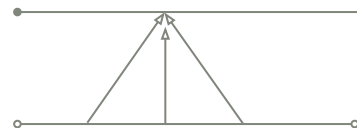
# Limitations of Vanilla GAN

In practice,

- ❖ Models do not converge to a local minima
- ❖ Mode collapse problem

Real data space

Generated data space



Because of,

- ❖ Most real life cases functions are necessarily convex and continuous and therefore not differentiable.
- ❖ Vanishing or exploding gradients

# Motivating Wasserstein

$$\delta(\mathbb{P}_r, \mathbb{P}_g) = \sup_{A \in \Sigma} |\mathbb{P}_r(A) - \mathbb{P}_g(A)| .$$

$$KL(\mathbb{P}_r \parallel \mathbb{P}_g) = \int \log \left( \frac{P_r(x)}{P_g(x)} \right) P_r(x) .$$

$$JS(\mathbb{P}_r, \mathbb{P}_g) = KL(\mathbb{P}_r \parallel \mathbb{P}_m) + KL(\mathbb{P}_g \parallel \mathbb{P}_m)$$

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

Notation	Description
$A$	Set with which you calculate the prob example $P_r(A)$
$\Sigma$	Total countable set space
$P_r$	Prob distribution for real data samples
$P_g$	Probability distribution representing generated data
$\gamma$	Scalar transport mass values/energies
$\Pi(P_r, P_g)$	Joint distribution of the transport mass value
$\inf$	Infimum
$\  \cdot \ $	Norm
$\sup$	Supremum
$\delta ( P_r, P_g )$	Total Variation (TV) distance
$W ( P_r, P_g )$	1-Wasserstein or Earth Mover distance

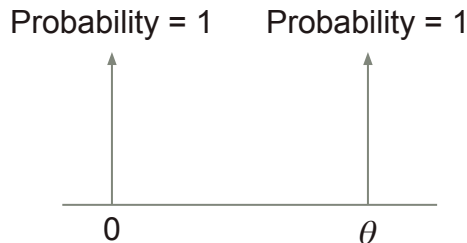
# Motivating Wasserstein

$$\delta(\mathbb{P}_r, \mathbb{P}_g) = \sup_{A \in \Sigma} |\mathbb{P}_r(A) - \mathbb{P}_g(A)| .$$

$$KL(\mathbb{P}_r \| \mathbb{P}_g) = \int \log \left( \frac{P_r(x)}{P_g(x)} \right) P_r(x) .$$

$$JS(\mathbb{P}_r, \mathbb{P}_g) = KL(\mathbb{P}_r \| \mathbb{P}_m) + KL(\mathbb{P}_g \| \mathbb{P}_m)$$

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



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

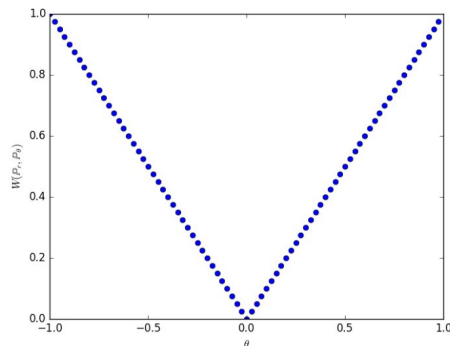
$$JS(\mathbb{P}_0, \mathbb{P}_\theta) = \begin{cases} \log 2 & \text{if } \theta \neq 0 , \\ 0 & \text{if } \theta = 0 , \end{cases}$$

$$KL(\mathbb{P}_\theta \| \mathbb{P}_0) = KL(\mathbb{P}_0 \| \mathbb{P}_\theta) = \begin{cases} +\infty & \text{if } \theta \neq 0 , \\ 0 & \text{if } \theta = 0 , \end{cases}$$

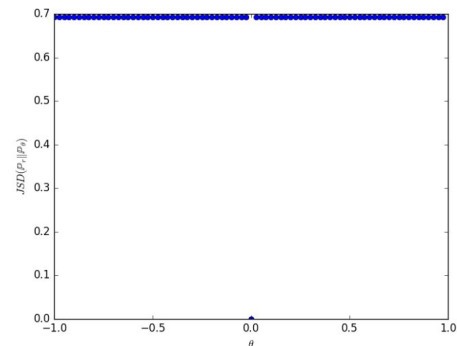
$$\text{and } \delta(\mathbb{P}_0, \mathbb{P}_\theta) = \begin{cases} 1 & \text{if } \theta \neq 0 , \\ 0 & \text{if } \theta = 0 . \end{cases}$$

# Motivating Wasserstein

- ❖ Wasserstein distance provides a **proportional value**, the further the distribution the greater the distance. In practical cases the generating sample distribution will **not overlap** the real data sample, wasserstein distance will be more practical than the other **metric measures** such as KL divergence or JSD or TV.
- ❖ Furthermore, Wasserstein distance or Earth mover distance is continuous and smoother function compared to the others.



EM distance



JSD

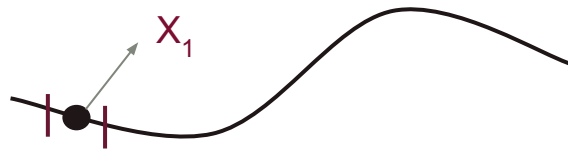
# Wasserstein and Continuity

What do you need for continuity for?

- **Smooth** gradient descent or ascent

What do we need for smooth gradient descent or ascent?

- Neural network  $G$  must be continuous in parameter  $\theta \rightarrow$  which is the case in EM distance
- $G$  must be locally Lipschitz



The slope of the curve at point  $X_n$  should be less than equal to  $K$ . ( $K$ -Lipschitz)

# Wasserstein Loss

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

- Note, for most practical probability distributions **EM distance calculation is intractable** because for that we would be **needing the joint probability  $\Pi(\mathbb{P}_r, \mathbb{P}_g)$** .
- So, we use **Kantorovich-Rubinstein duality** to get us a new formulae for EM distance that works.

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

- Note here the **supremum** is over **1-Lipschitz functions**
- But this can be easily changed to **k-Lipschitz functions** with

$$K \cdot W(\mathbb{P}_r, \mathbb{P}_g)$$

Notation	Description
$\sup$	Supremum
$f(x)$	K-Lipschitz function
$\ f\ _{L \leq 1}$	1-Lipschitz function(s)
$P_\theta$	Prob distribution representing generated data
$q$	Probability distribution
$K$	Constant value for Lipschitz constraint

# Wasserstein Loss

## Vanilla GAN Objective function

$$\min_{\theta} \max_w \mathbb{E}_{x \sim p_{data}} [\log D_w(x)] + \mathbb{E}_{z \sim p_{noise}} [\log(1 - D_w(G_{\theta}(z)))]$$

## Wasserstein GAN

$$\min_{\theta} \max_w V(G, C) = \min_{\theta} \max_w [ \mathbb{E}_{x \sim p_{data}} [C_w(x)] - \mathbb{E}_{z \sim p_{noise}} [(C_w(G_{\theta}(z)))] ]$$

Where  $C_w$  should follow the K-Lipschitz constraint.

# Weight Clipping

How do we ensure or **enforce Lipschitz constraint**?

- **Gradient/weight clipping** or gradient penalty

## ❖ **Solution:**

- **Weight clipping** is basically to cut off the value of the weights if it goes beyond a **threshold value** and then use the threshold value instead. This has a parameter that controls the amount of clipping.  $[-c \leq \text{weight} \leq c]$

## ❖ **Problem:**

- This poses another problem. This reduces the optimality of the critic.



# WGAN vs Vanilla GAN

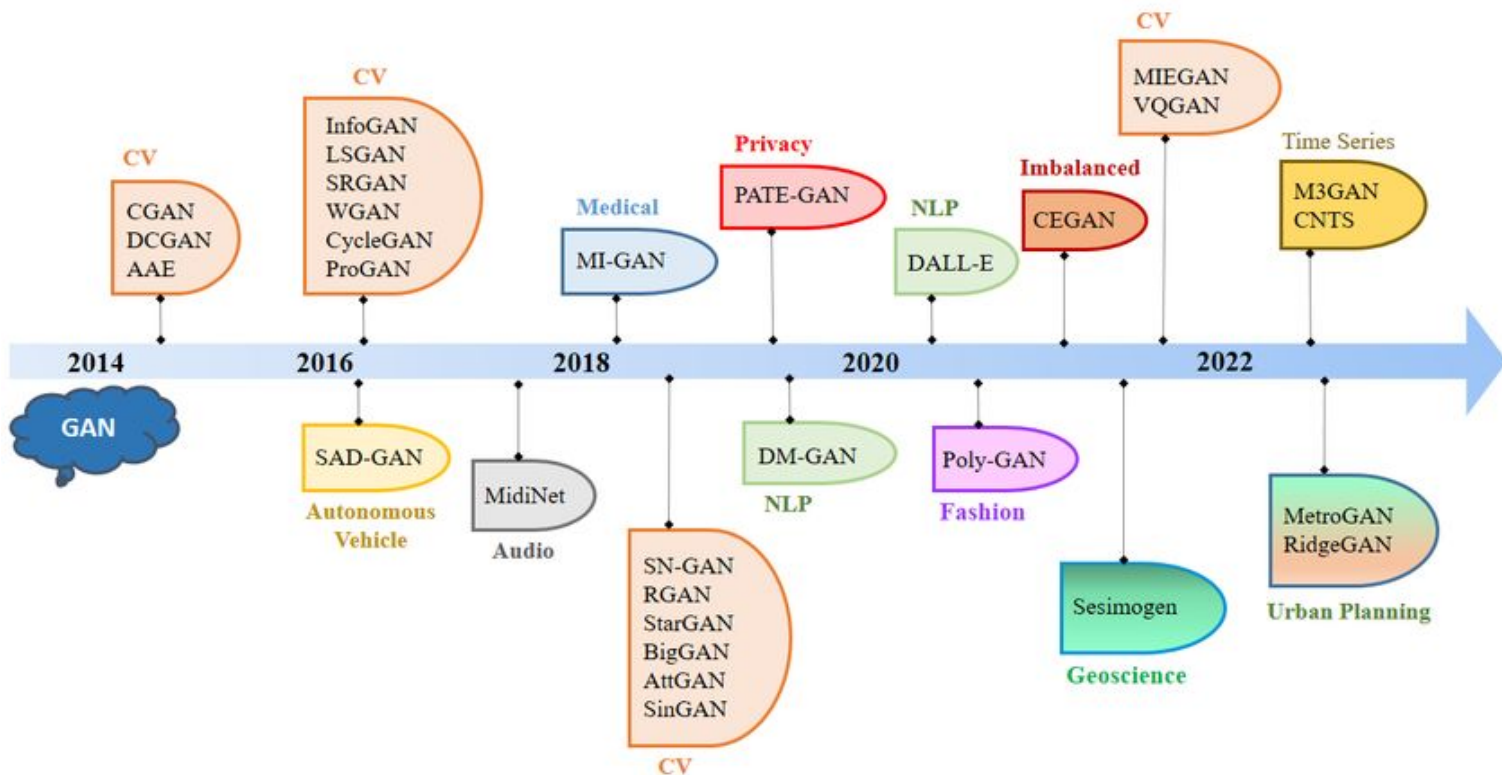
## WGAN

- Discriminator acts as a critic and gives out real values (regression)
- Wasserstein motivated loss
- Loss contains no log
- Discriminator is updated 5 times for each update in generator
- Root Mean Squared Propagation is used
- Includes clipping of parameter values

## Vanilla GAN

- Discriminator gives out 0 or 1 (classification)
- JSD motivated loss
- Loss contains log
- Discriminator is updated once for each update in generator
- Adam is used
- Includes no weight clipping

# Timeline of GAN based architecture

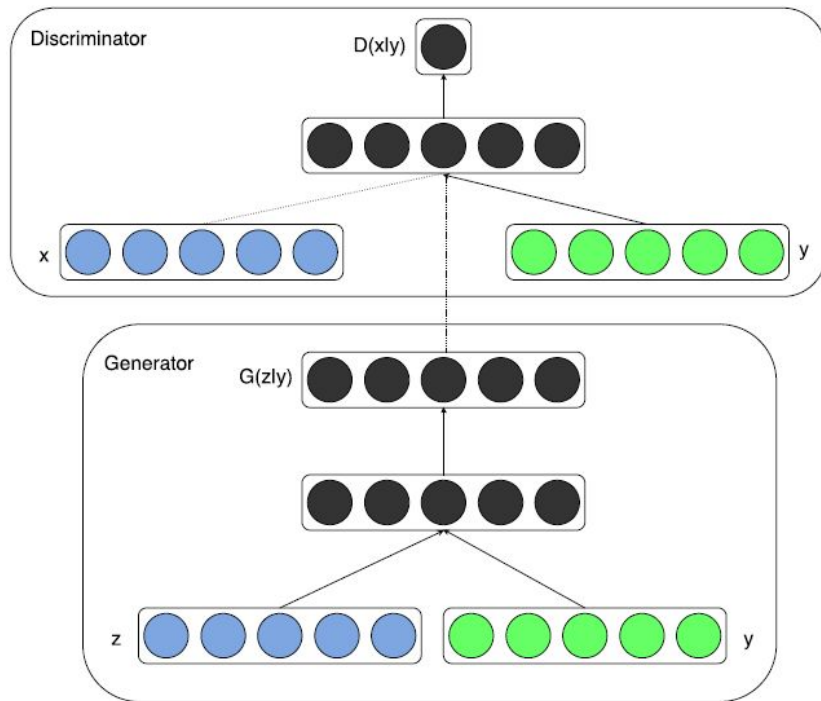


# Conditional GAN (CGAN)

- ❖ In an unconditioned generative model, there is no control on modes of the data being generated
- ❖ In CGAN, the generator learns to generate a fake sample with a specific condition (such as a label associated with an image or more detailed tag) rather than a generic sample from unknown noise distribution
- ❖ Simply feeding the data 'y' and conditioning both the generator and discriminator

Learning Loss of CGAN:

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

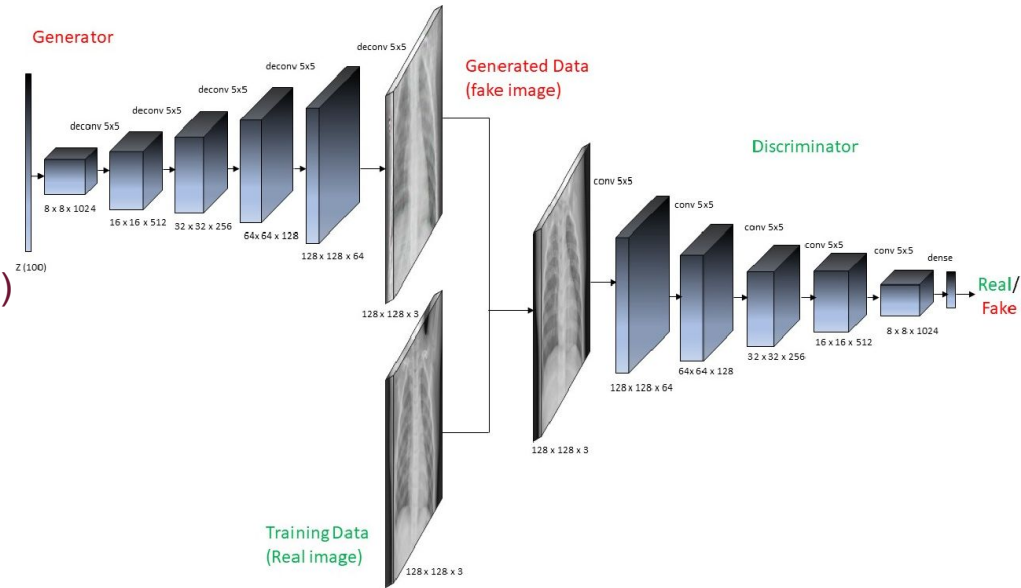


## Generated MNIST digits, each row conditioned on one label



# Deep Convolutional GAN (DCGAN)

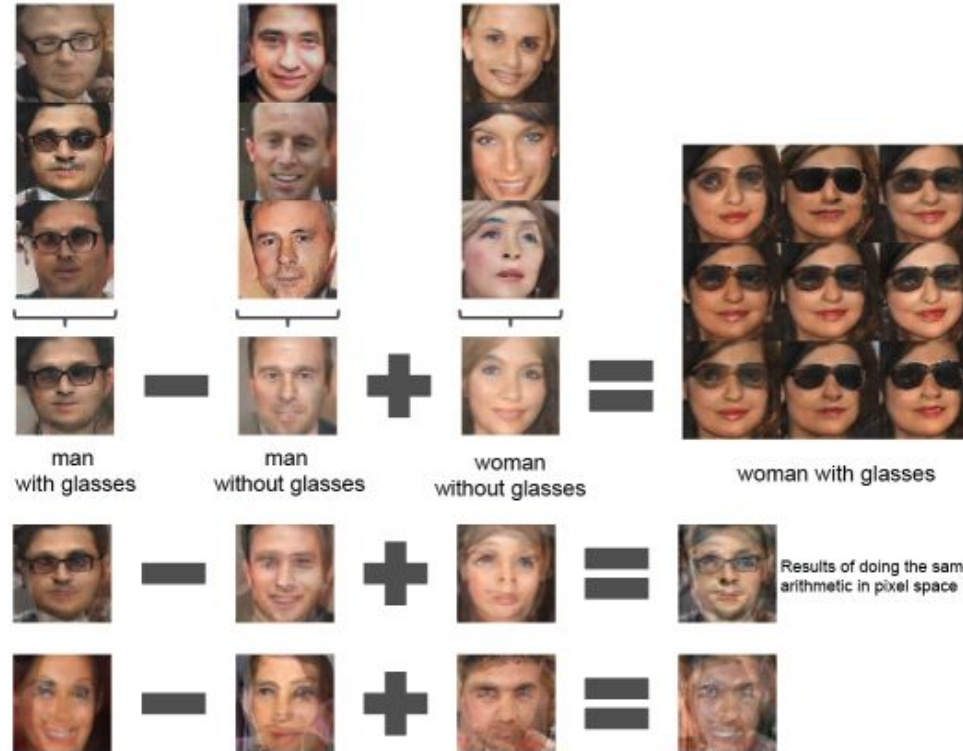
- ❖ Generator is an upsampling network with fractionally-strided convolutions
- ❖ Discriminator is a convolutional network
- ❖ Eliminating fully connected layers on top of Convolutional Layers (instead of max polling)
- ❖ Batch Normalization
- ❖ Leaky ReLUs



# Vector Arithmetic For Visual Concepts

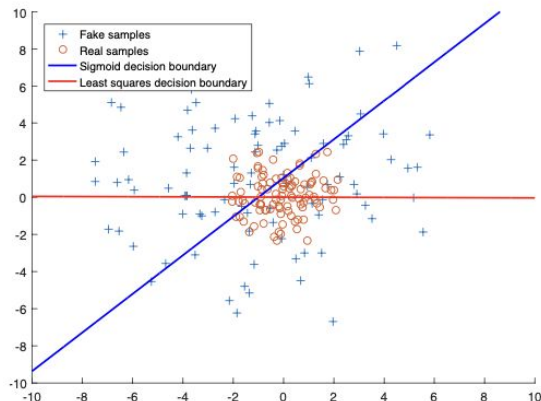
Samples  
from the  
model

Average Z  
vectors, do  
arithmetic

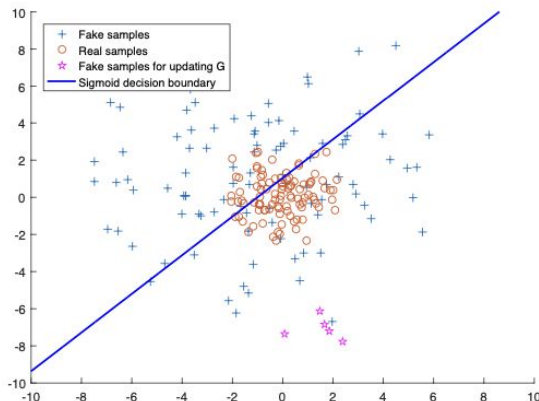


# Least Square GAN (LSGAN)

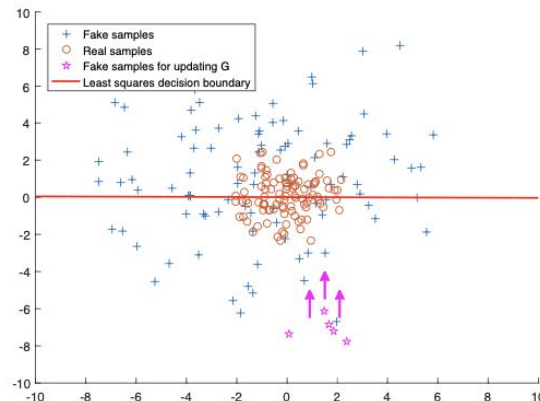
**Idea:** Use loss function that provides smooth and non-saturation gradient in Discriminator



(a)



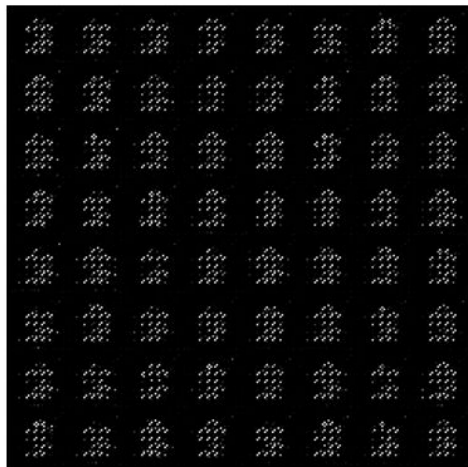
(b)



(c)

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

# Training a GAN: MNIST Data



Epoch 1



Epoch 5



Epoch 10



# Conclusion

GANs,

- ❖ Don't work with an explicit density function
- ❖ Game-theoretic approach: learn to generate from training distribution through two-player game

Pros:

- ❖ Beautiful, state-of-the-art samples!

Cons:

- ❖ Trickier / more unstable to train

Active areas of research:

- ❖ Better loss functions, more stable training (Wasserstein GAN, LSGAN, many others)
- ❖ Conditional GANs, GANs for all kinds of applications

# Explosion of GANs



LSGAN, Zhu 2017.



Wasserstein GAN,  
Arjovsky 2017.  
Improved Wasserstein  
GAN, Gulrajani 2017.



Progressive GAN, Karras 2018.