

Do you like to be successful?

Able to see the big picture

Are you able to recognise a  
scientific  
**GEM**

# How to recognise good work? suggestions please

item#1 1st of its kind

item#2 solve problem

item#3 learning process

effort put in

value added

# All these works resulted in Nobel Prizes

Description	Reason for rejection
Enrico Fermi's seminal paper on weak interaction, 1933	It contained speculations too remote from reality to be of interest to the reader
Hans Krebs' paper on the citric acid cycle, AKA the Krebs cycle, 1937	Nature publishing had a huge backlog
Murray Gell-Mann's work on classifying the elementary particles, 1953	He did not use a name editor like for his discovery
The invention of the radioimmunoassay, 1955	reviewers were skeptical that humans could make antibodies small enough to bind to things like insulin
The discovery of quasicrystals, 1984	It was rejected on the grounds that it will not interest physicists
Paper outlining nuclear magnetic resonance (NMR) spectroscopy, 1966	rejected twice by the J. Chem. Phys. to be finally accepted and published in the Review of Scientific Instruments

- ❖ Is your idea so radical that your professor thinks you are out of your mind?
- ❖ Do you discover nature?
- ❖ Are you trying to break conventional wisdom?
- ❖ If you are doing deep learning for healthcare, are you listening to good doctors?  
engineers building a tennis racket must ask the tennis players

# Generative Adversarial Networks

This (GAN), . . . is the most interesting idea in the last 10 years in ML, in my opinion.

... Yann LeCun

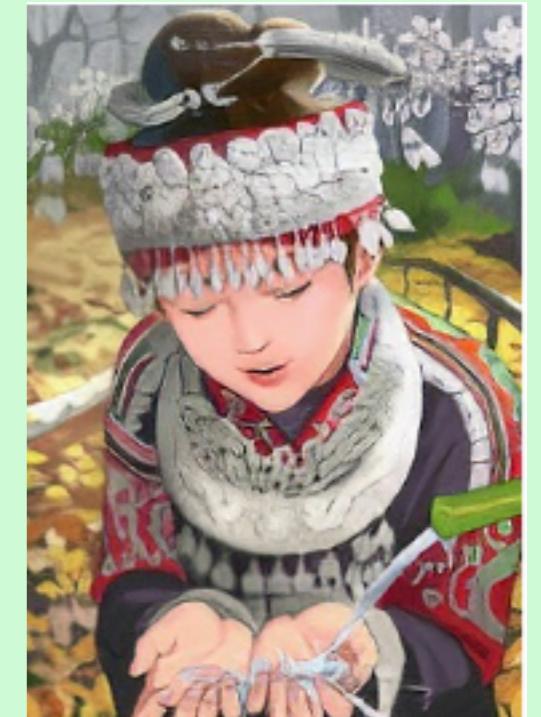
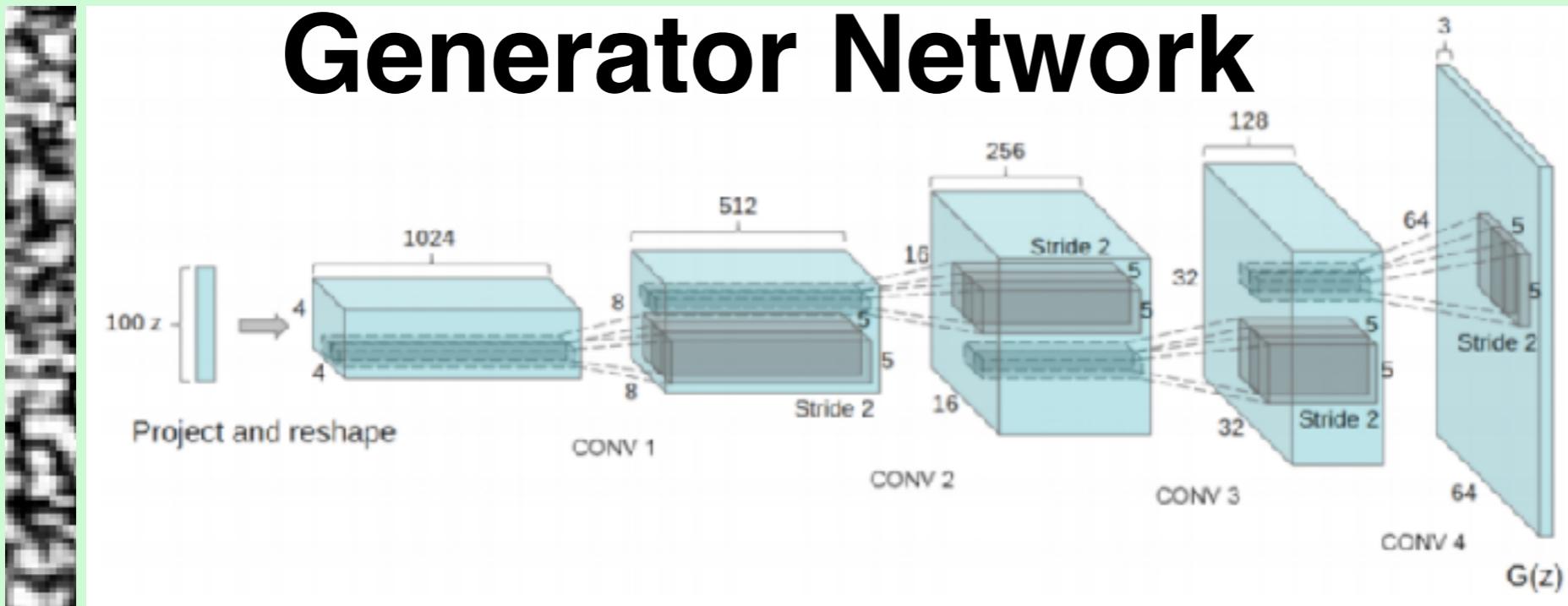
# Generative



# Generative

sampled  
noise  
input

image  
generated

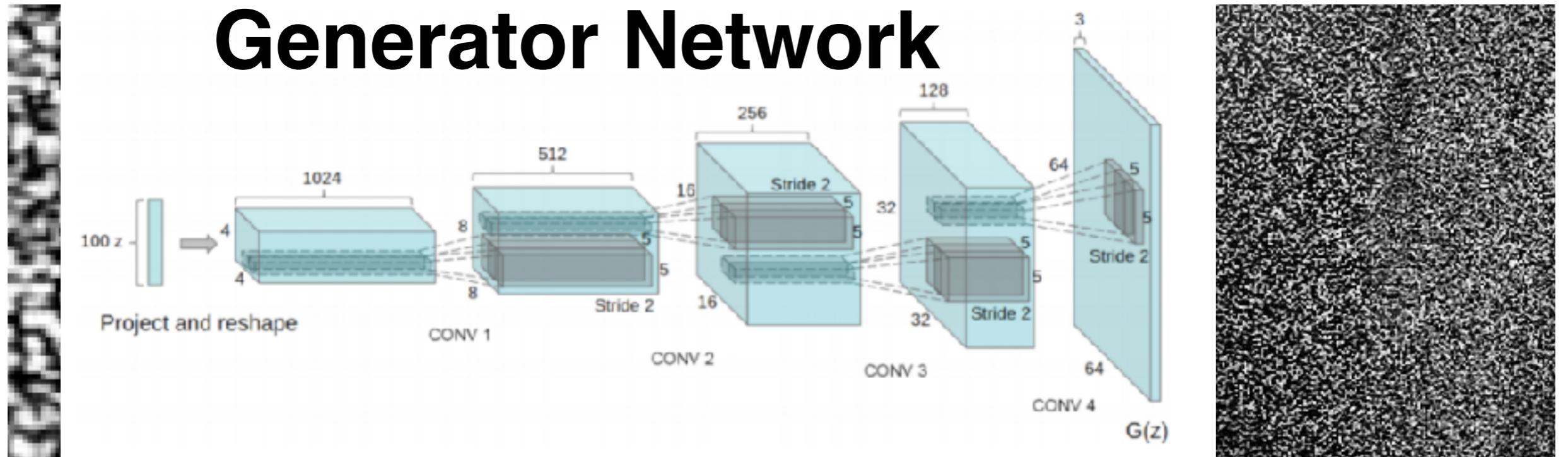


<https://www.linkedin.com/pulse/unsupervised-representation-learning-deep-generative-diego>

# Adversarial



When the generator network is not trained well

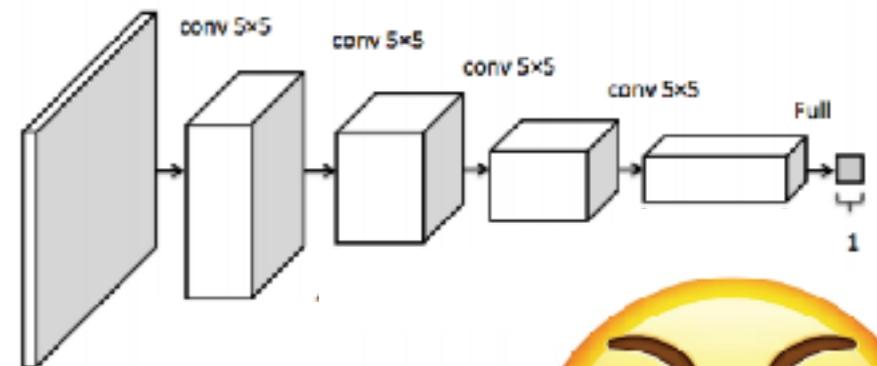
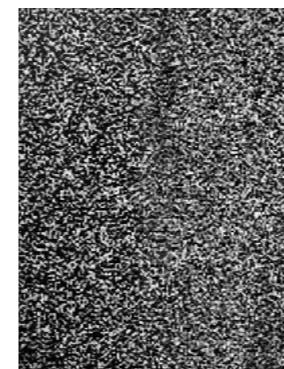


get nonsense instead  
of a nice picture

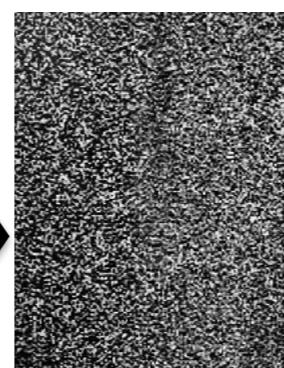
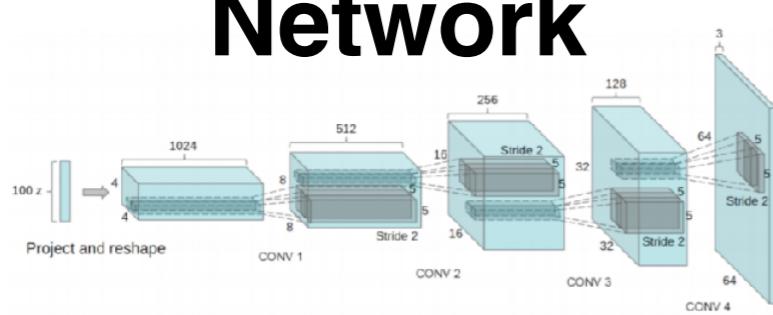


Training the generator via an **adversarial**  
the discriminator network

## Discriminator Network

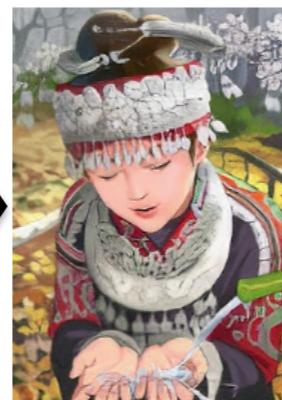
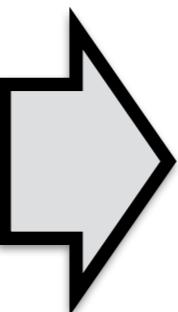
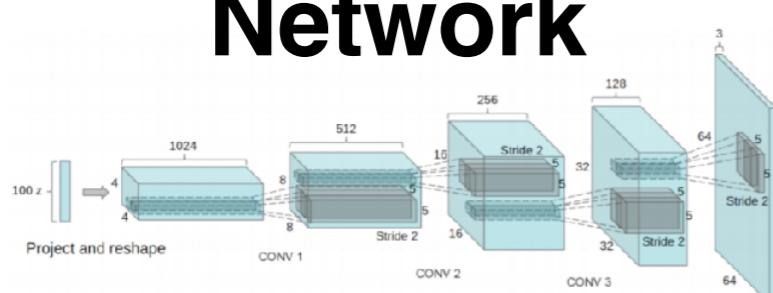


## Generator Network

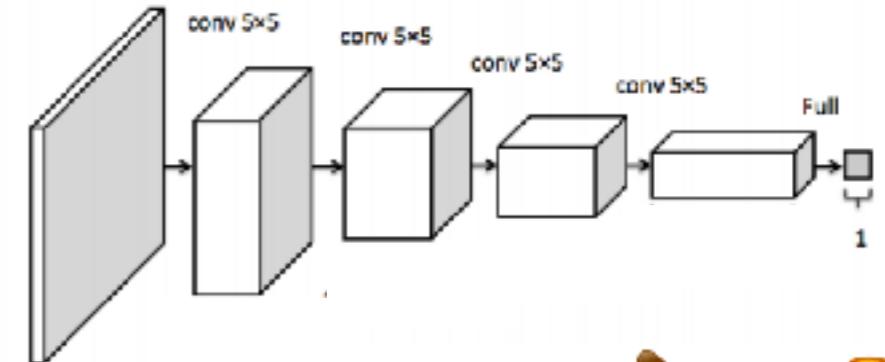


Training the generator via an **adversarial**  
the discriminator network

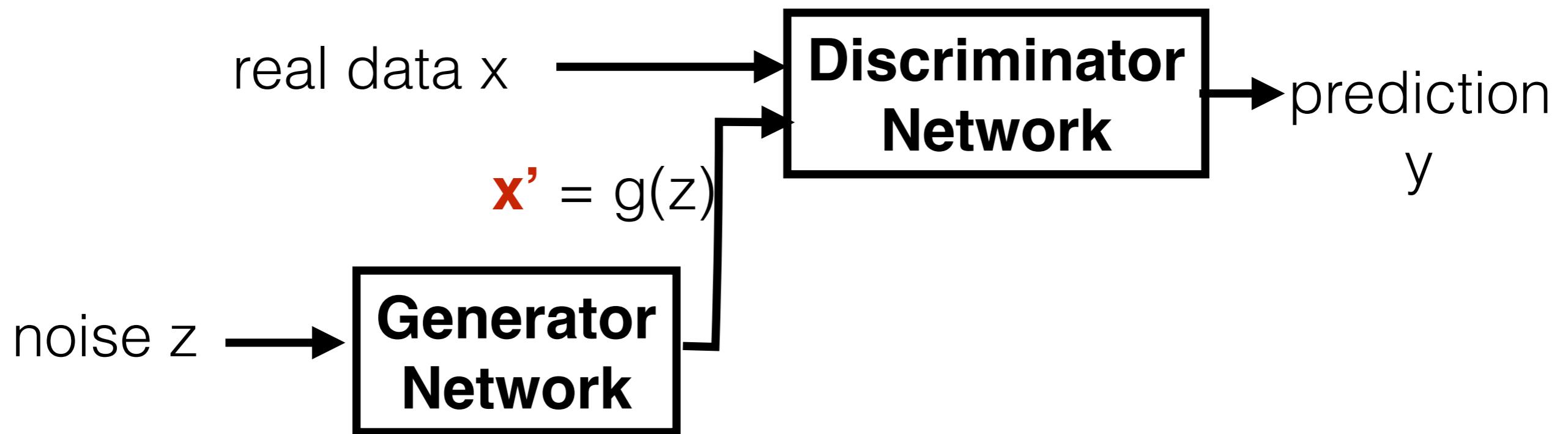
## Generator Network



## Discriminator Network

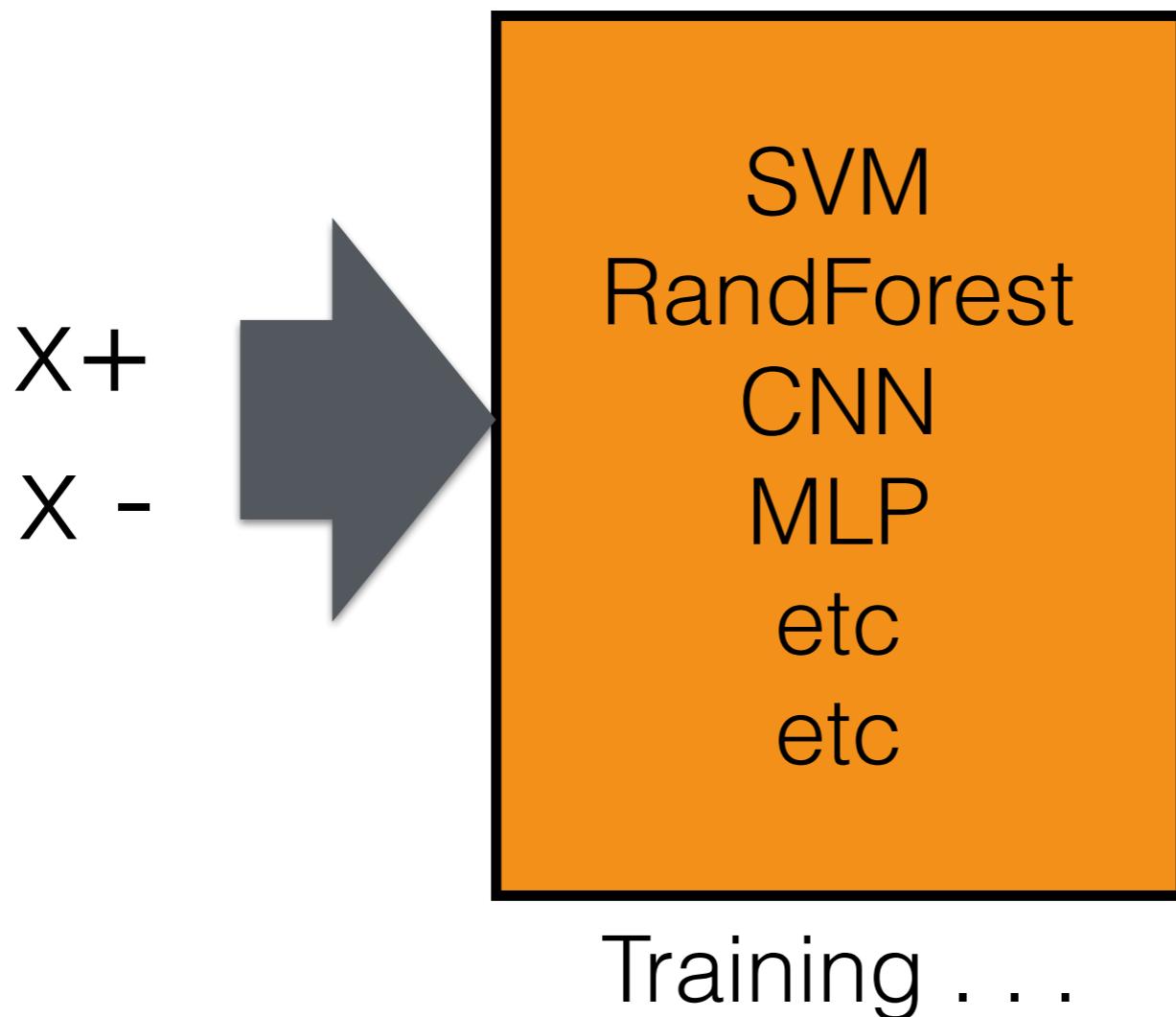


# The original GAN diagram



# “Classical Supervised Learning”

## The **passive learner**



# “Classical Supervised Learning”

## **The passive learner**

- ♣ It ‘eats’ what you gave to it
- ♣ No data = ‘dead learner’
- ♣ Bad data = ‘bad learner’
- ♣ Has totally no control over data
- ♣ Retrain when data change a little bit

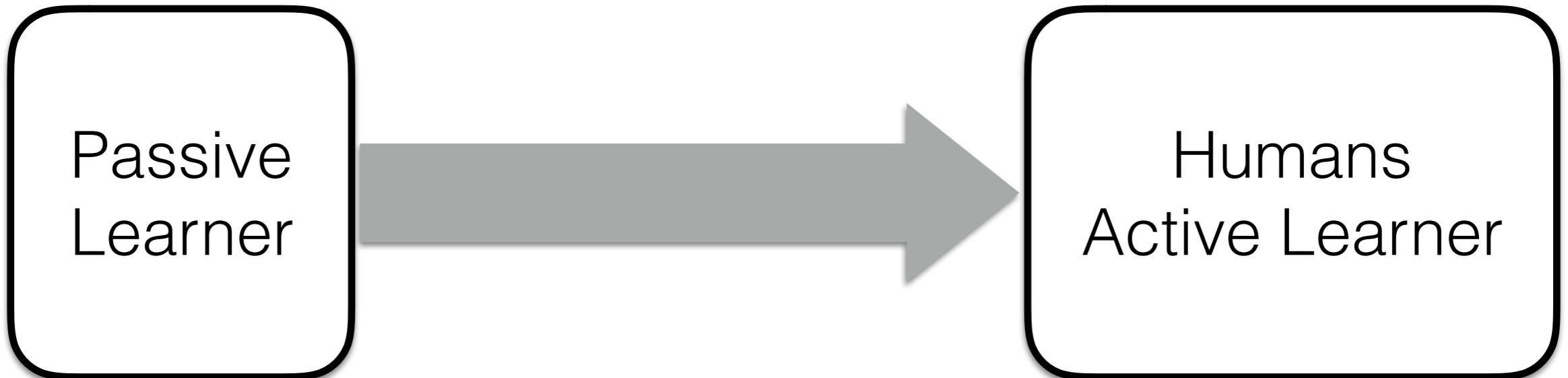
Is this an intelligent machine?

# Humans

## The active learner

- ♣ Select to pay attention or not pay attention
- ♣ Creative and perform thought experiments
- ♣ Generalises well, just little bit of data is ok
- ♣ We perform experiments to generate data

**Learn by analogy :: 举一反三**



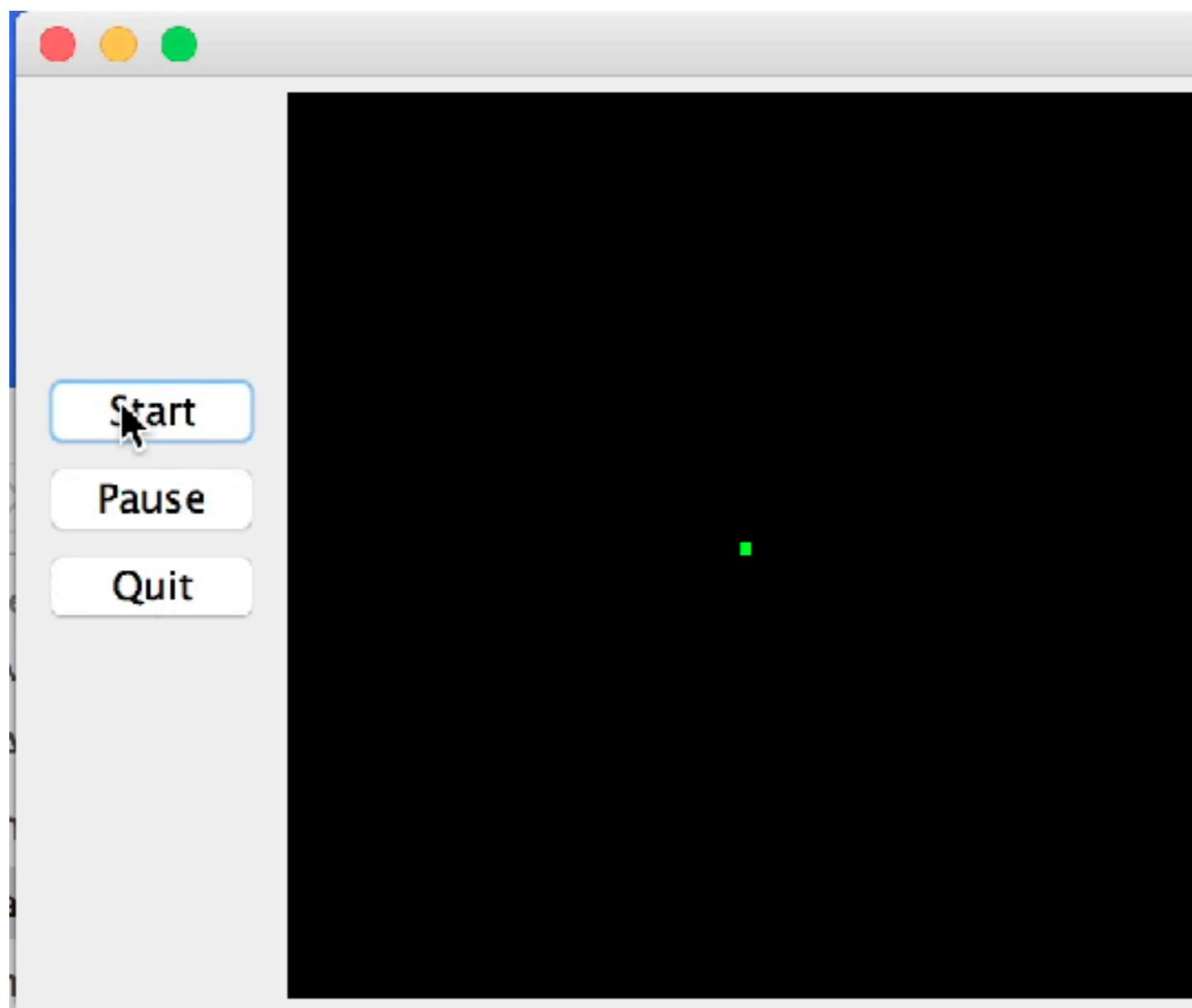
# where is GAN?

- X Select to pay attention or not pay attention
- X Creative and perform thought experiments
- ? Generalises well, just little bit of data is ok
- ✓ Performs experiments to generate data

# Generative models and their problems

All generative models boils down to sampling data from a distribution

$$(x_1, x_2) \sim a_1 N(\mu_1, \sigma_1) + a_2 N(\mu_2, \sigma_2)$$



# Sampling black / white pixels

$$(x_1, x_2, \dots, x_N) \sim \frac{\exp(-E(x_1, \dots, x_N)/T)}{Z} \quad x_i = \pm 1$$
$$Z = \sum_{x_1, x_2, \dots, x_N}^{2^N} \exp(-E(x_1, \dots, x_N)/T)$$

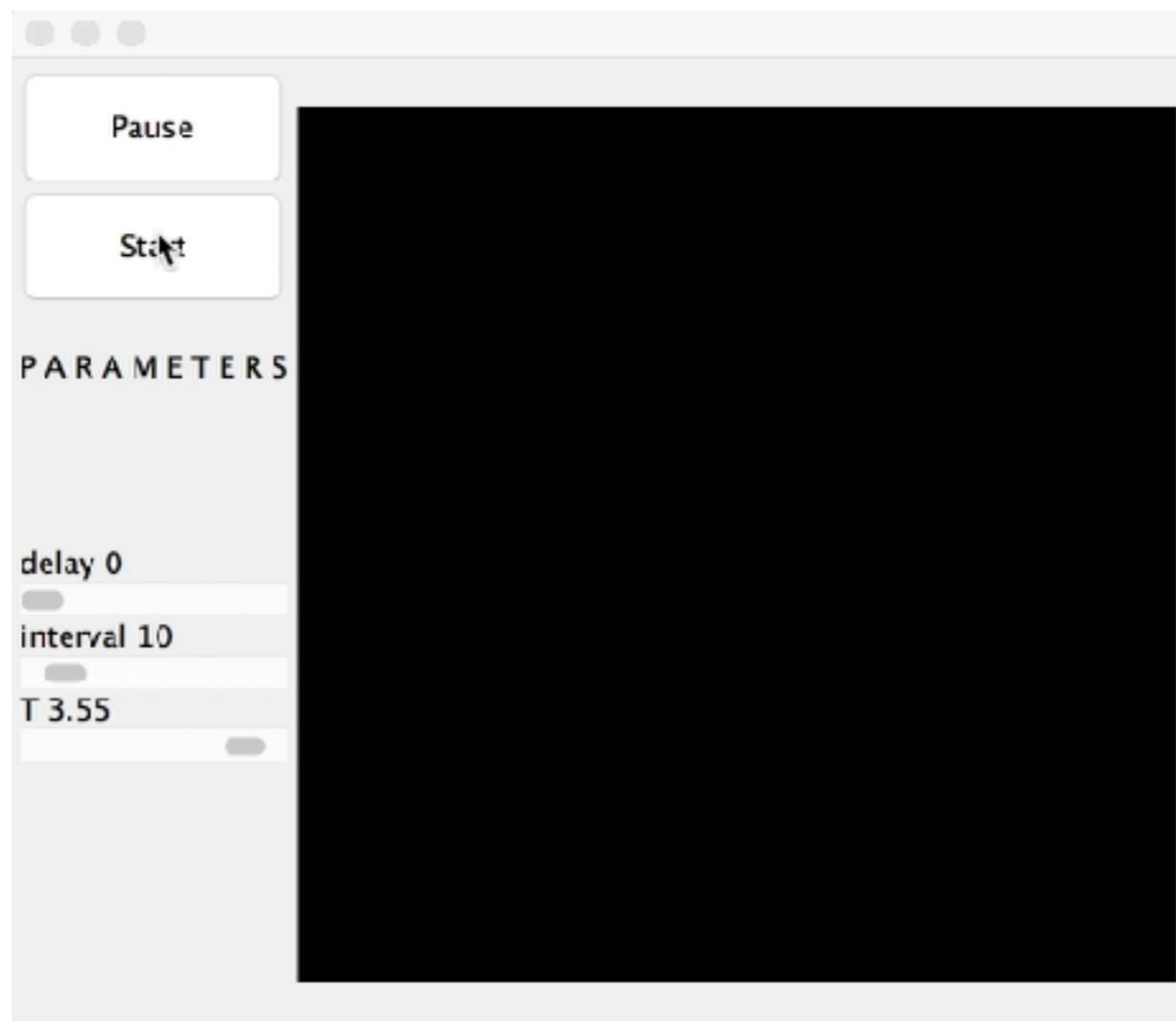
$E(x_1, \dots, x_N) = \sum$  number of neighboring  $x_i$  with the same sign

$x1$	$x6$	$x11$	$x16$	$x21$	$x26$
$x2$	$x7$	$x12$	$x17$	$x22$	$x27$
$x3$	$x8$	$x13$	$x18$	$x23$	$x28$
$x4$	$x9$	$x14$	$x19$	$x24$	$x29$
$x5$	$x10$	$x15$	$x20$	$x25$	$x30$

# Sampling black / white pixels

$$(x_1, x_2, \dots, x_N) \sim \frac{\exp(E(x_1, \dots, x_N)/T)}{Z} \quad x_i = \pm 1$$
$$Z = \sum_{x_1, x_2, \dots, x_N}^{2^N} \exp(E(x_1, \dots, x_N)/T)$$

$E(x_1, \dots, x_N) = \sum$  number of neighboring  $x_i$  with the same sign



How to sample is the key problem in generative models

This one is easy

$$(x_1, x_2) \sim a_1 N(\mu_1, \sigma_1) + a_2 N(\mu_2, \sigma_2)$$

This one is harder but can do with advanced methods

$$(x_1, x_2, \dots, x_N) \sim \frac{\exp(-E(x_1, \dots, x_N)/T)}{Z}$$

In general if we can build a distribution, theoretical or empirical distribution will do, we can find a way to sample

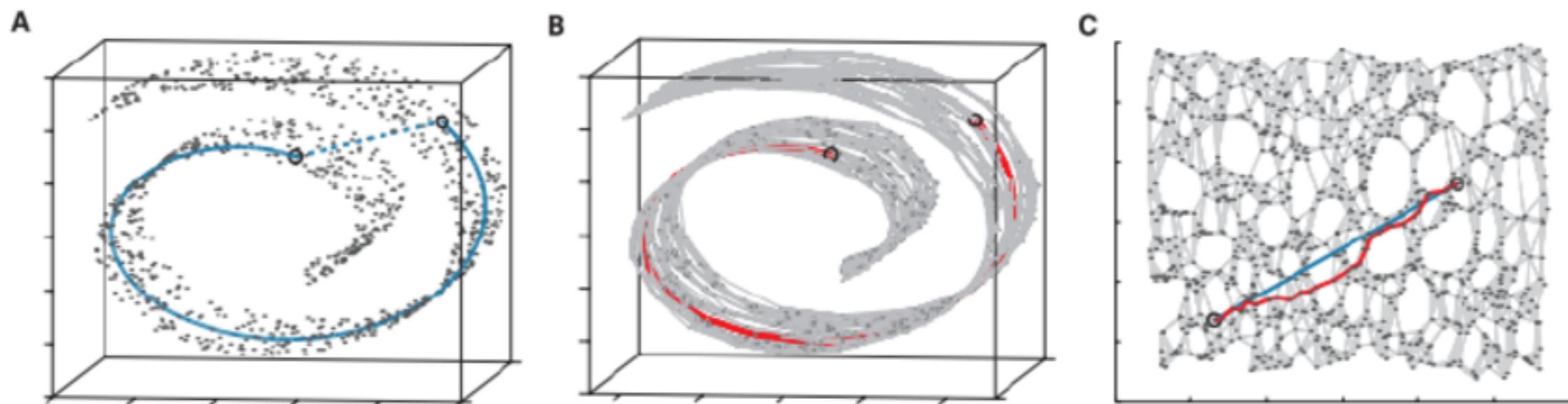
The difficulty is building a distribution  
We don't have the slightest idea how the distribution looks like

How can a distribution be build in a high dimension space where data is very sparse?

e.g. images with one mega pixels of greyscale [0,255], is embedded in a volume  $255^{\{1,000,000\}}$

Image data are certainly always very very sparse, it is believed that they lie in some very complex manifold.

Many real world data too occupy infinitesimal volume of feature space and form very complex manifold.



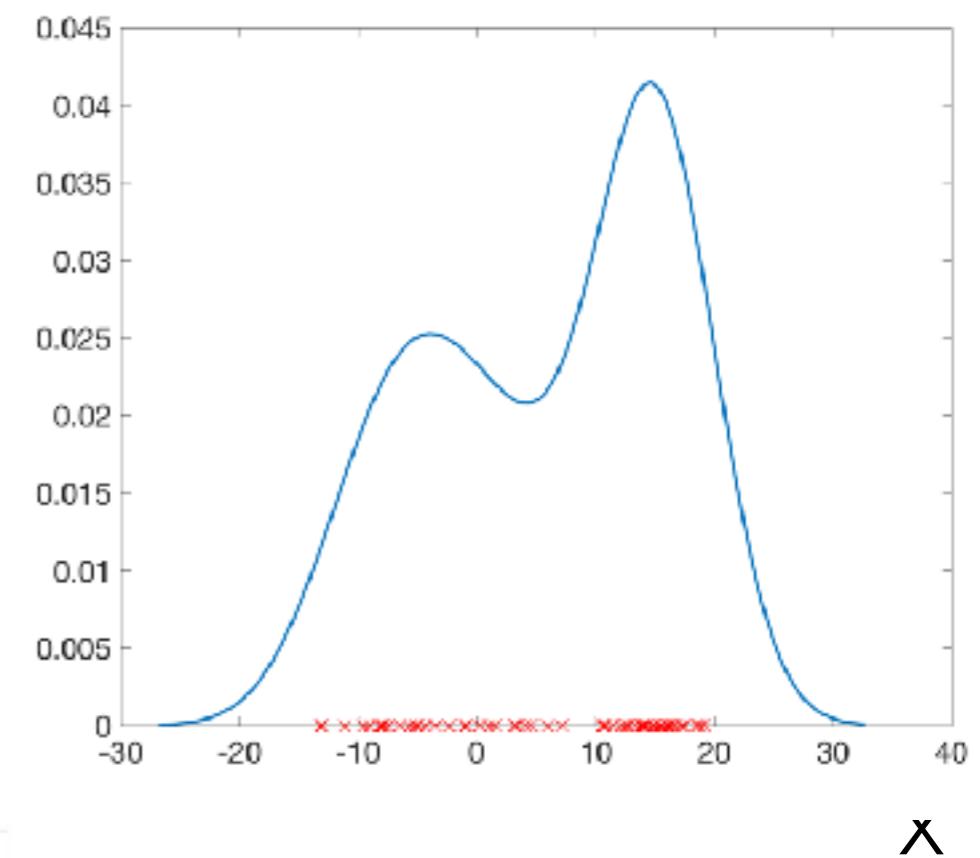
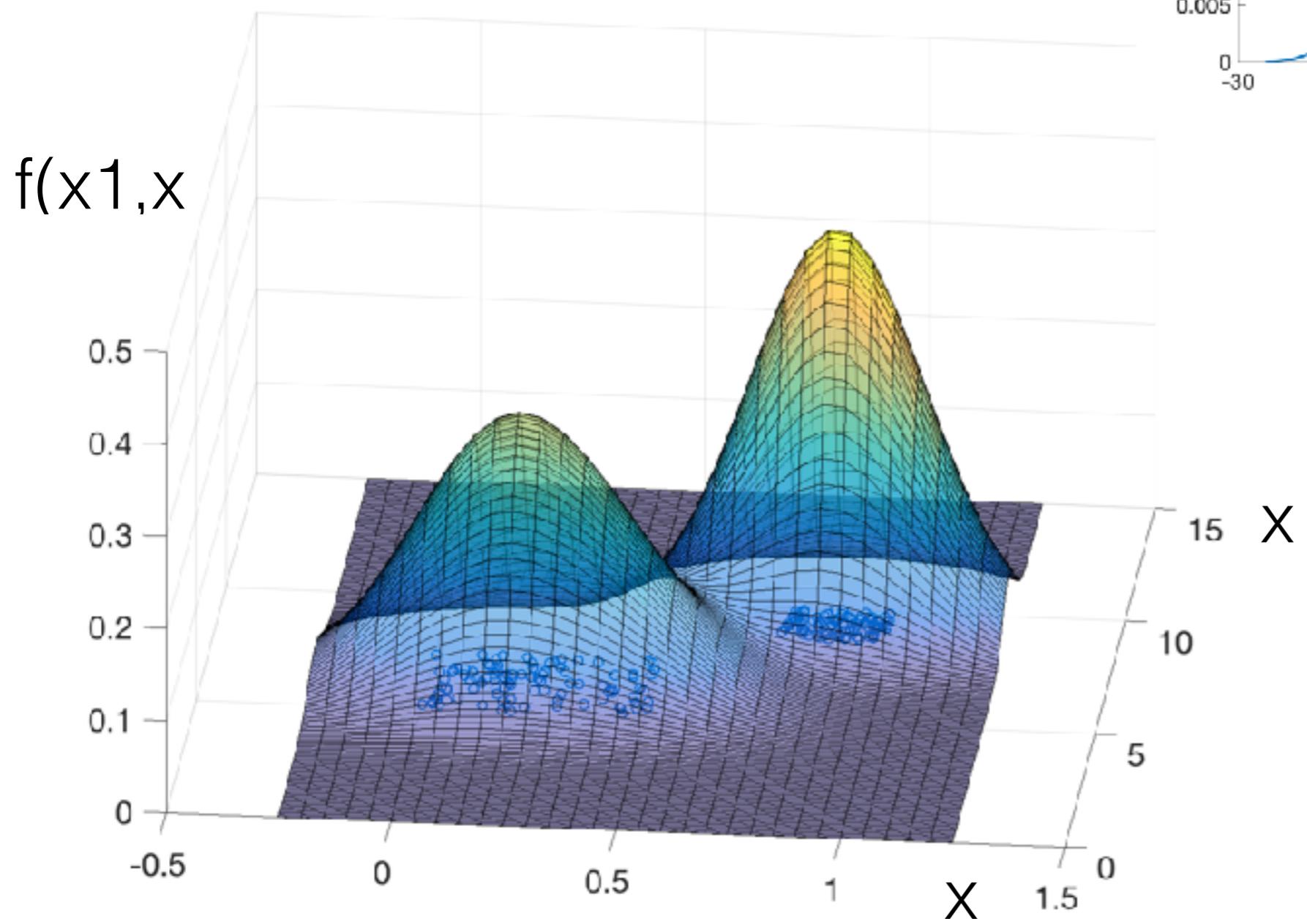
[https://www.researchgate.net/figure/286513346\\_fig1\\_Figure-10-The-swiss-roll-dataset-a-The-euclidean-distance-blue-dashed-line-of-two](https://www.researchgate.net/figure/286513346_fig1_Figure-10-The-swiss-roll-dataset-a-The-euclidean-distance-blue-dashed-line-of-two)

Given a data set  $x_i$

Distribution can be constructed by finding a function  $f(x)$  that is large ‘near’ data points and small elsewhere by minimising a cost function over  $f(x)$

$$C = - \int_{x \text{ 'near' } x_i} f(x') dx' + \int_{x \text{ everywhere else}} f(x') dx'$$

$$C = - \int_{x \text{ 'near' } x_i} f(x') dx' + \int_{x \text{ everywhere else}} f(x') dx'$$



Given a data set  $x_i$

Distribution can be constructed by finding a function  $f(x)$  that is large `near' data points and small elsewhere by minimising a cost function over  $f(x)$

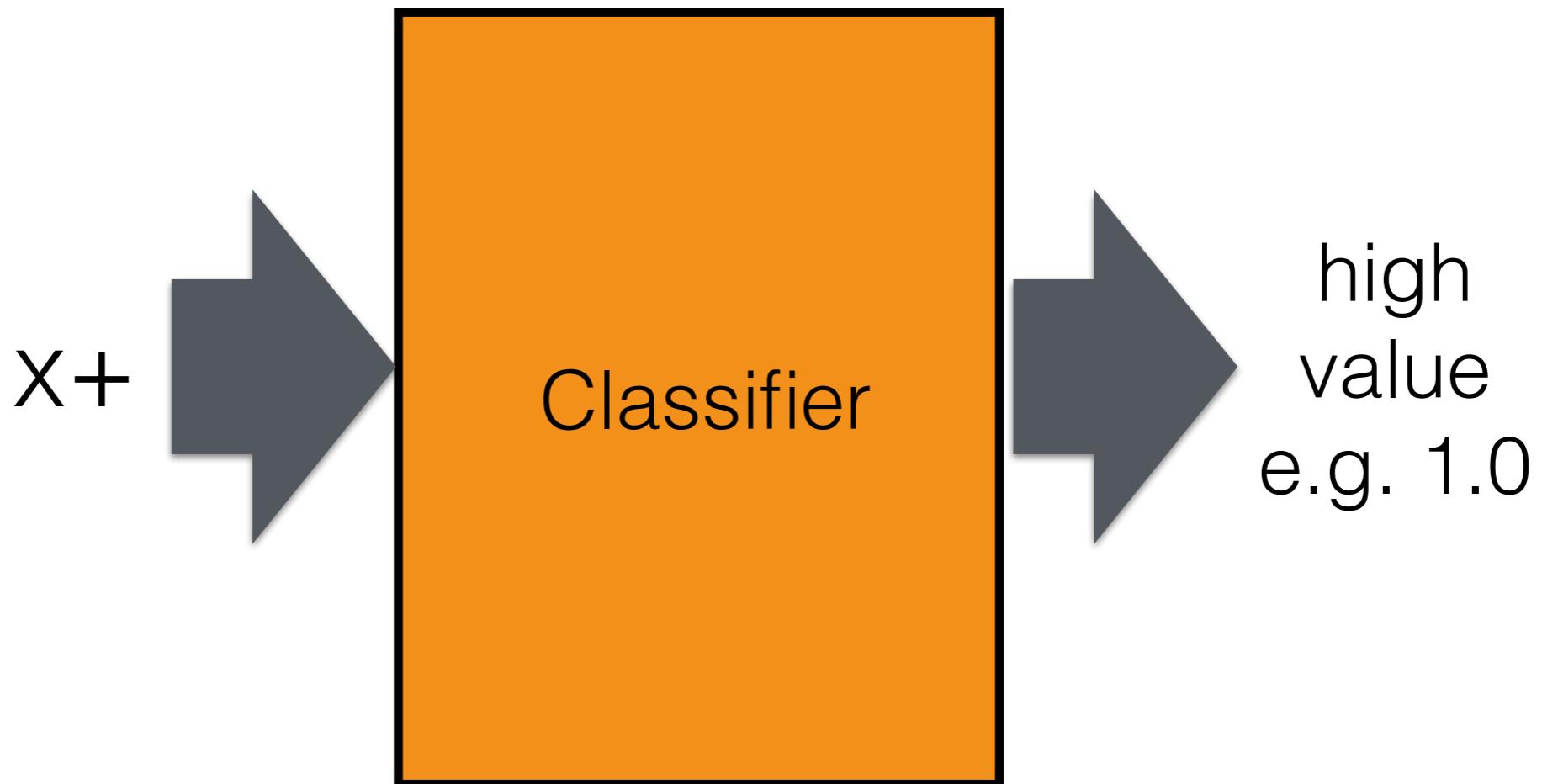
$$C = - \int_{x \text{ 'near' } x_i} f(x') dx' + \int_{x \text{ everywhere else}} f(x') dx'$$

Two difficulties here:

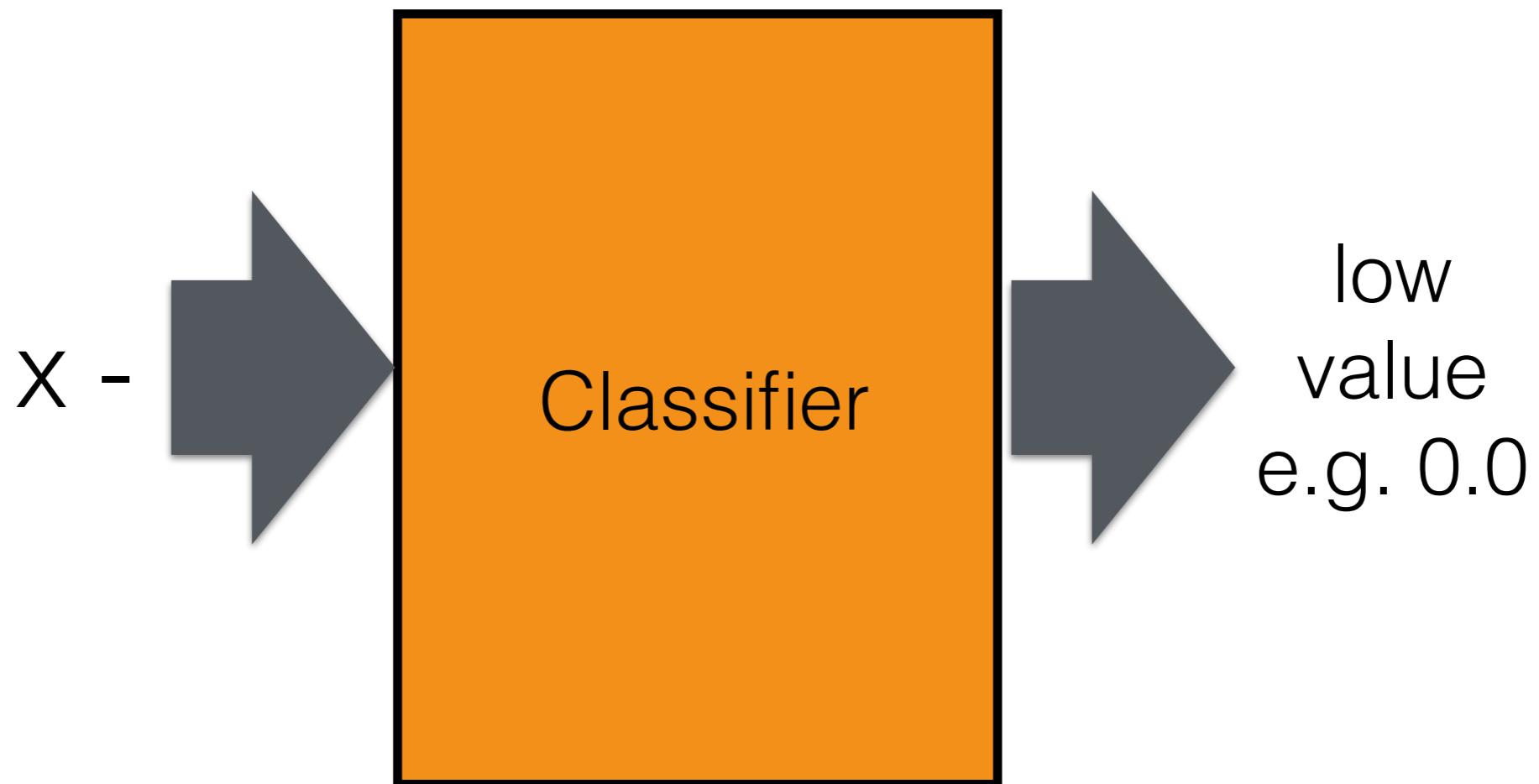
- ♣ Defining `near' - regularisation issue
- ♣ High dimensional integrals, 1,000,000 dimension for image data

Neural networks offer a pseudo-solution

Lets start by reviewing the classical  
supervise learning on two classes



Lets start by reviewing the classical  
supervise learning on two classes





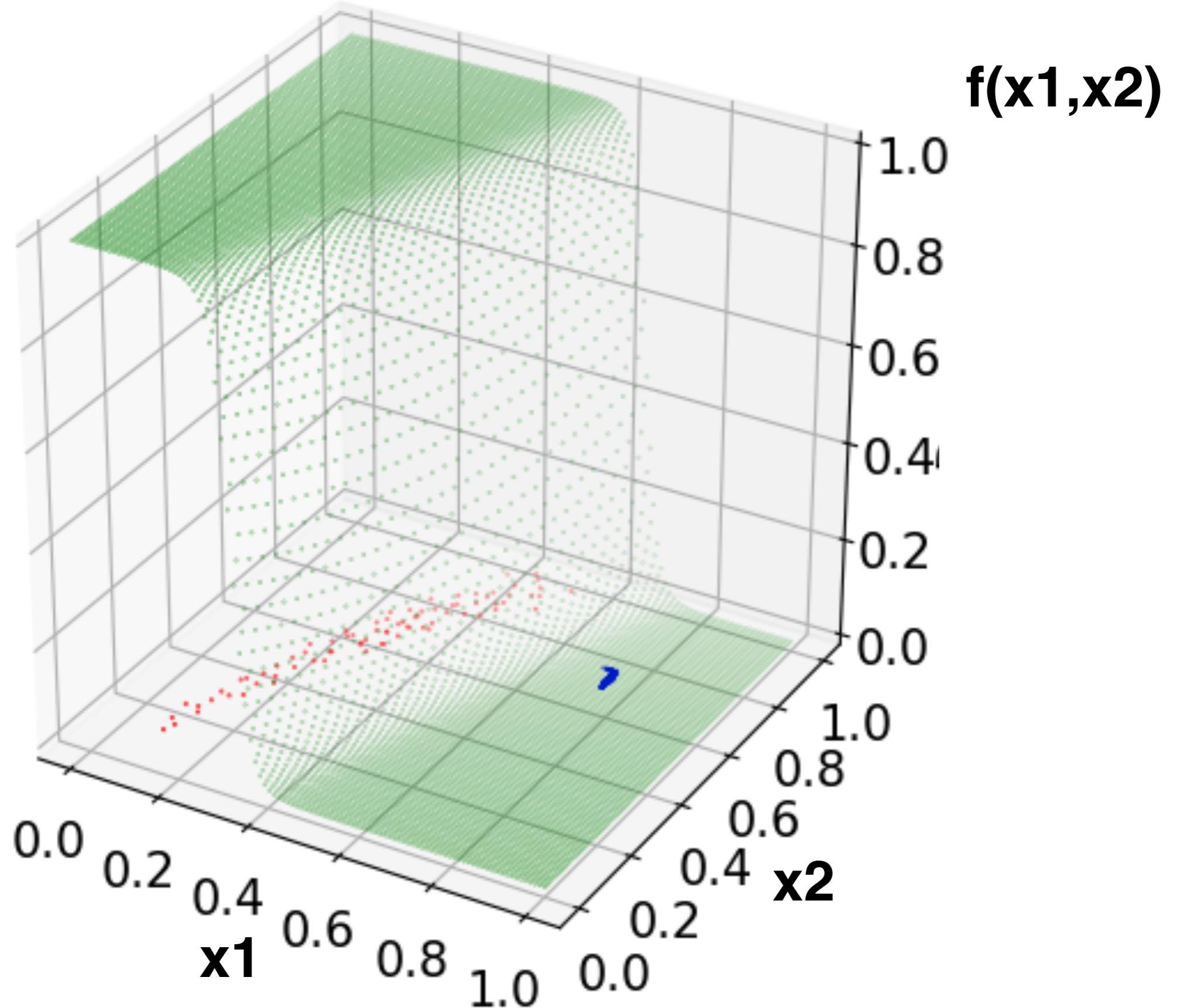
Minimize a cost function with respect to  $\mathcal{W}$

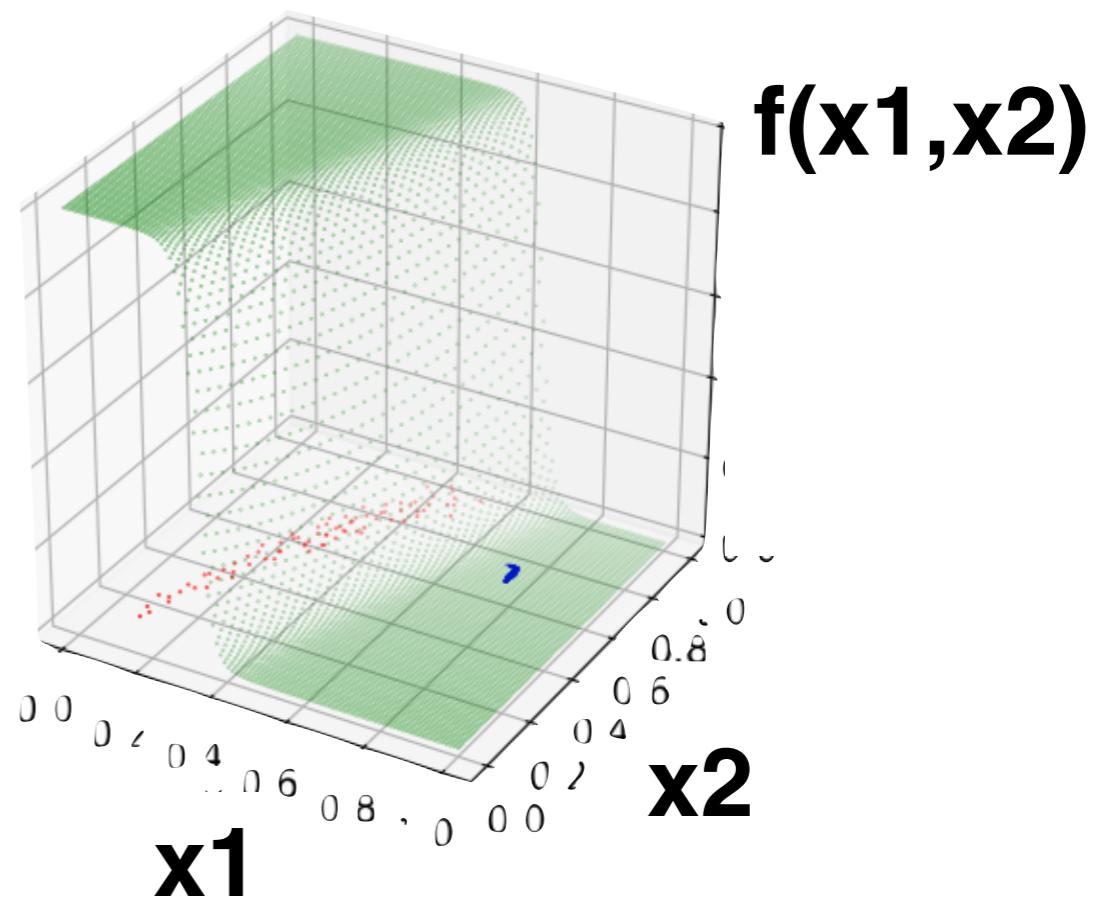
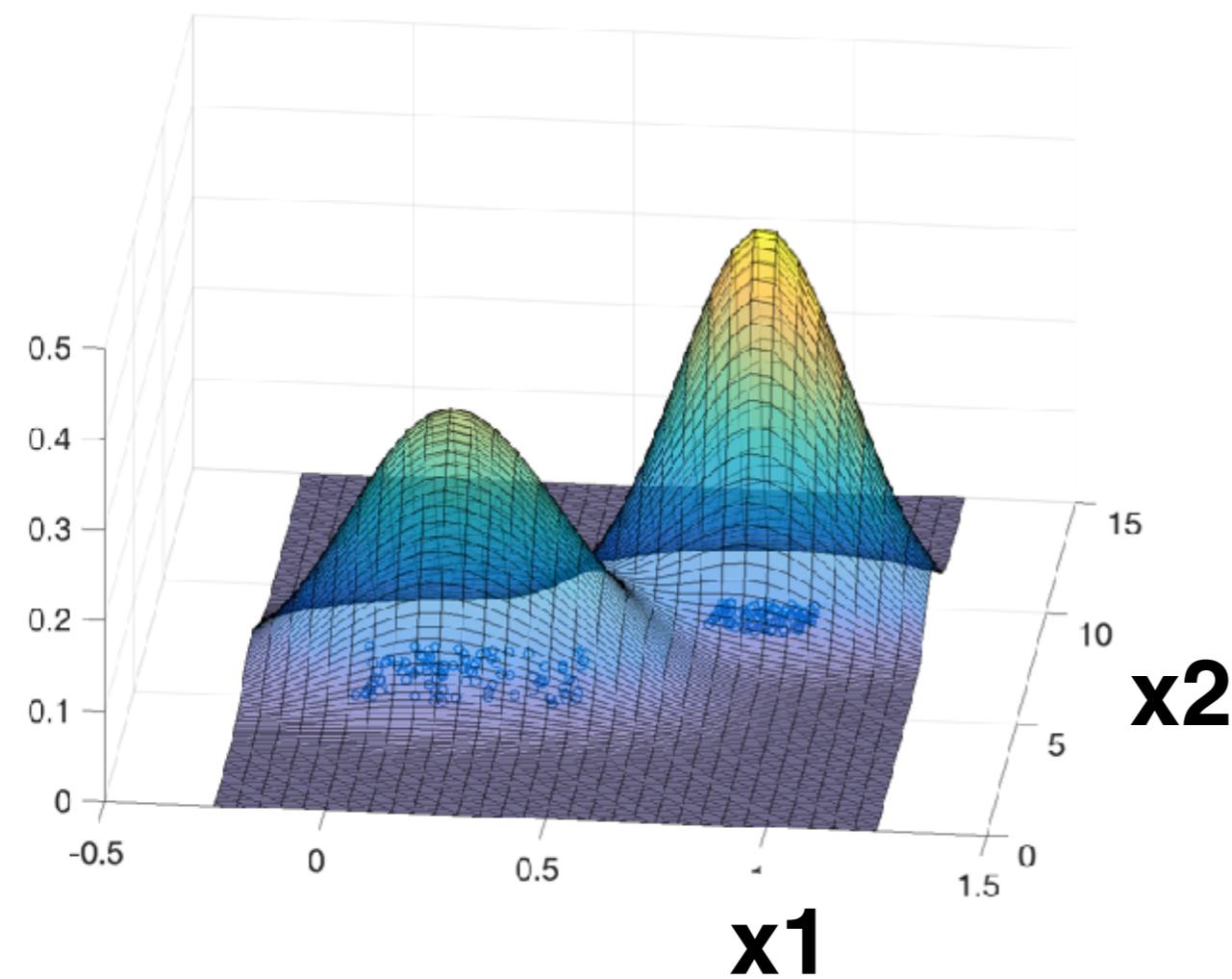
$$C(w) = -\frac{1}{N^+} \sum_i^{N^+} f_w(x_i) + \frac{1}{N^-} \sum_j^{N^-} f_w(x_j)$$

An example

$$C(w) = - \left( \frac{1}{N^+} \sum_i \log(D_w(x_i)) + \frac{1}{N^-} \sum_j \log(1 - D_w(x_j)) \right)$$

$$D_w(x_i) \in (0, 1)$$





# Idea of GAN (Generative Adversarial Networks)

- ❖ Given a set of example real data
- ❖ Generate “fake” data that looks like real data

Generator is a neural network parametrised by  $\theta$

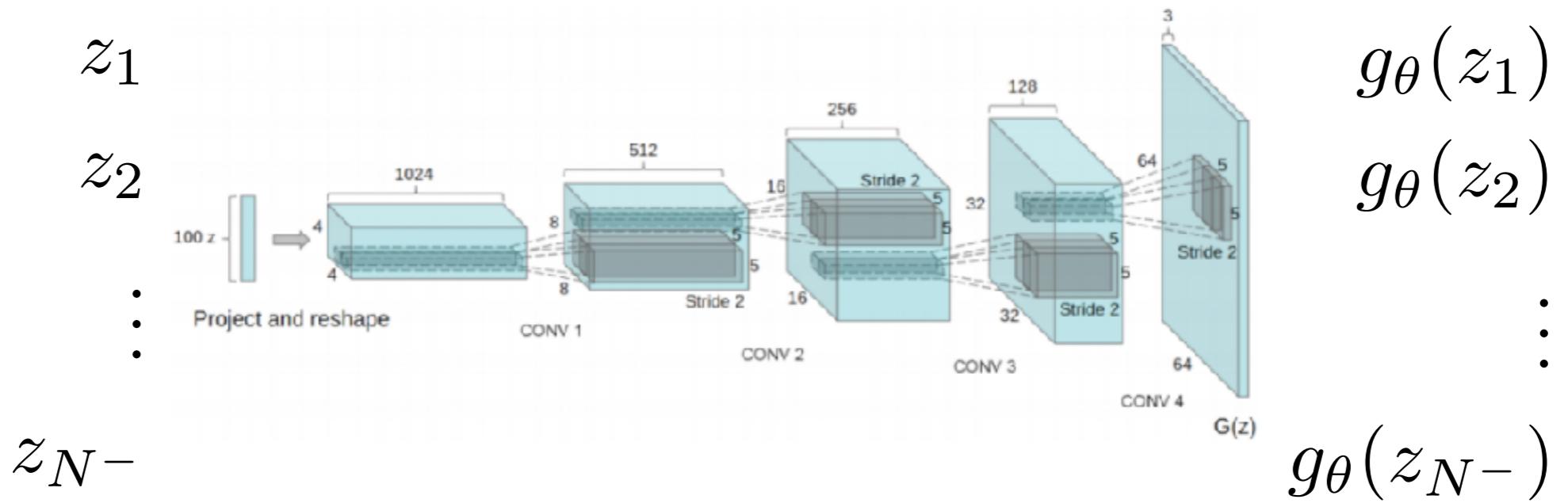
Sample a set of random noise

$$z_i \quad i = 1, \dots N^-$$

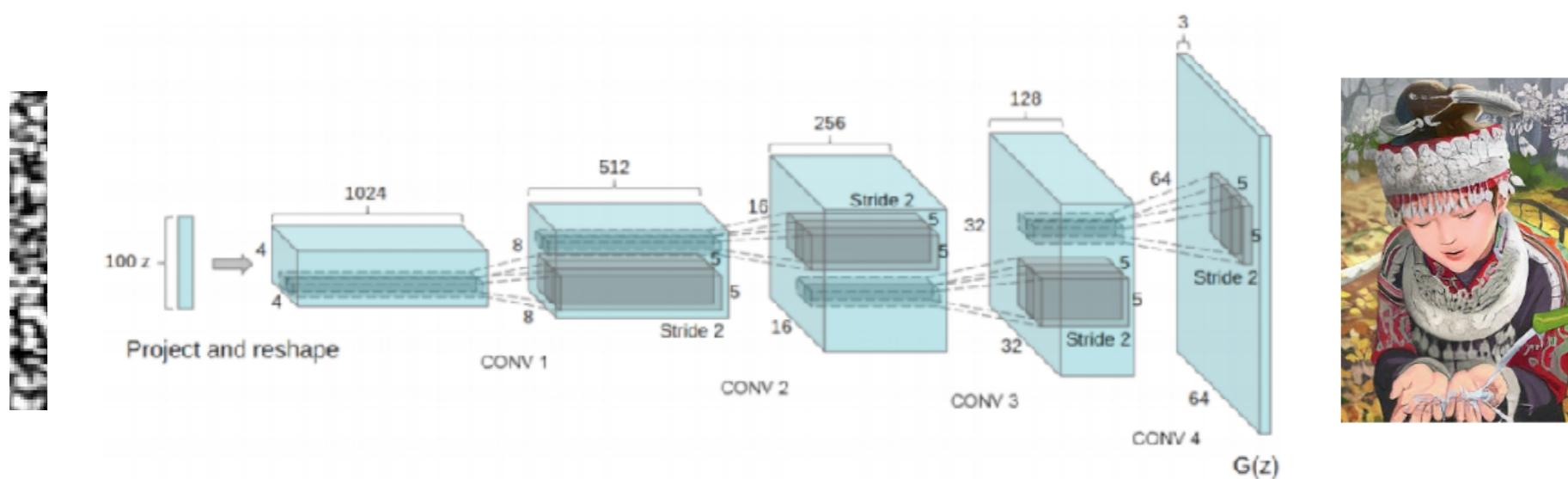
Generate set of fake data from sampled noise

$$g_\theta(z_i) \quad i = 1, \dots N^-$$

# Generator Network



If we adjust  $\theta$  well,  $g_\theta(z_i)$  will look like this

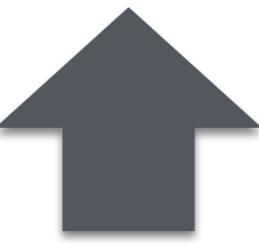


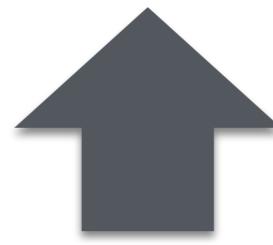
# Training Generator network using Discriminator network

Discriminator network needs positive and negative training examples



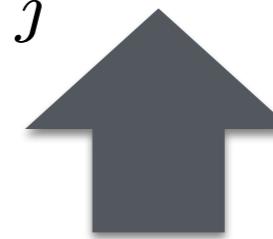
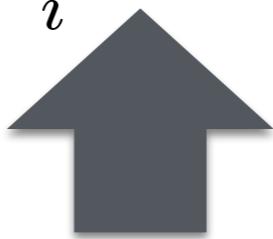
$$C(w) = -\frac{1}{N^+} \sum_i^{N^+} f_w(x_i) + \frac{1}{N^-} \sum_j^{N^-} f_w(x_j)$$

  
put real data  
examples here

  
put “fake”  
generated data here

## First train the Discriminator network

$$C(w, \theta) = -\frac{1}{N^+} \sum_i^{N^+} f_w(x_i) + \frac{1}{N^-} \sum_j^{N^-} f_w(g_\theta(z_j))$$



put real data  
examples here

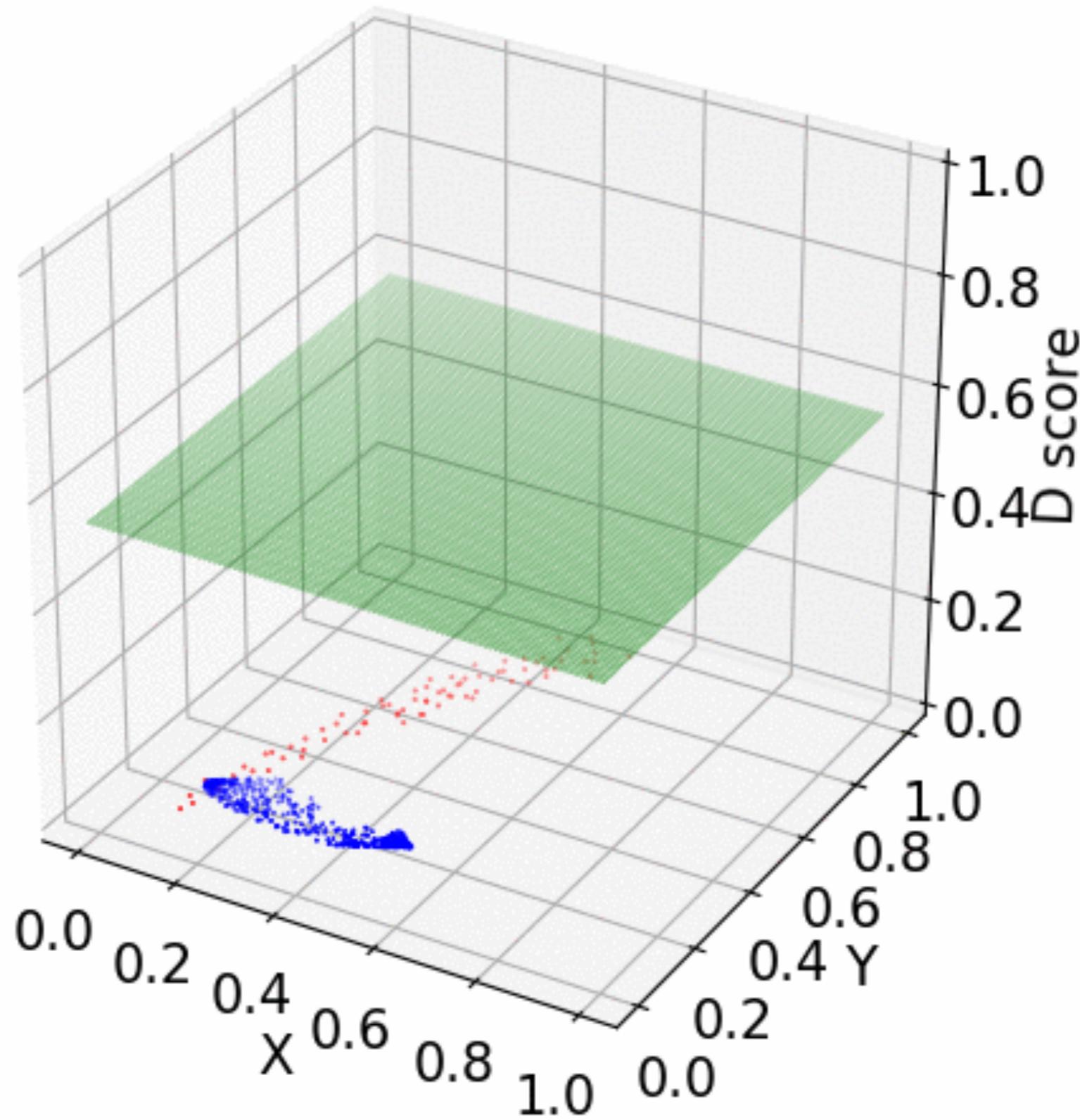
put “fake”  
generated data here

$$w^* = \arg \min_w C(w, \theta) \quad \text{keeping } \theta \text{ constant}$$

Discriminator network is trained to grade generated data as negative

It criticised any imperfections of the generator

Discriminator :  $x \rightarrow 10 \rightarrow 10 \rightarrow 10 \rightarrow \text{out}$



## Next train the Generator network

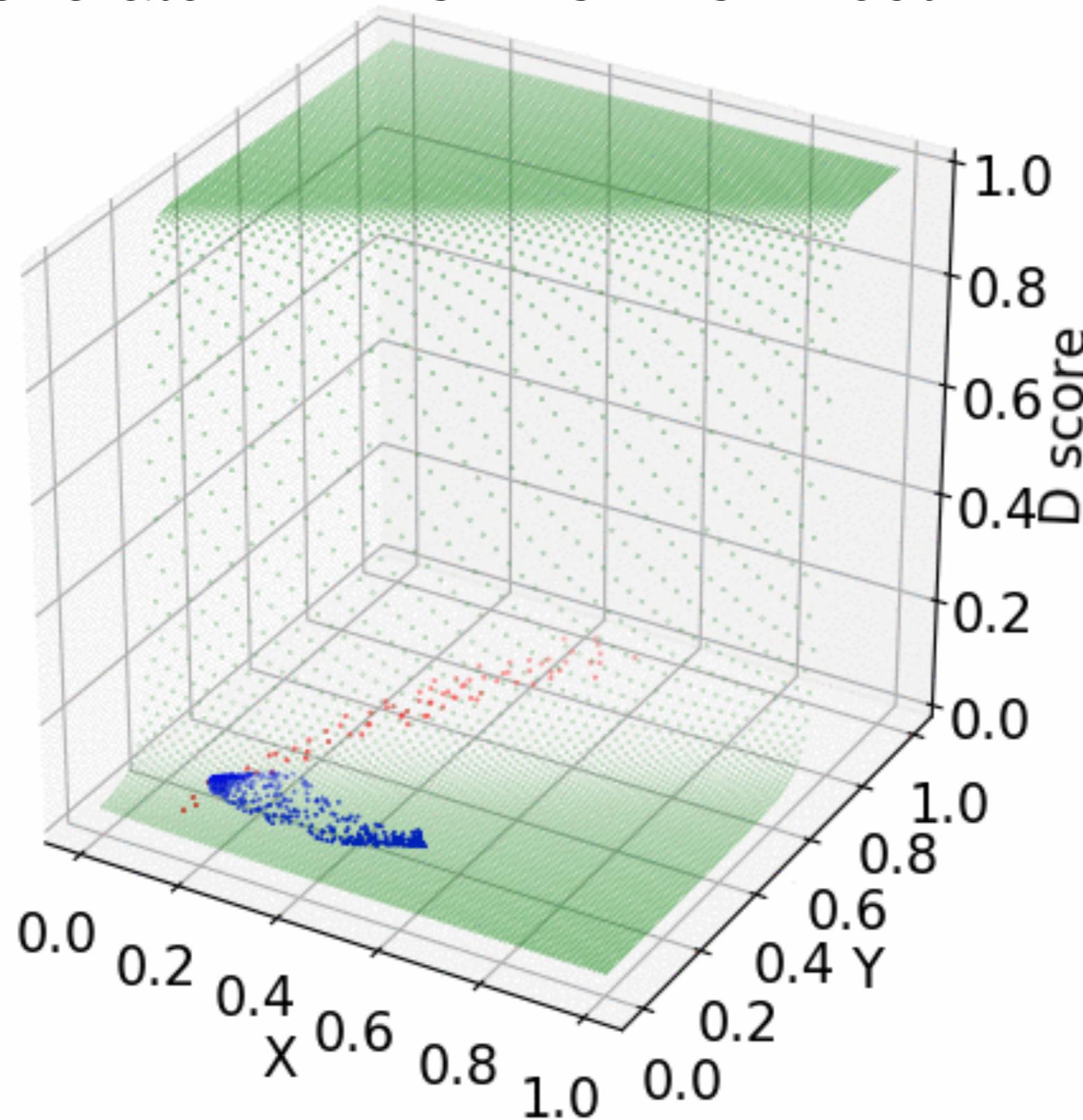
$$C(w, \theta) = -\frac{1}{N^+} \sum_i^{N^+} f_w(x_i) + \frac{1}{N^-} \sum_j^{N^-} f_w(g_\theta(z_j))$$

↑  
put real data  
examples here      ↑  
put “fake”  
generated data here

$$w^* = \arg \min_w C(w, \theta) \quad \text{keeping } \theta \text{ constant}$$

$$\theta^* = \arg \max_\theta C(w, \theta) \quad \text{keeping } w \text{ constant}$$

Generator :  $x \rightarrow 5 \rightarrow 5 \rightarrow 5 \rightarrow \text{out}$



## Cost function use in original GAN paper

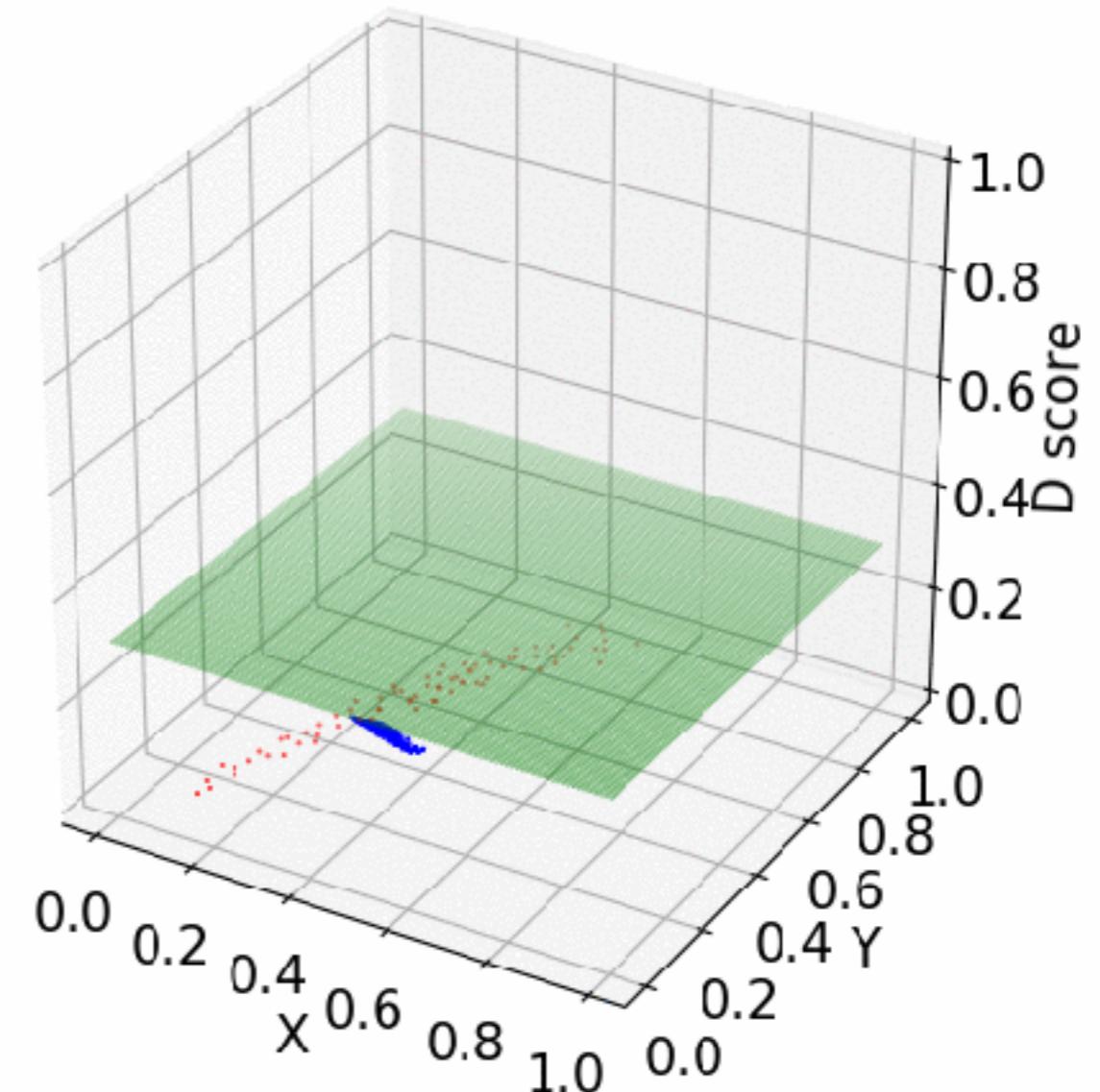
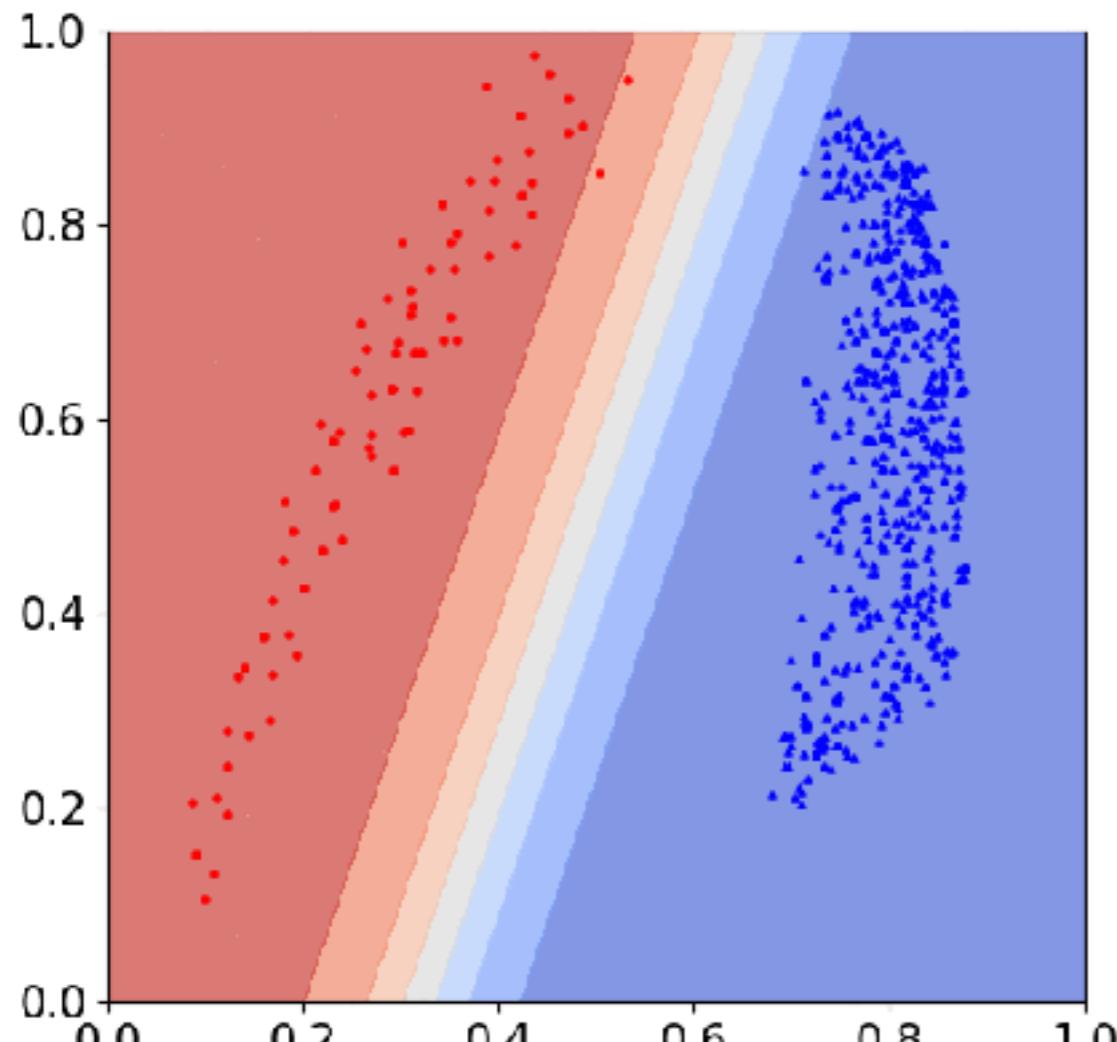
$$C(w, \theta) = - \left( \frac{1}{N^+} \sum_i \log(D_w(x_i)) + \frac{1}{N^-} \sum_j \log(1 - D_w(g_\theta(z_j))) \right)$$

$$w^* = \arg \min_w C(w, \theta) \quad \text{keeping } \theta \text{ constant}$$

$$\theta^* = \arg \max_\theta C(w, \theta) \quad \text{keeping } w \text{ constant}$$

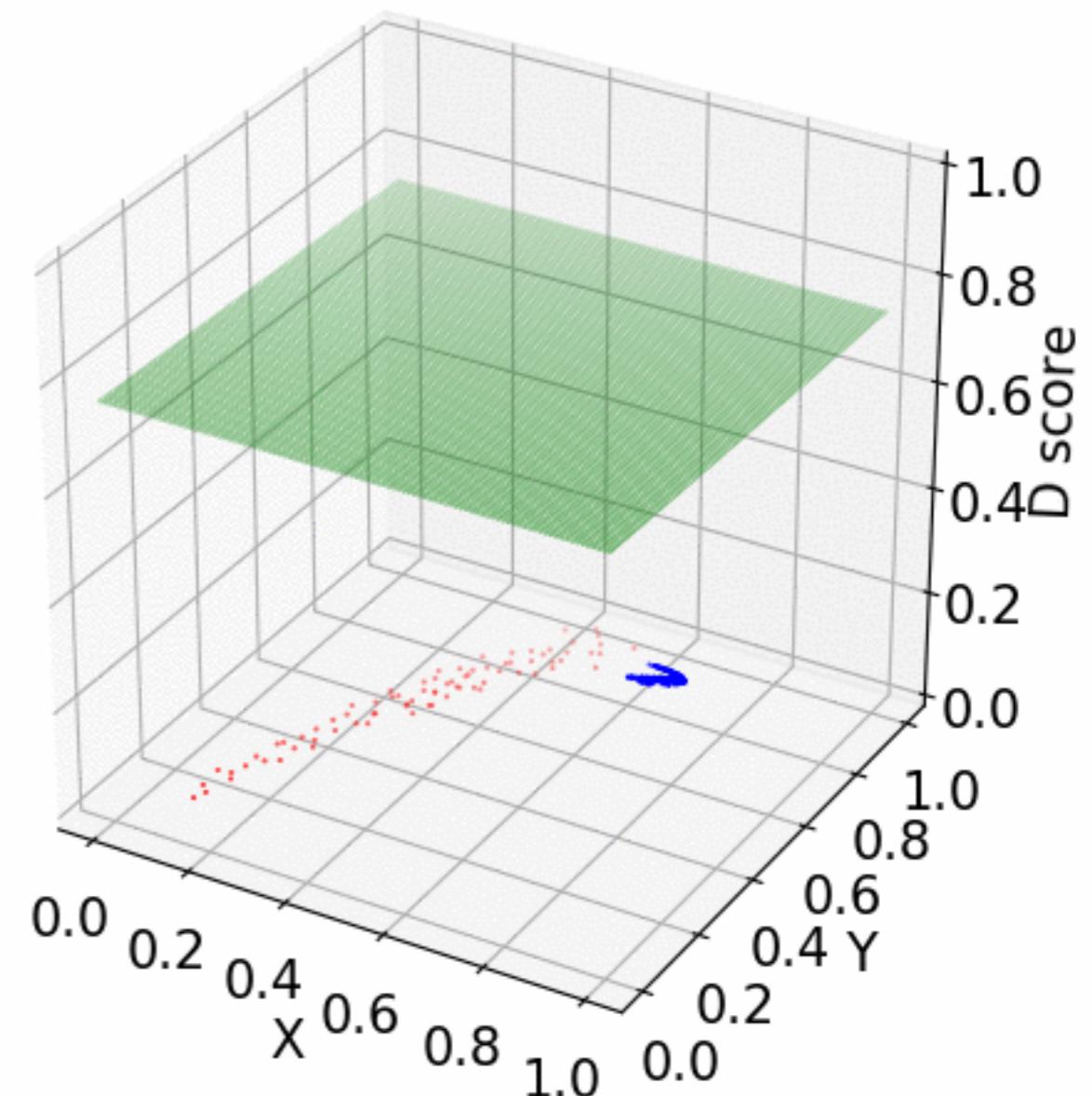
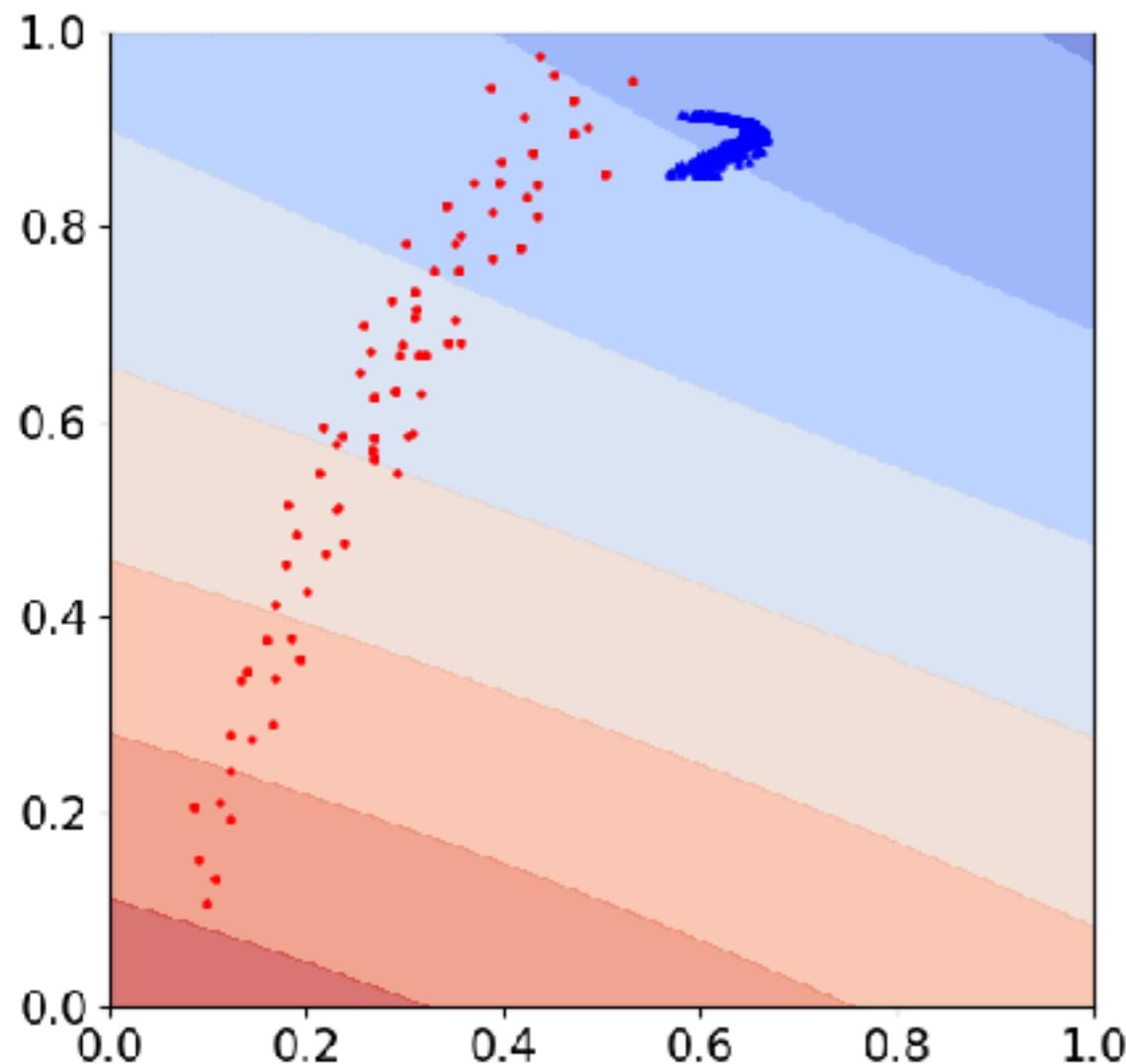
1. Randomly initialise  $\theta$  and  $\mathcal{W}$
2. Use  $g_\theta(z_i)$  to generate “fake” data
3. Keep  $\theta$  constant and optimise  $\mathcal{W}$  for some n iterations
4. Keep  $\mathcal{W}$  constant and optimise  $\theta$  for some k iterations

Discriminator :  $x \rightarrow 10 \rightarrow 10 \rightarrow 10 \rightarrow \text{out}$   
Generator :  $x \rightarrow 5 \rightarrow 5 \rightarrow 5 \rightarrow \text{out}$



Discriminator :  $x \rightarrow 10 \rightarrow 10 \rightarrow 10 \rightarrow \text{out}$

Generator :  $x \rightarrow 5 \rightarrow 5 \rightarrow 5 \rightarrow \text{out}$



# Shortcomings of GAN

When two adversarial of equal strength  
competes both can improve





Discriminator

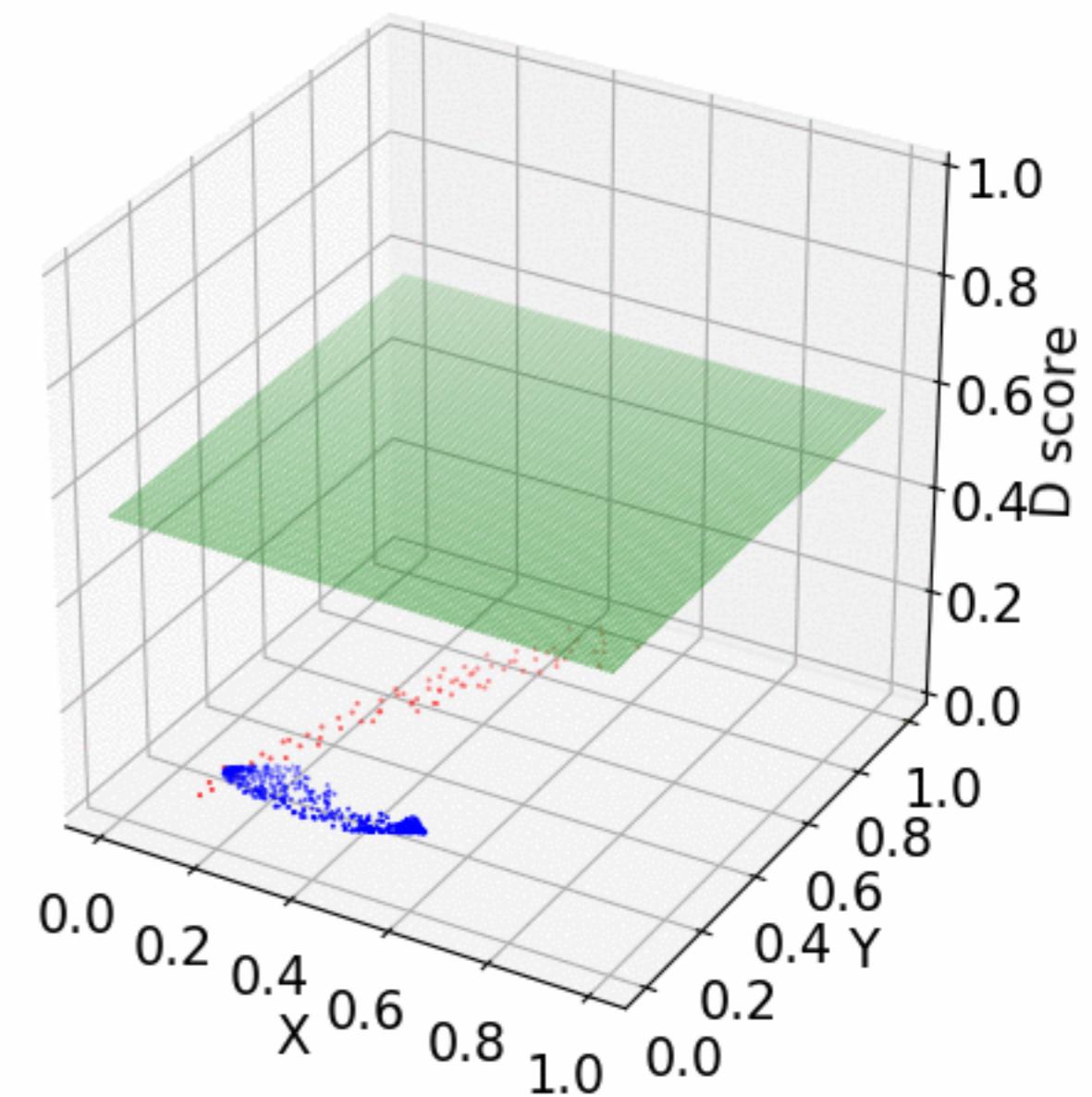
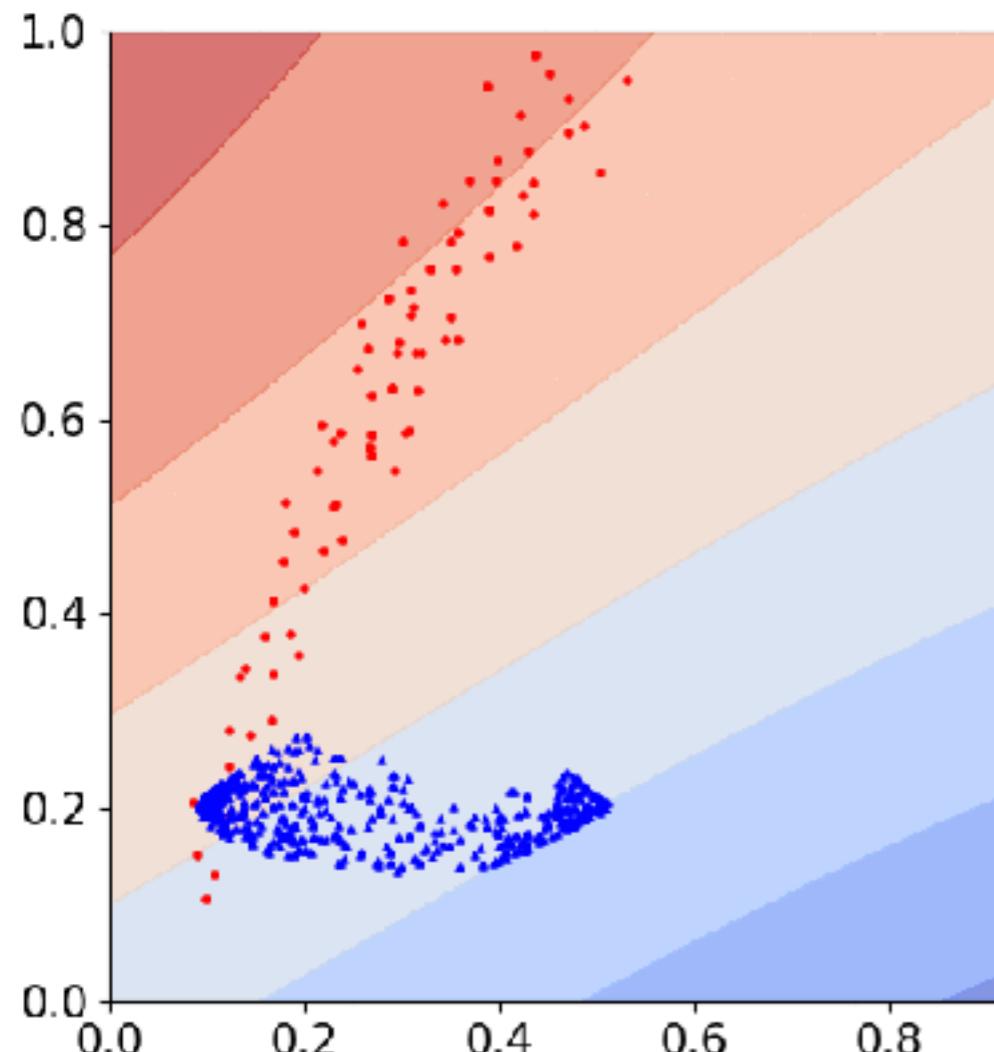
Generator



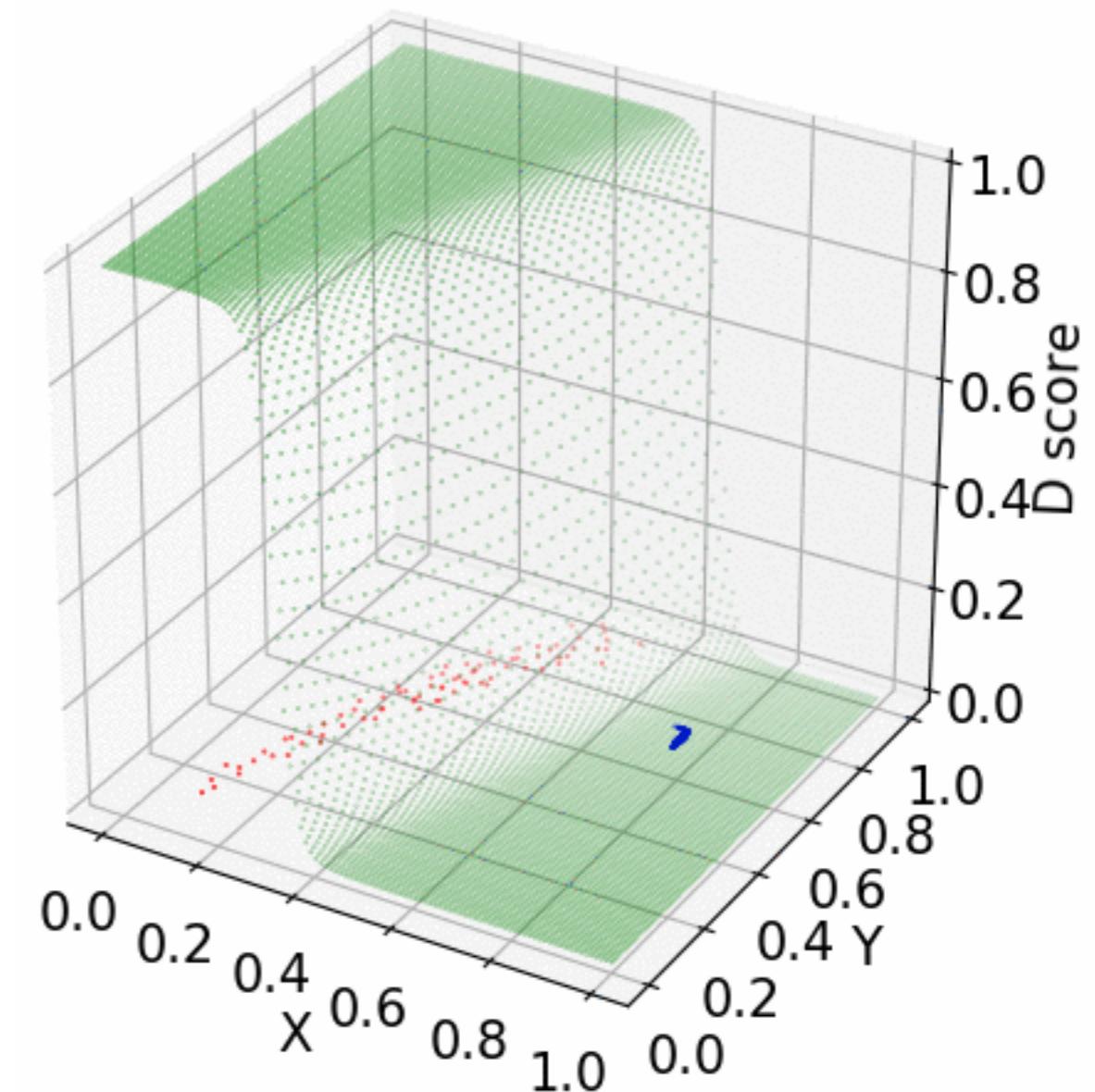
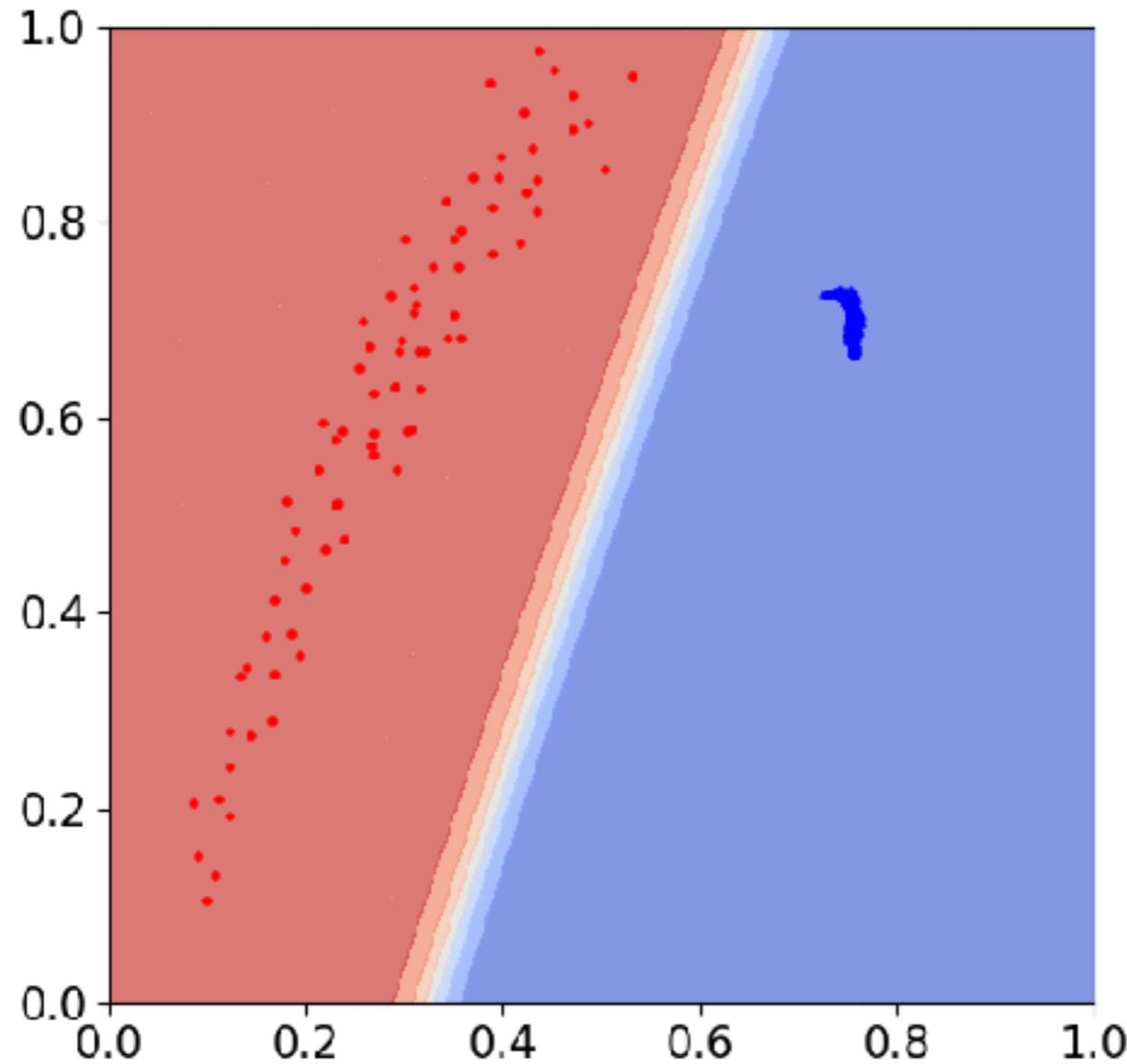
# Discriminator fail

Discriminator :  $x \rightarrow 10 \rightarrow 10 \rightarrow 10 \rightarrow \text{out}$

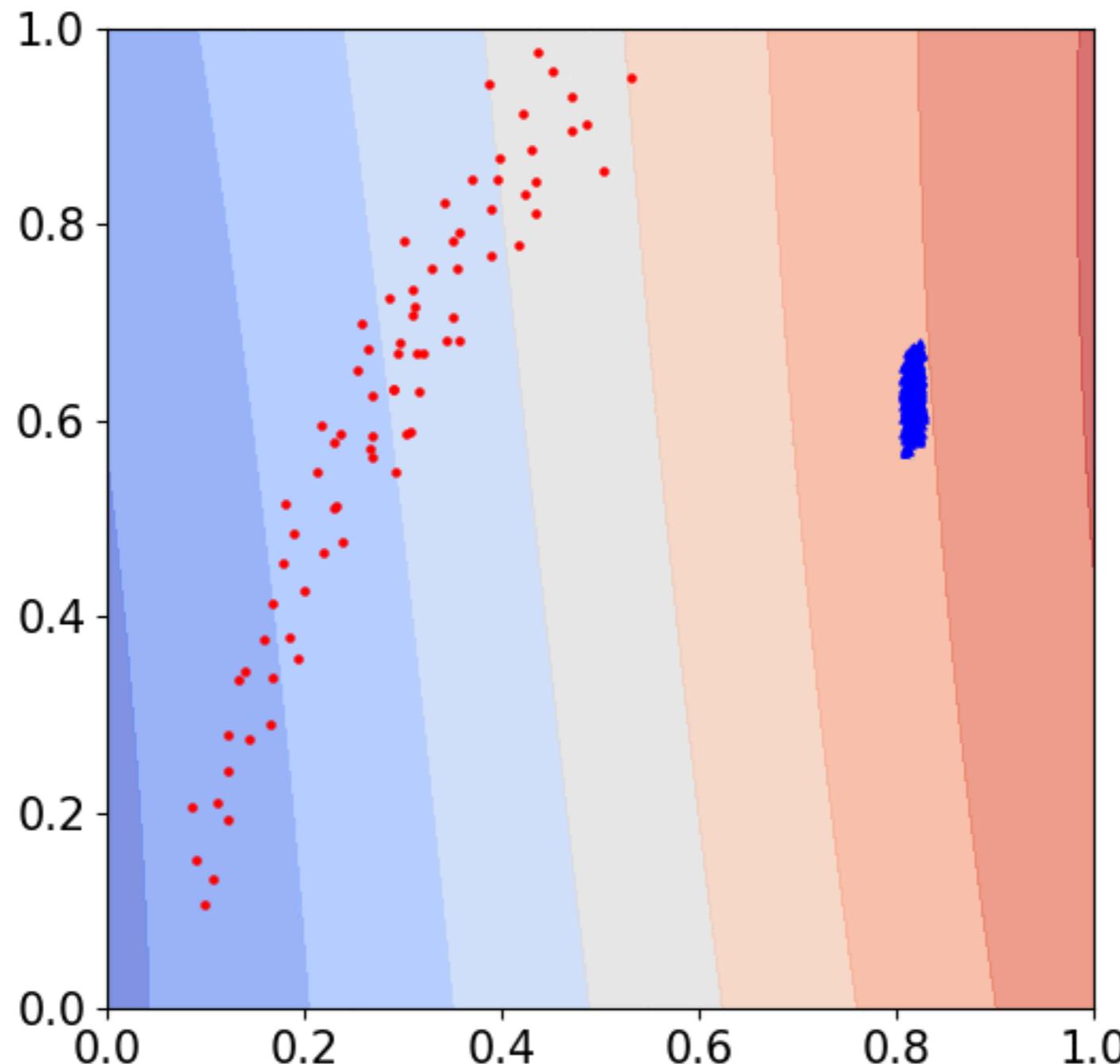
Generator :  $x \rightarrow 5 \rightarrow 5 \rightarrow 5 \rightarrow \text{out}$



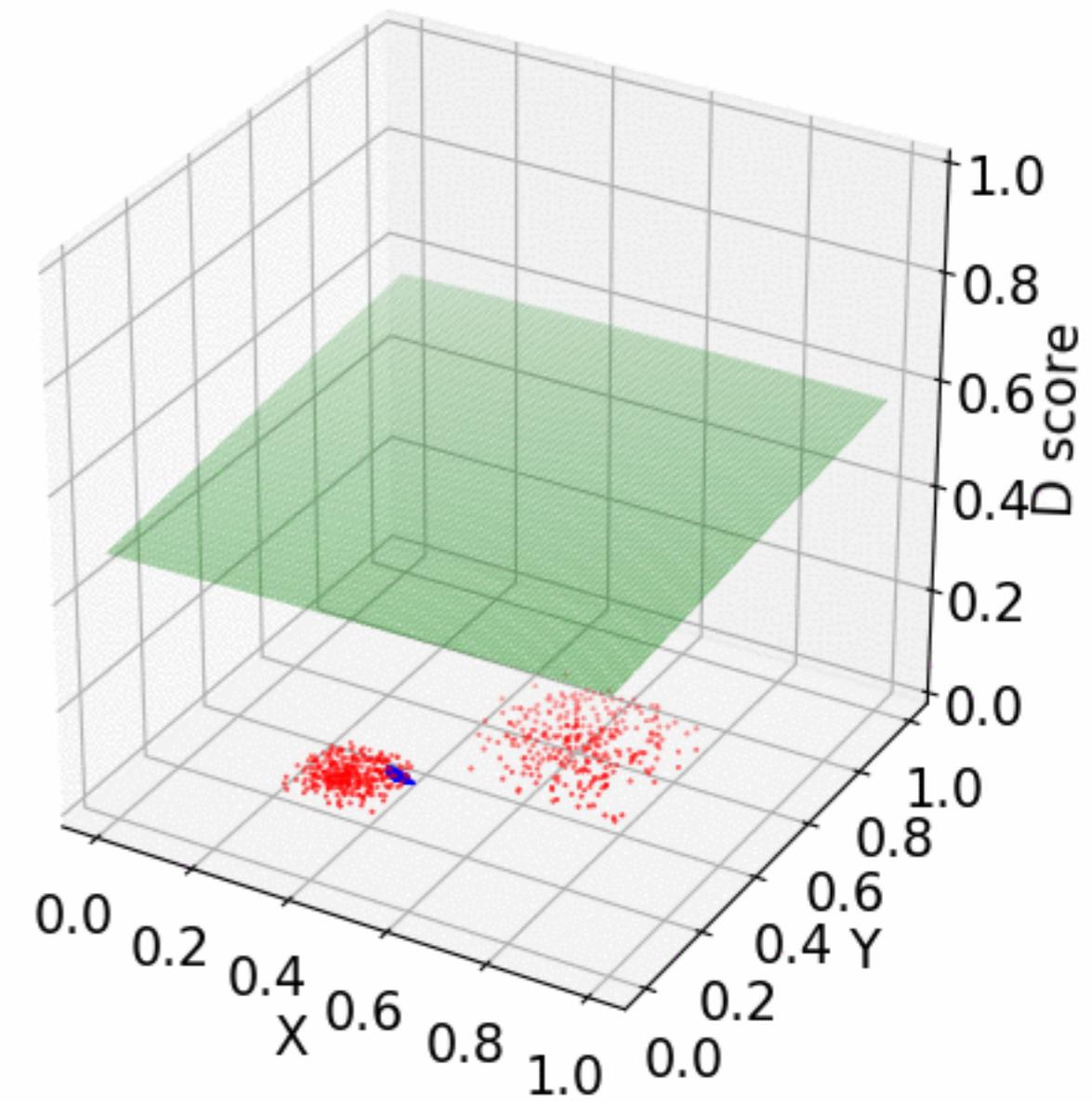
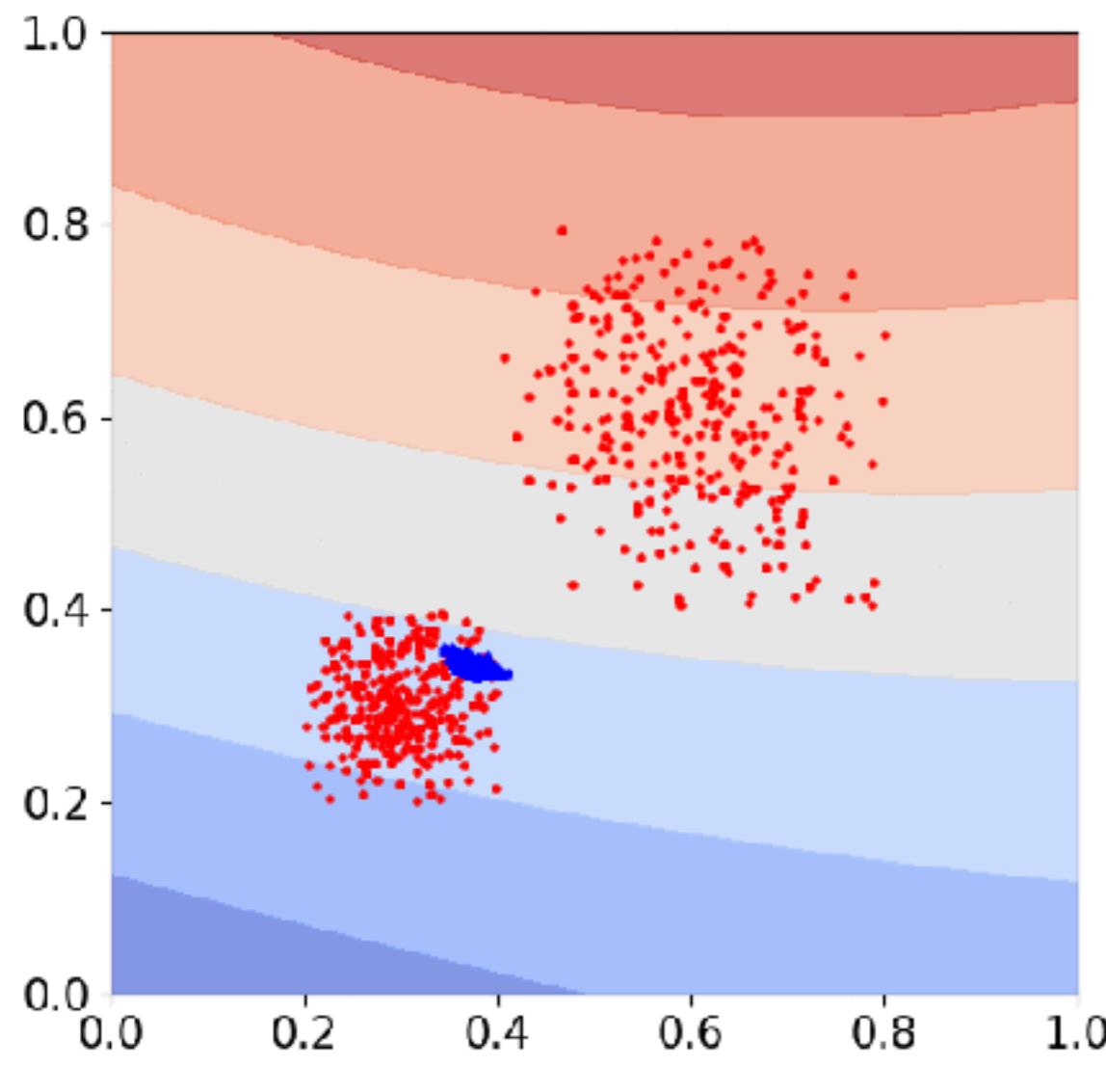
# Generator fail due to zero gradient



# Generator fail due to zero gradient (surrounded)



# Mode collapse



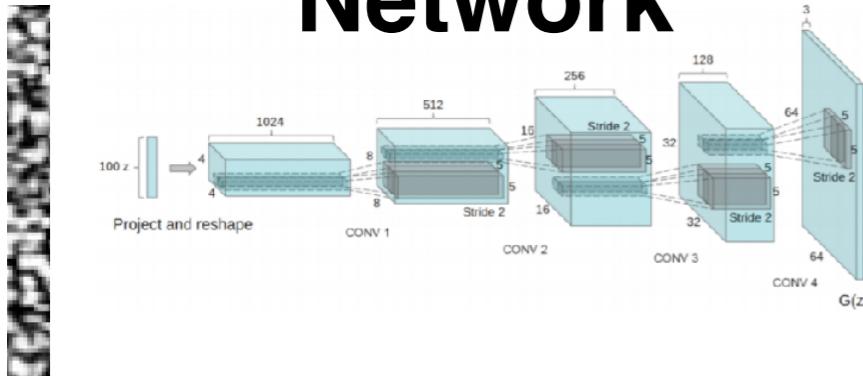
# Improved Techniques for Training GANs

T. Salimans, I Goodfellow *et al*

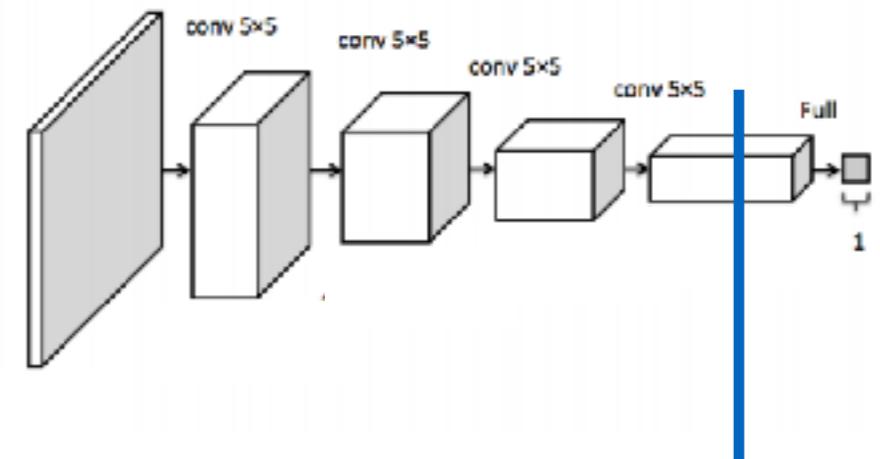
# 1. Feature matching

for training the generator

## Generator Network



## Discriminator Network



$$\vec{f}(x)$$

Cost function for generator

$$C_g(w, \theta) = \left\| \frac{1}{N} \sum_i \vec{f}_w(x_i) - \frac{1}{M} \sum_j \vec{f}_w(g_\theta(z_j)) \right\|_2^2$$

$$\theta^* = \arg \min_{\theta} C_g(w, \theta)$$

# Discriminator cost function use in original GAN paper

$$C(w, \theta) = - \left( \frac{1}{N^+} \sum_i \log(D_w(x_i)) + \frac{1}{N^-} \sum_j \log(1 - D_w(g_\theta(z_j))) \right)$$

$$\theta^* = \arg \max_{\theta} C(w, \theta) \quad \text{keeping } w \text{ constant}$$

## 2. Minibatch discrimination

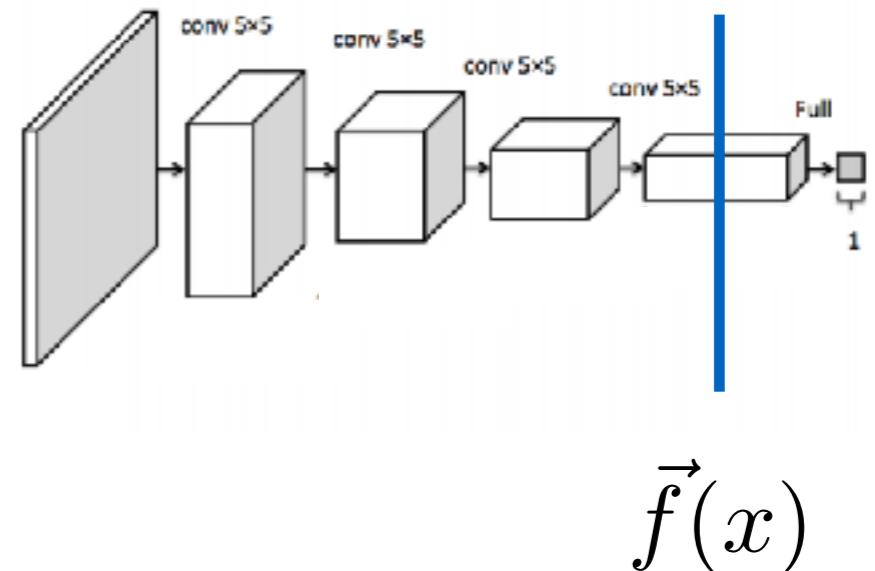
for training the discriminator

$$\vec{f}(x_i) \rightarrow \mathbf{T} \rightarrow \vec{o}(x_i)$$

Read paper for details

$\vec{o}$  contains information  
on how this data point  
cluster with other data  
points in the minibatch

### Discriminator Network



Concatenate features to create a longer feature

$$(\vec{f}(x_i), \vec{o}(x_i))$$

Put this new feature back into subsequent layers of discriminator

### 3. Historical averaging

Add a term in cost function

$$\left\| \theta - \frac{1}{t} \sum_i \theta_i \right\|^2$$

Effect, at every step, theta does not change too much in magnitude and direction

## 4. One-sided label smoothing

Changes in discriminator training

Instead of asking discriminator to output ‘1’ for real data, ask discriminator to output a softer value like ‘0.9’ for real data.

Fake data still has label ‘0’

Encourage a softer decision boundary, e.g. don’t need to saturate out the sigmoid and make gradient zero

## 5. Virtual batch normalisation

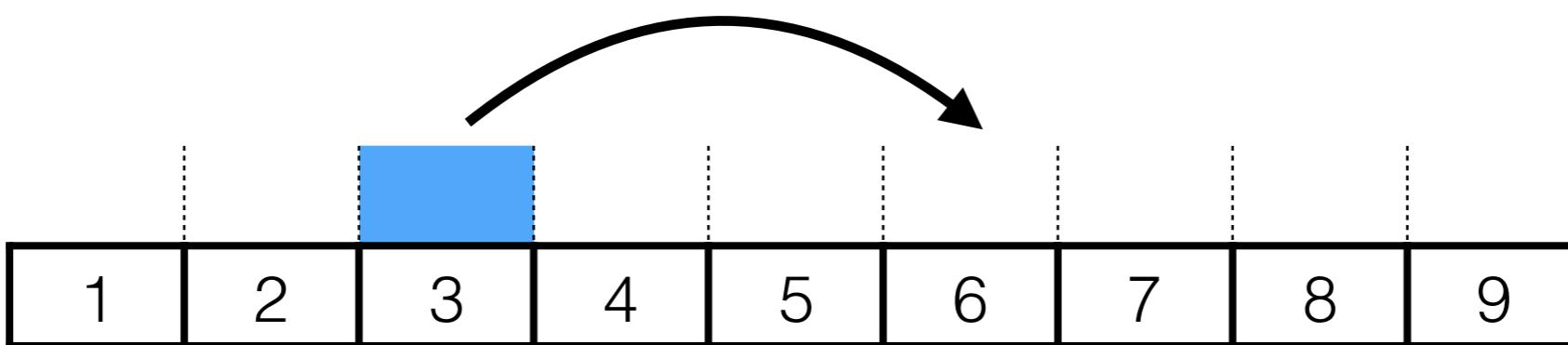
Instead of normalising differently for each minibatch, use a reference batch so that normalisation is consistent across minibatches.

What I would do is to use one statistics for normalisation across all data, e.g. compute z-score using one mean and one variance

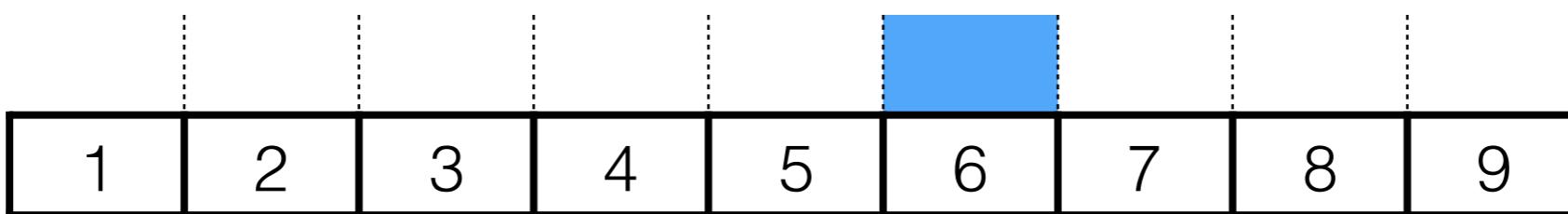
What paper proposed is to recalculate mean and variance for each data point, which will increase computational time a lot

# Wasserstein GAN

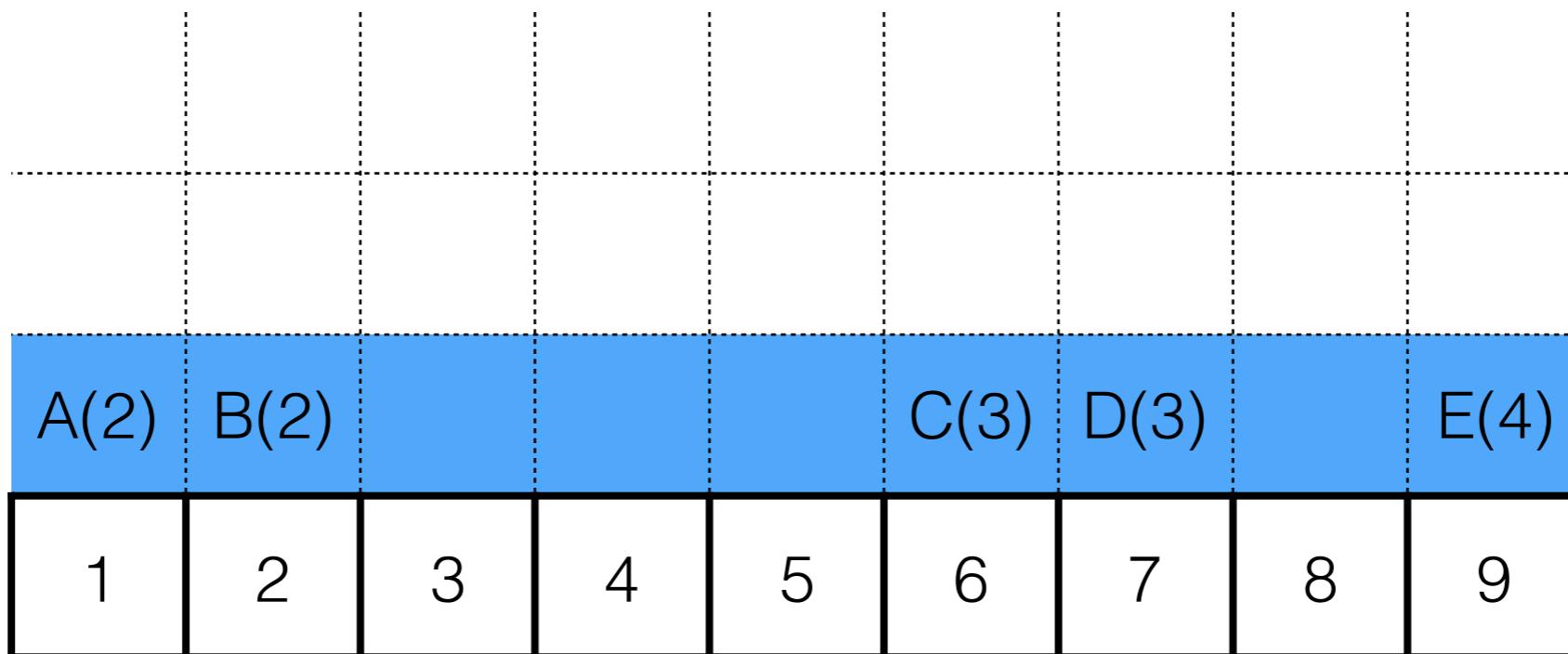
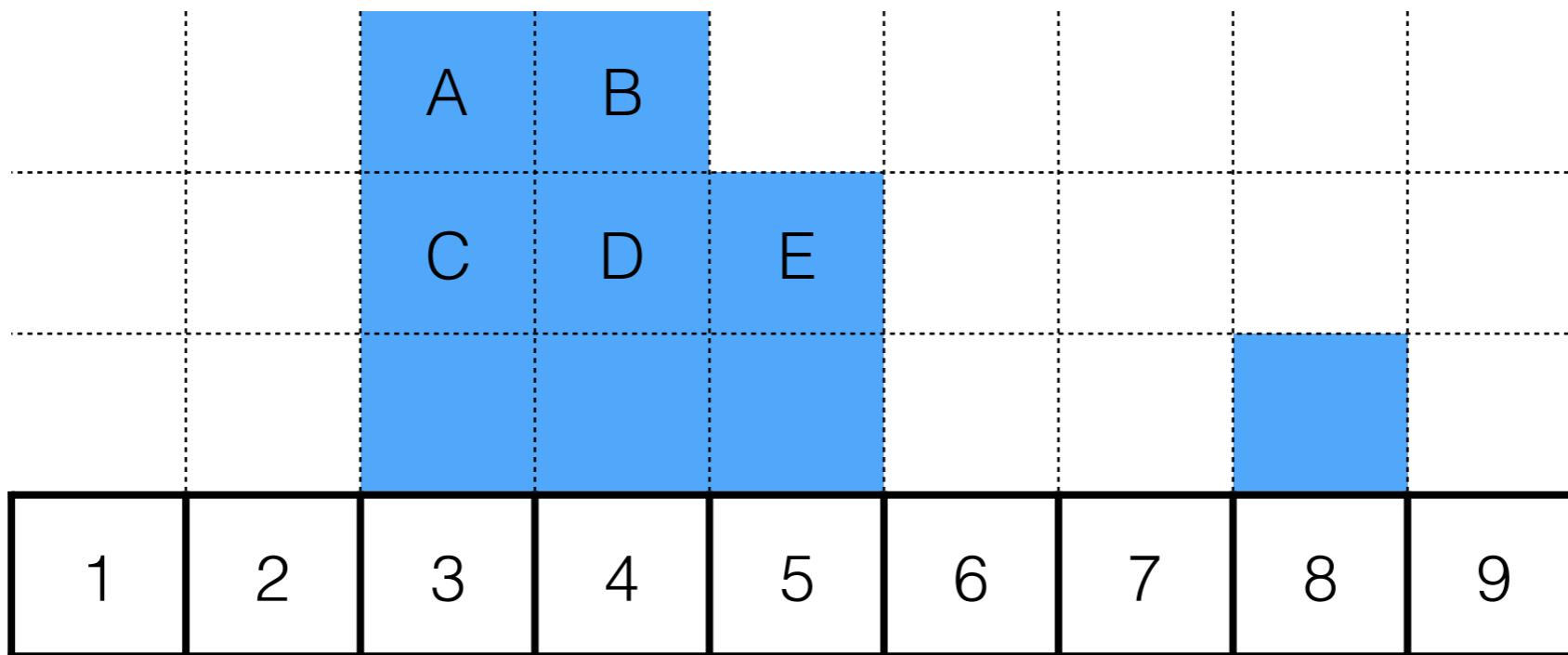
# Earth Mover's Distance



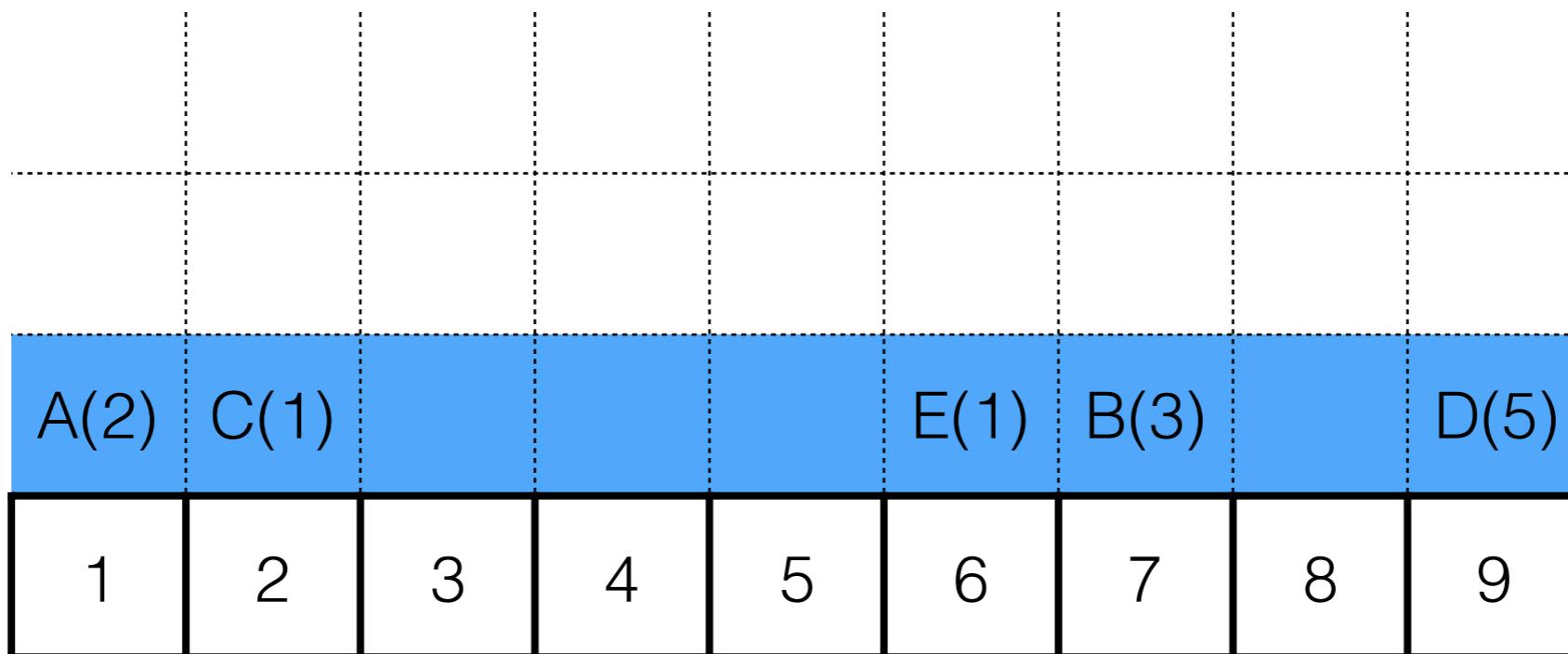
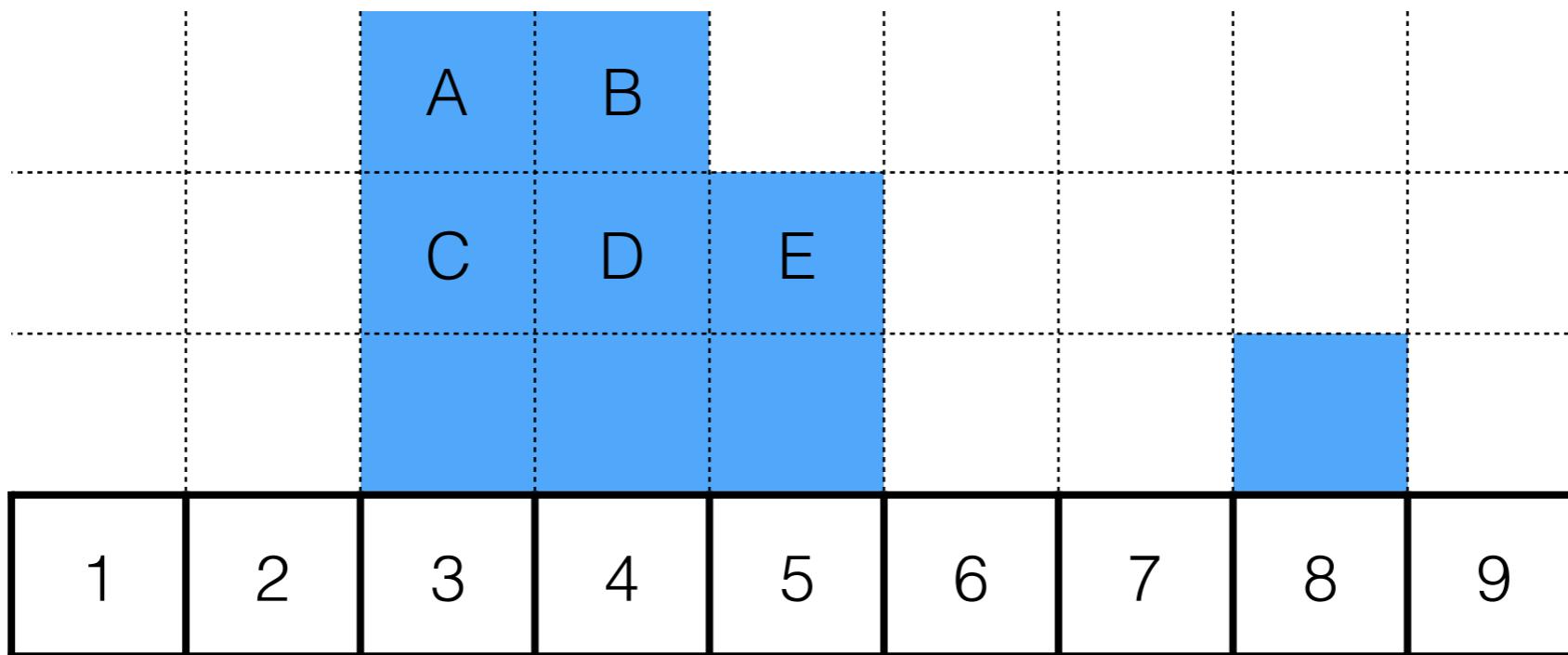
cost = 3



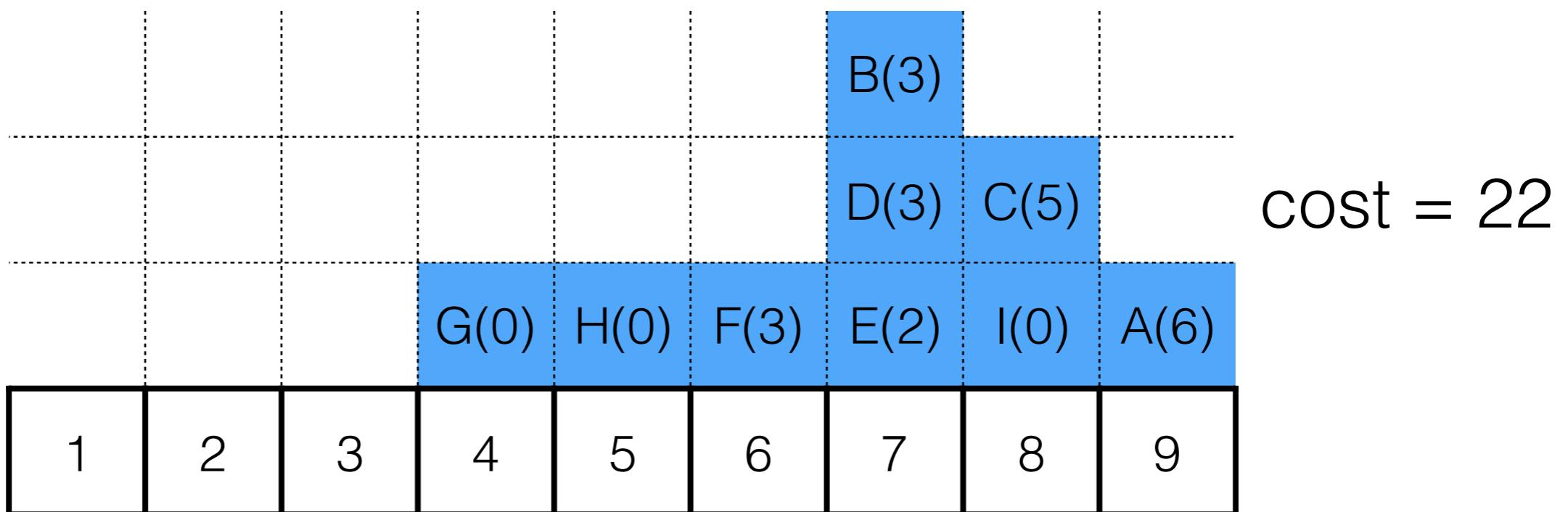
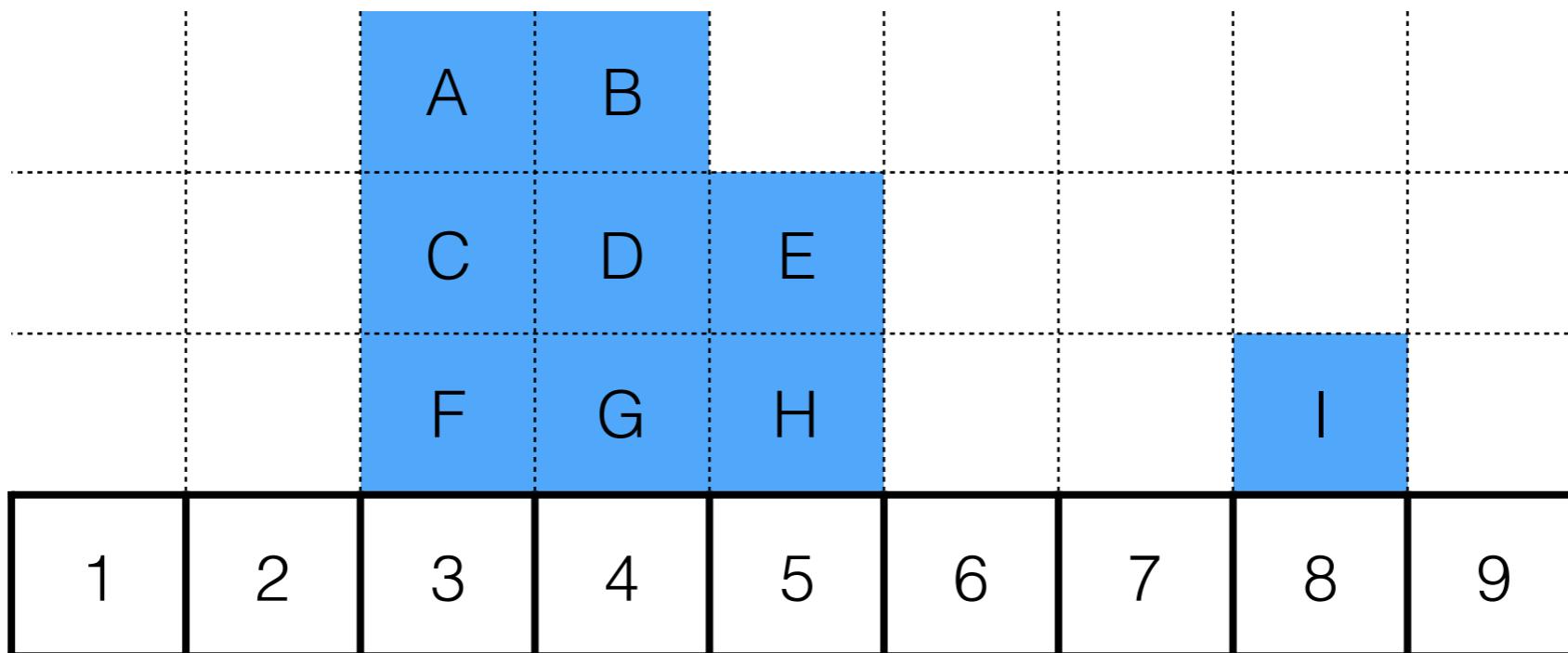
# Earth Mover's Distance



# Earth Mover's Distance



# Earth Mover's Distance



## Distribution of real and fake data

$P_r(x)$  is the distribution of real data, it is constant

$P_\theta(g_\theta(z))$  is the distribution of generated data  
it changes when generator change

$$W(P_r(x), P_\theta(g_\theta(z)))$$

is the Earth Mover's Distance between  
real and generated data

## Idea of Wasserstein GAN

No discriminator (in the paper they use the word ‘critic’)

Given  $P_r(x)$

Initialize  $g_\theta(z_i)$  and generate  $P_\theta(g_\theta(z))$

Adjust generator such that

$W(P_r(x), P_\theta(g_\theta(z)))$  is minimised

# Kantorovich-Rubinstein duality

*I did not prove it*

$$W(P_r, P_\theta) = \sup_{\|f\| < 1} \left( \frac{1}{N} \sum_i f(x_i) - \frac{1}{M} \sum_j f(g_\theta(z_j)) \right)$$

over 1-Lipschitz functions  
(loosely speaking gradient does not exceed one)

Sample from  $P_r(x)$

Sample from  $P_\theta(g_\theta(z))$

The diagram illustrates the Kantorovich-Rubinstein duality formula. The formula is enclosed in a large black bracket. Three arrows point upwards from the text below to the corresponding terms in the formula: one arrow points to the first term in the sum over N, another to the first term in the sum over M, and a third to the entire term involving  $g_\theta$ .

# Kantorovich-Rubinstein duality

*I did not prove it*

This is a constrained optimisation problem!

# Representation of 1-Lipschitz function

One way to parameterise a function is to use a neural network. By adjusting weights, we create different functions.

When network is small, class of function obtained is limited.

When network is deep, we can represent more functions.

$$W(P_r, P_\theta) = \sup_{\|f\| < 1} \left( \frac{1}{N} \sum_i f(x_i) - \frac{1}{M} \sum_j f(g(z_j)) \right)$$

Represent  $f(\cdot)$  by a neural network  $f_w(\cdot)$

Use gradient ascend to find  $W(P_r(x), P_\theta(g_\theta(z)))$

# Algorithm

Initialize  $w$  and  $\theta$

Generates  $P_\theta(g_\theta(z))$

(1) Maximize by adjusting  $w$

$$W(P_r, P_\theta) = \sup_{\|f\| < 1} \left( \frac{1}{N} \sum_i f_w(x_i) - \frac{1}{M} \sum_j f_w(g(z_j)) \right)$$

(2) Minimize Earth Mover's distance

$$\arg \min_{\theta} W(P_r, P_\theta)$$

(3) Repeat (1) & (2)

---

**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
```

---

---

**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.

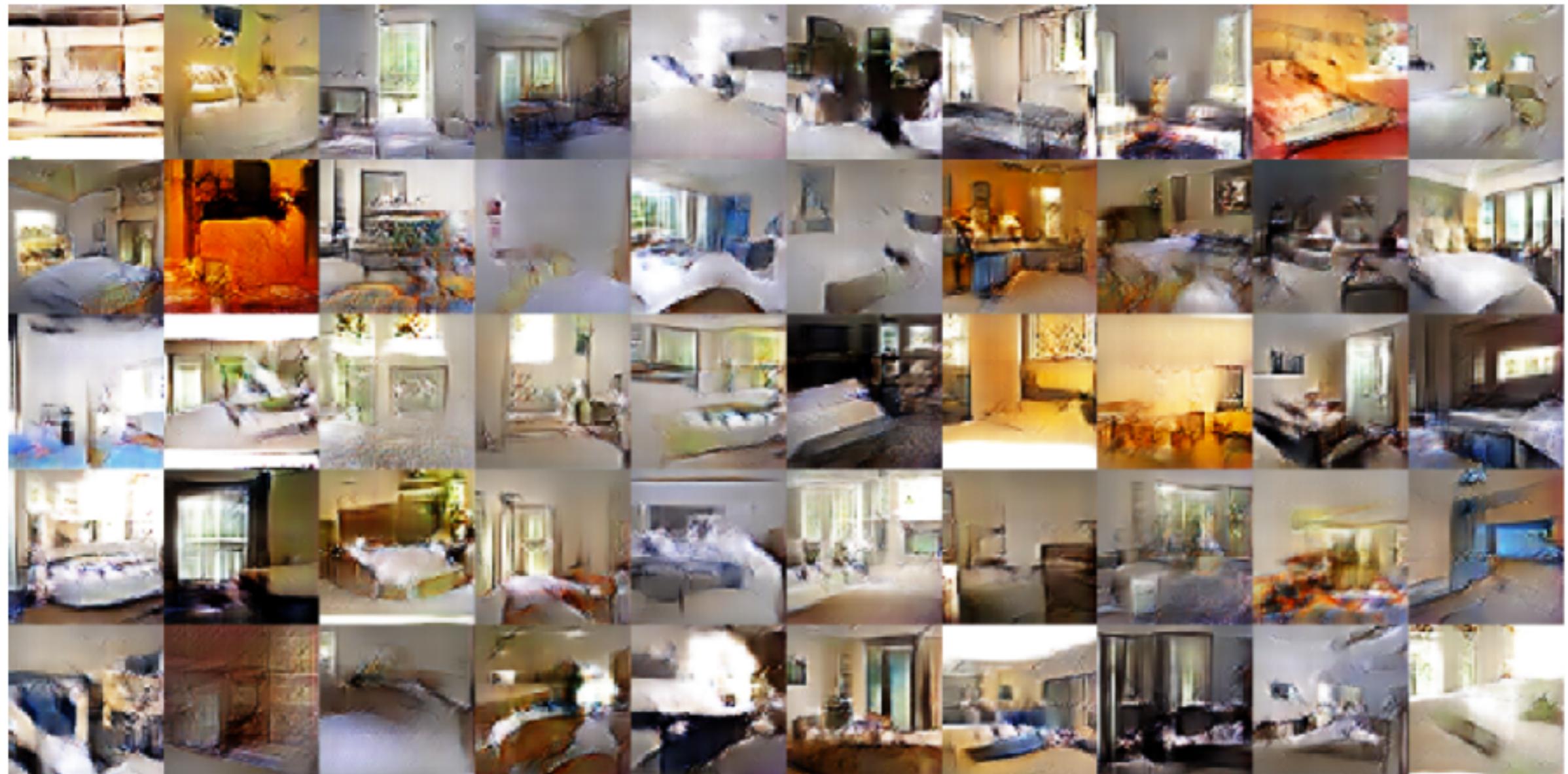
- 1: **while**  $\theta$  has not converged **do**
- 2:   **for**  $t = 0, \dots, n_{\text{critic}}$  **do**
- 3:     Sample  $\{x^{(i)}\}_{i=1}^m \sim \mathbb{P}_r$  a batch from the real data.
- 4:     Sample  $\{z^{(i)}\}_{i=1}^m \sim p(z)$  a batch of prior samples.
- 5:      $g_w \leftarrow \nabla_w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

---

**Assignment #3: There is something strange about this algorithm. What is it?**

# Deep Convolutional GAN

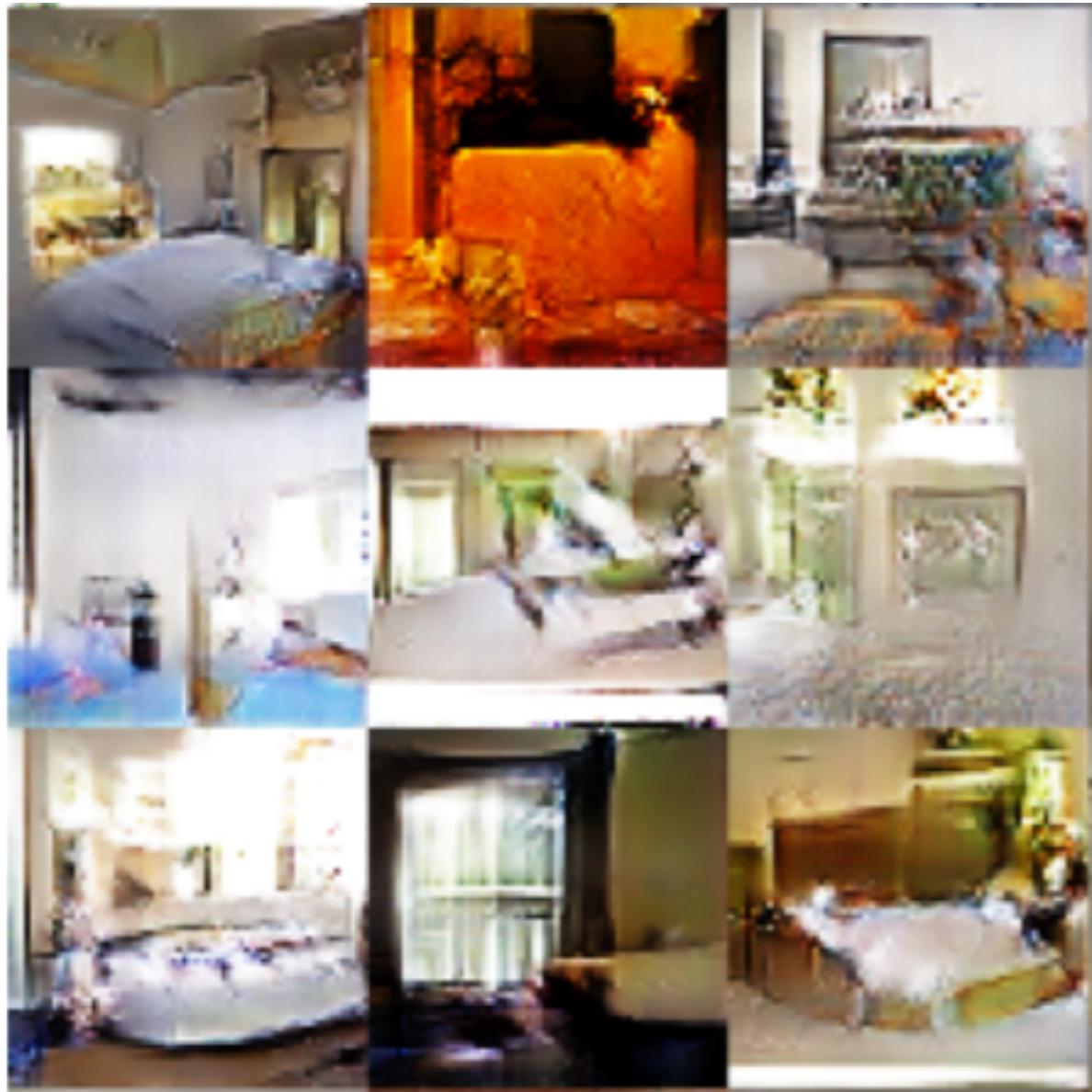
A. Radford, L. Metz, S. Chintala  
Unsupervised Representation Learning with Deep Convolutional GAN, ICLR2016



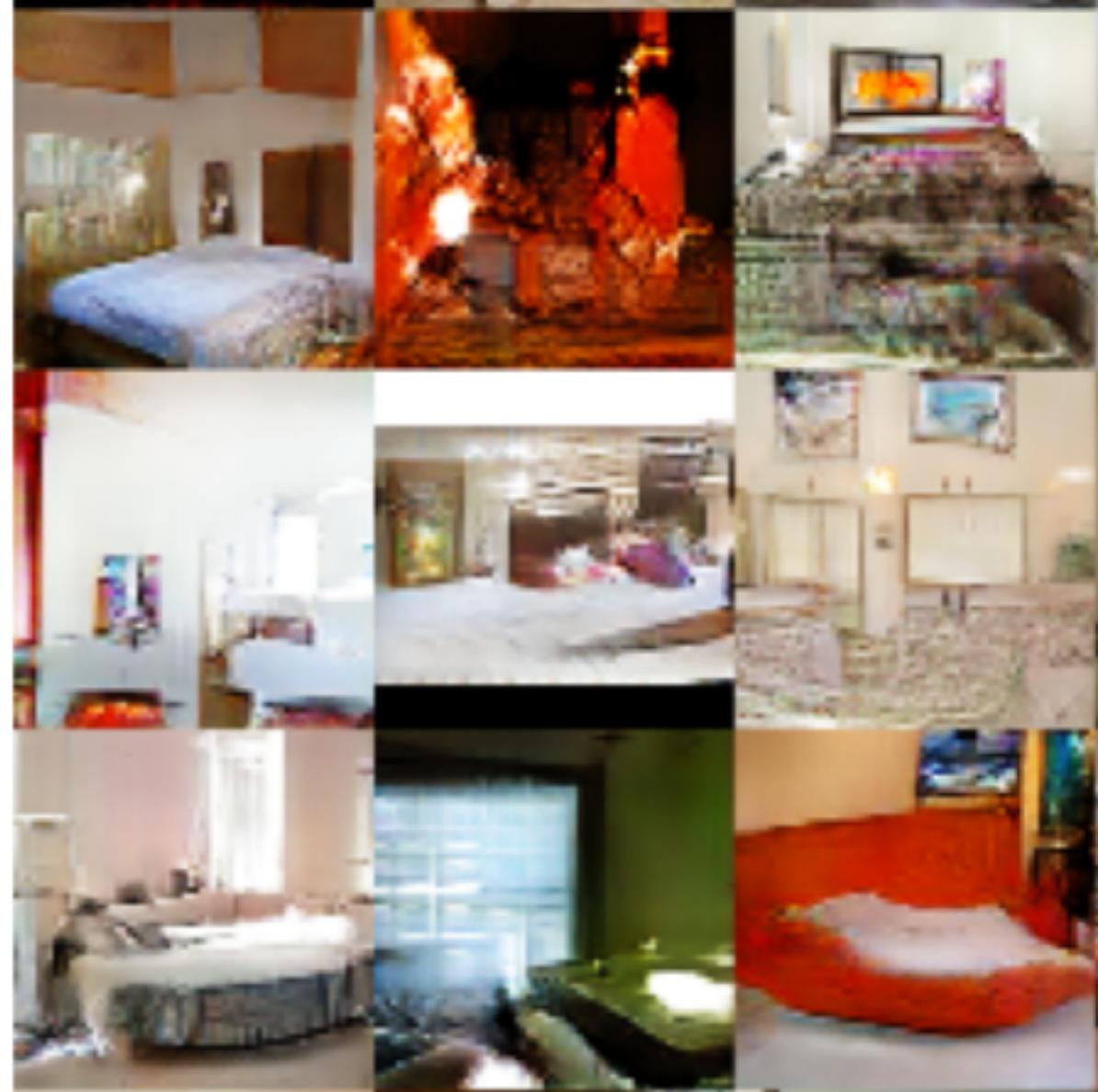
one epoch



5 epoch



one epoch



5 epoch

# Generator

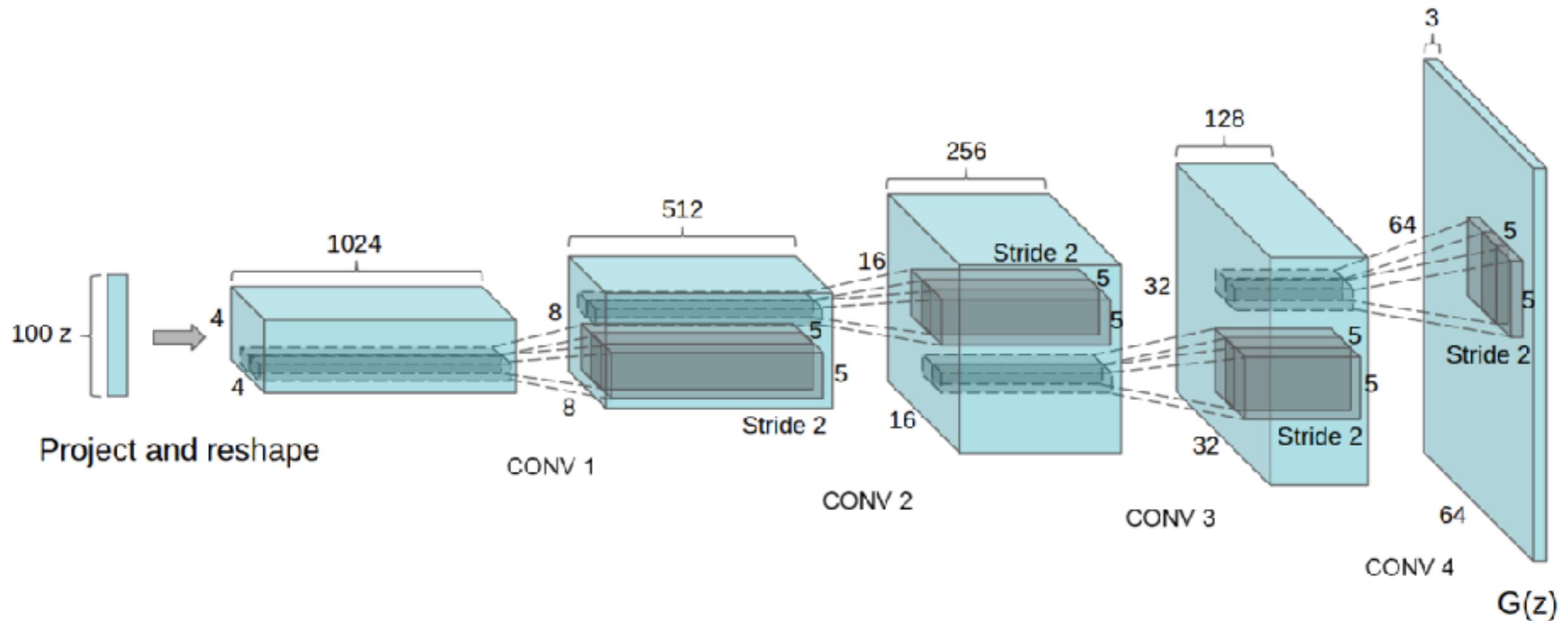
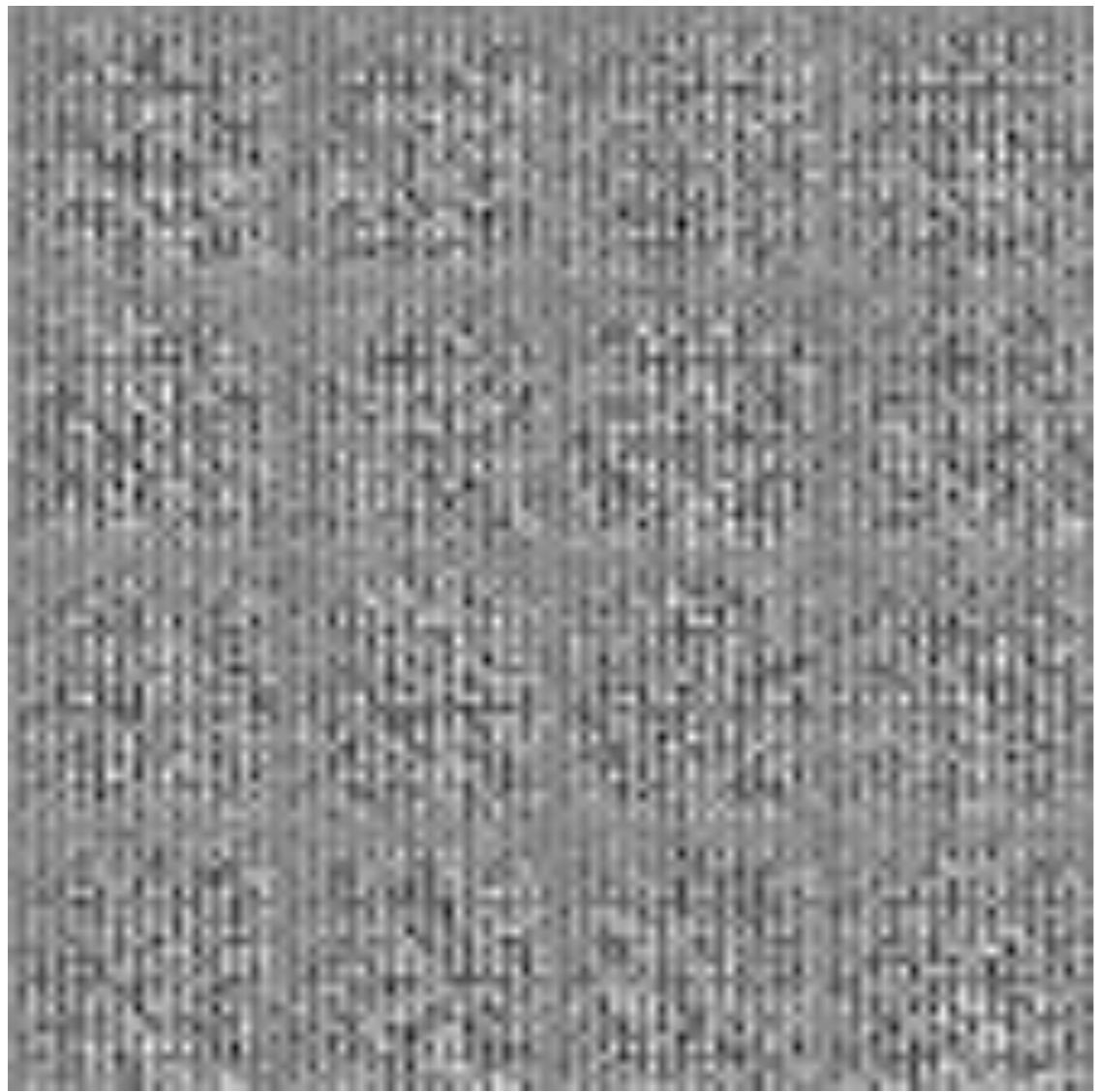


Figure 1: DCGAN generator used for LSUN scene modeling. A 100 dimensional uniform distribution  $Z$  is projected to a small spatial extent convolutional representation with many feature maps. A series of four fractionally-strided convolutions (in some recent papers, these are wrongly called deconvolutions) then convert this high level representation into a  $64 \times 64$  pixel image. Notably, no fully connected or pooling layers are used.

## Architecture guidelines for stable Deep Convolutional GANs

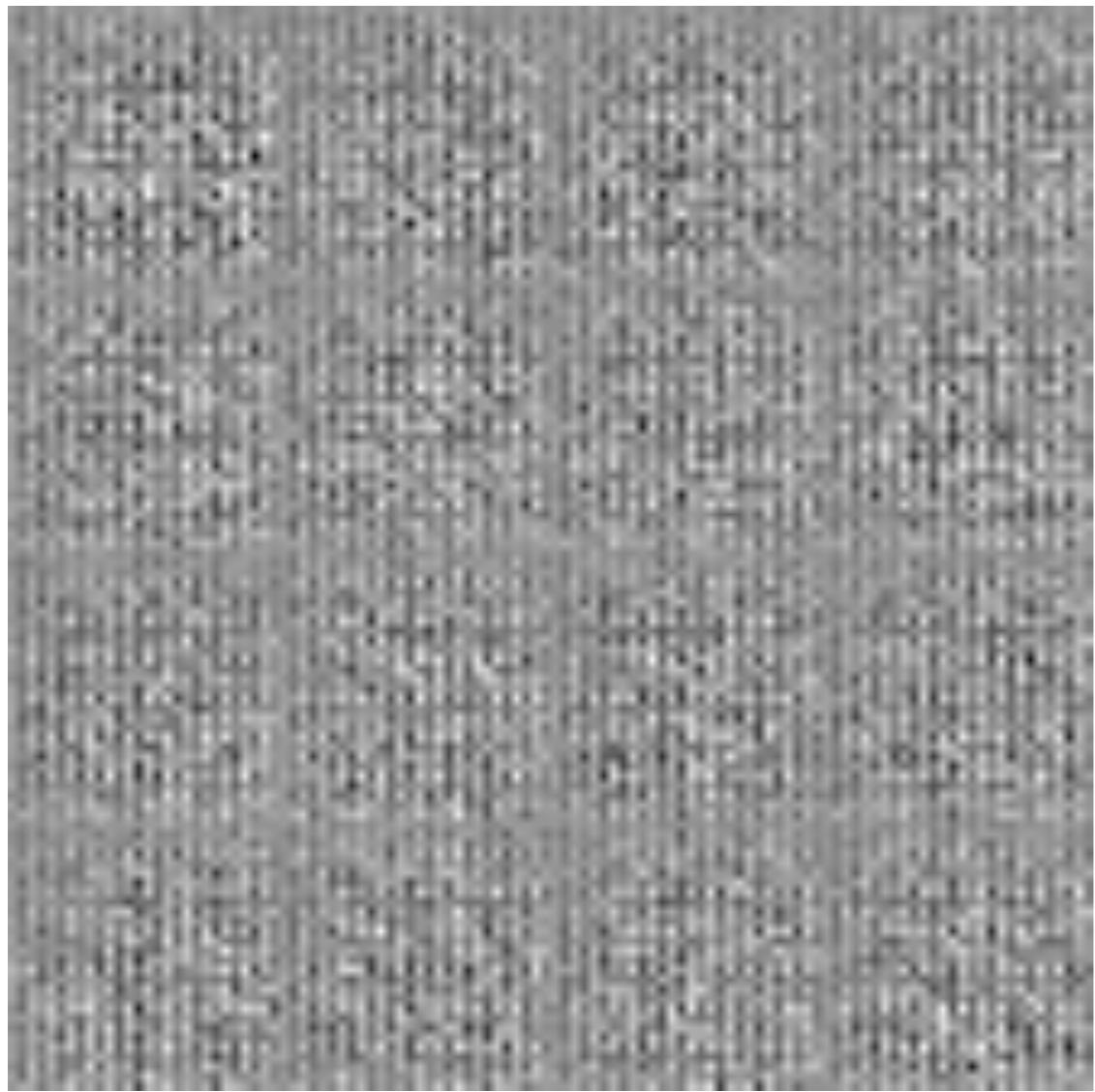
- Replace any pooling layers with strided convolutions (discriminator) and fractional-strided convolutions (generator).
- Use batchnorm in both the generator and the discriminator.
- Remove fully connected hidden layers for deeper architectures.
- Use ReLU activation in generator for all layers except for the output, which uses Tanh.
- Use LeakyReLU activation in the discriminator for all layers.



Every 10 iterations



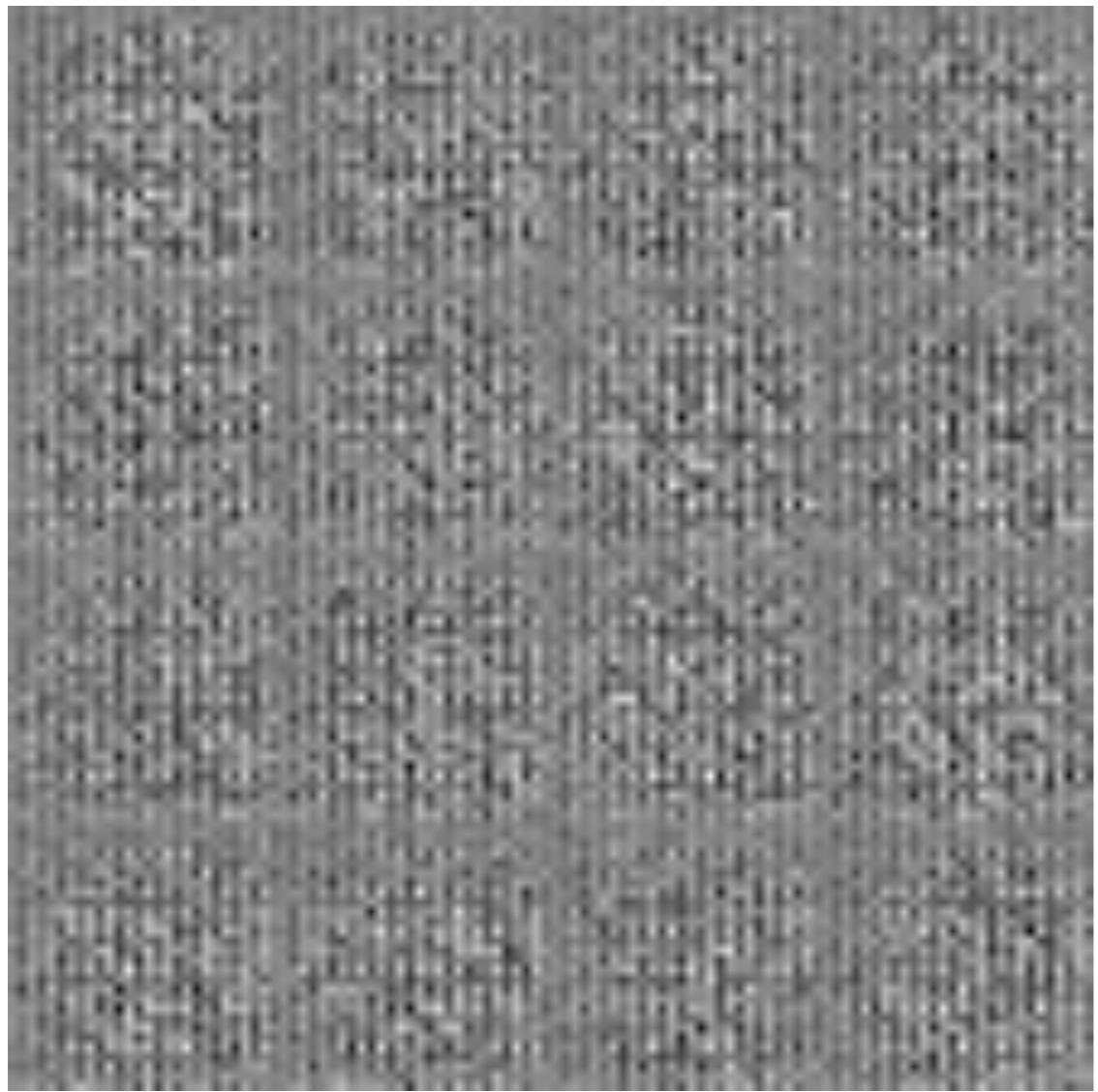
Final generated digits



Every 10 iterations



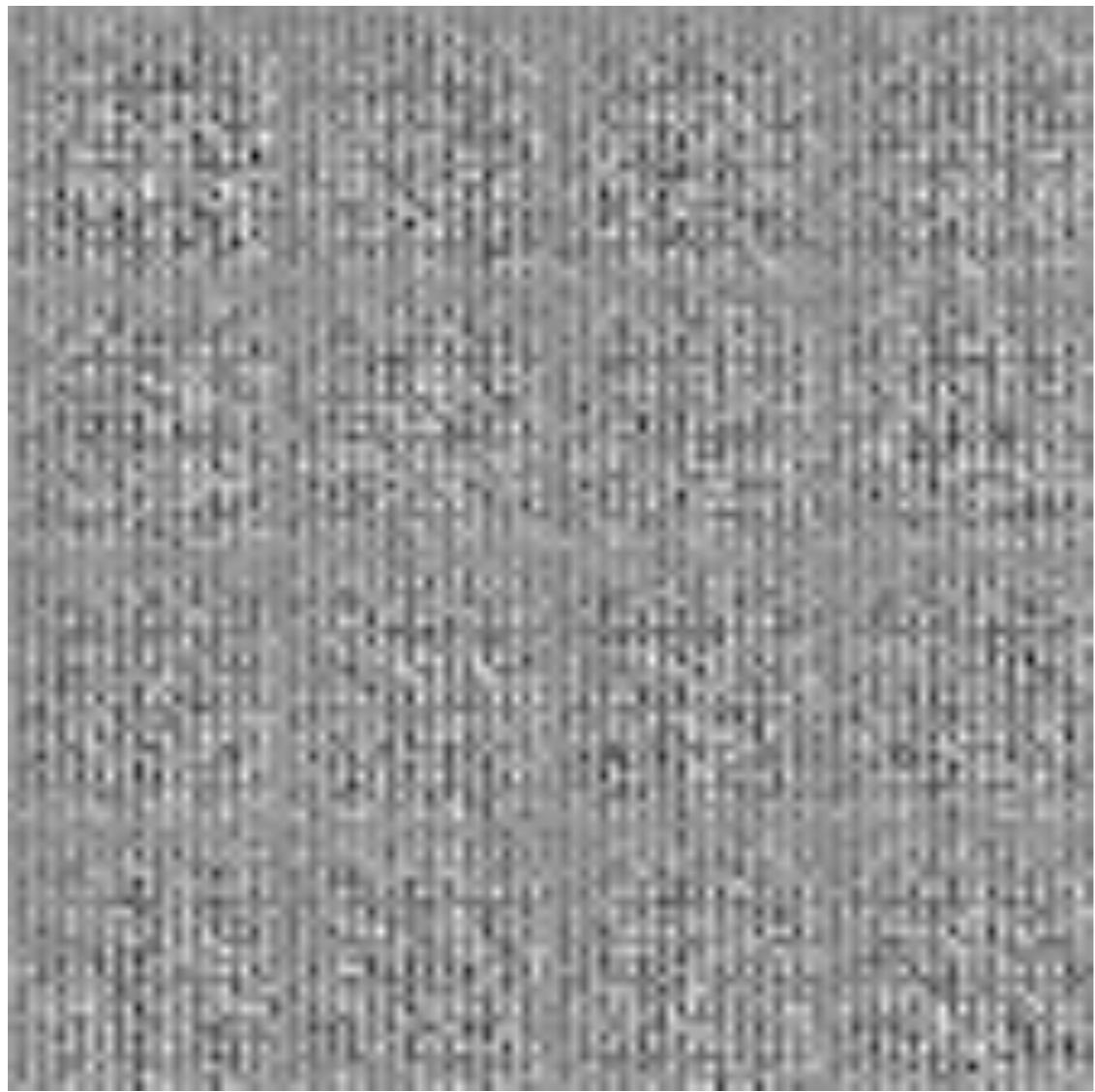
Final generated digits



Every 10 iterations



Final generated digits



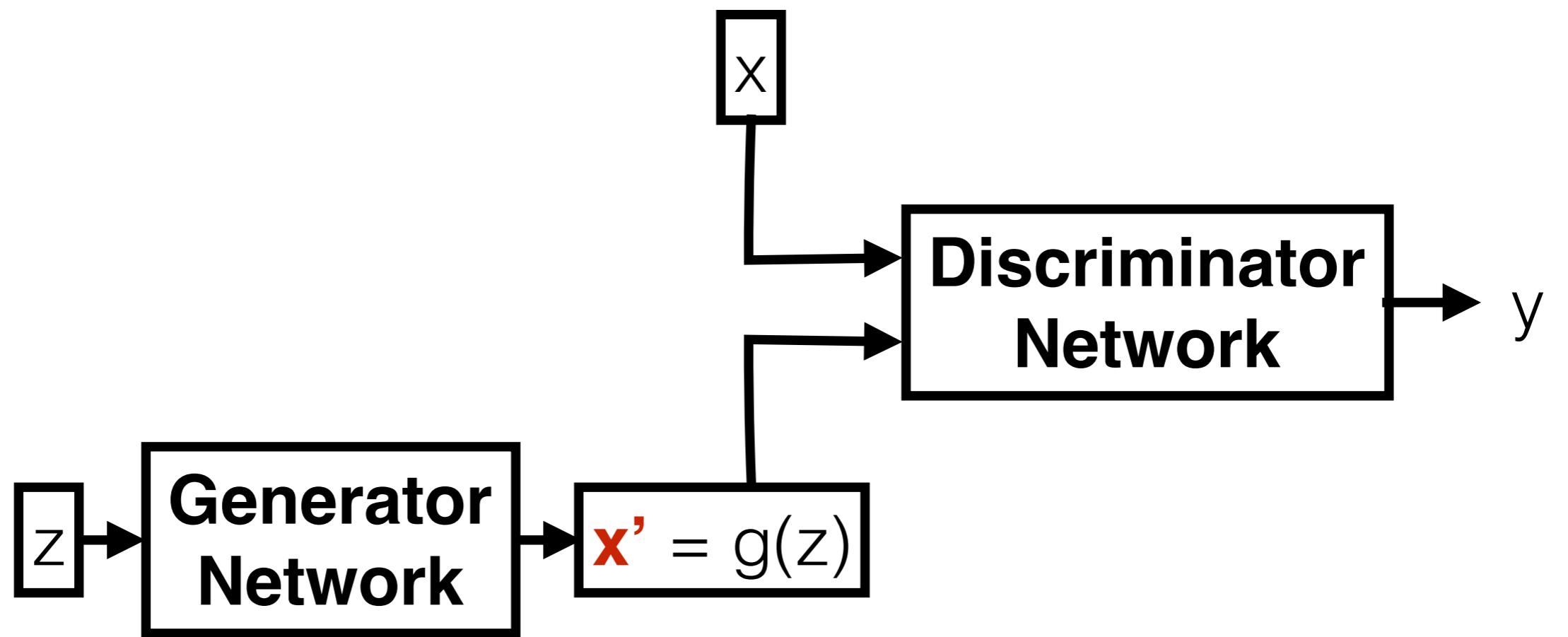
Every 10 iterations



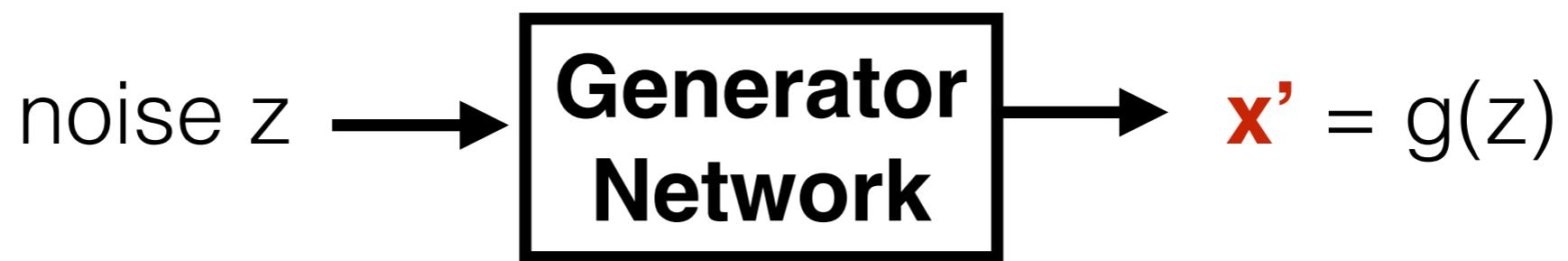
Final generated digits

BiGAN

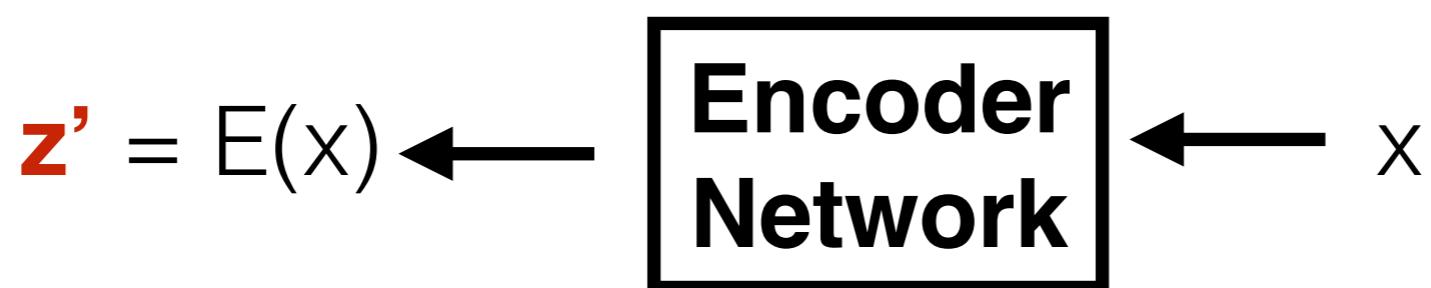
Bidirectional GAN

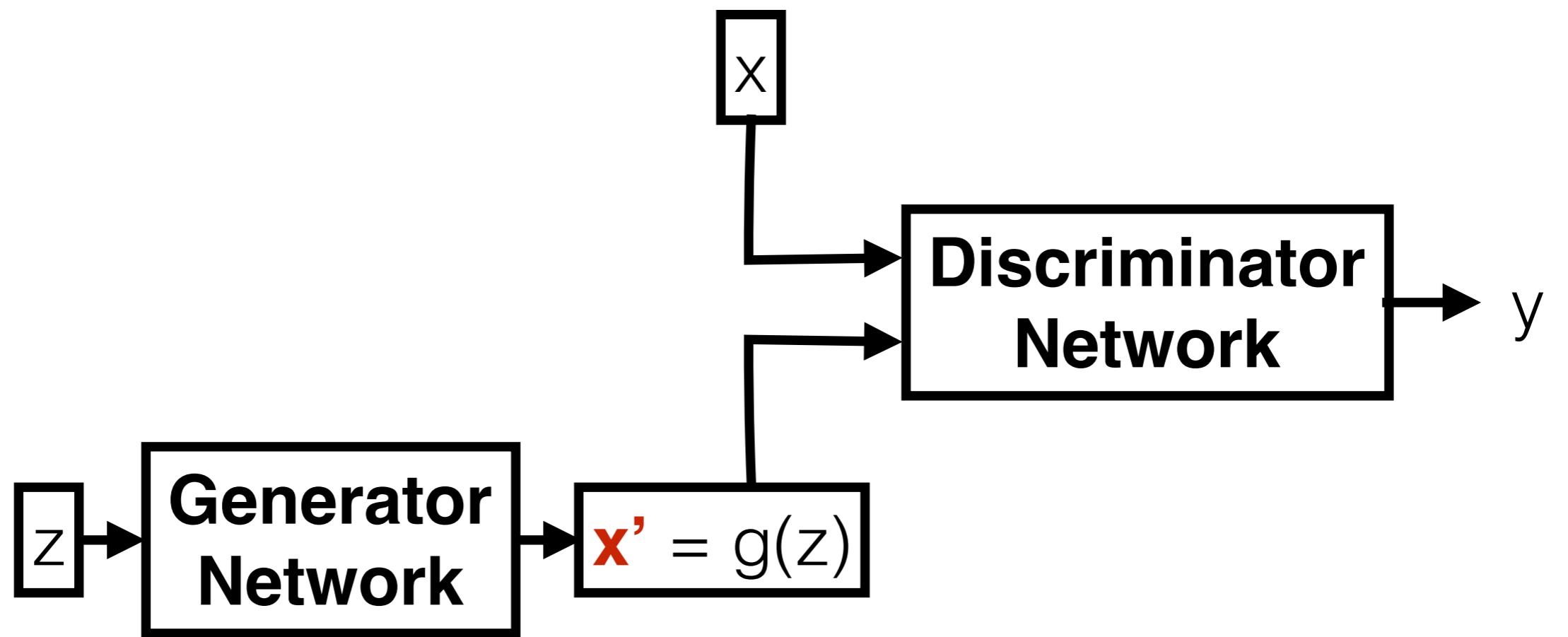


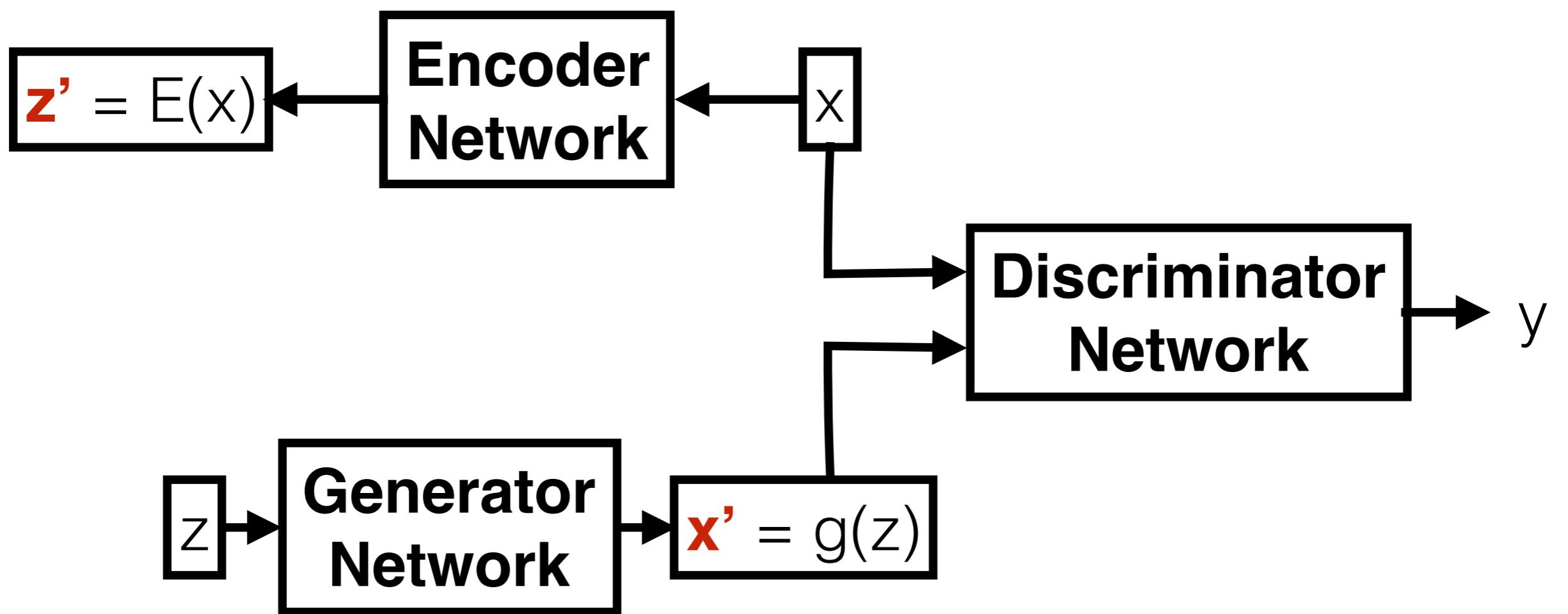
Give  $z$ , generate  $x'$

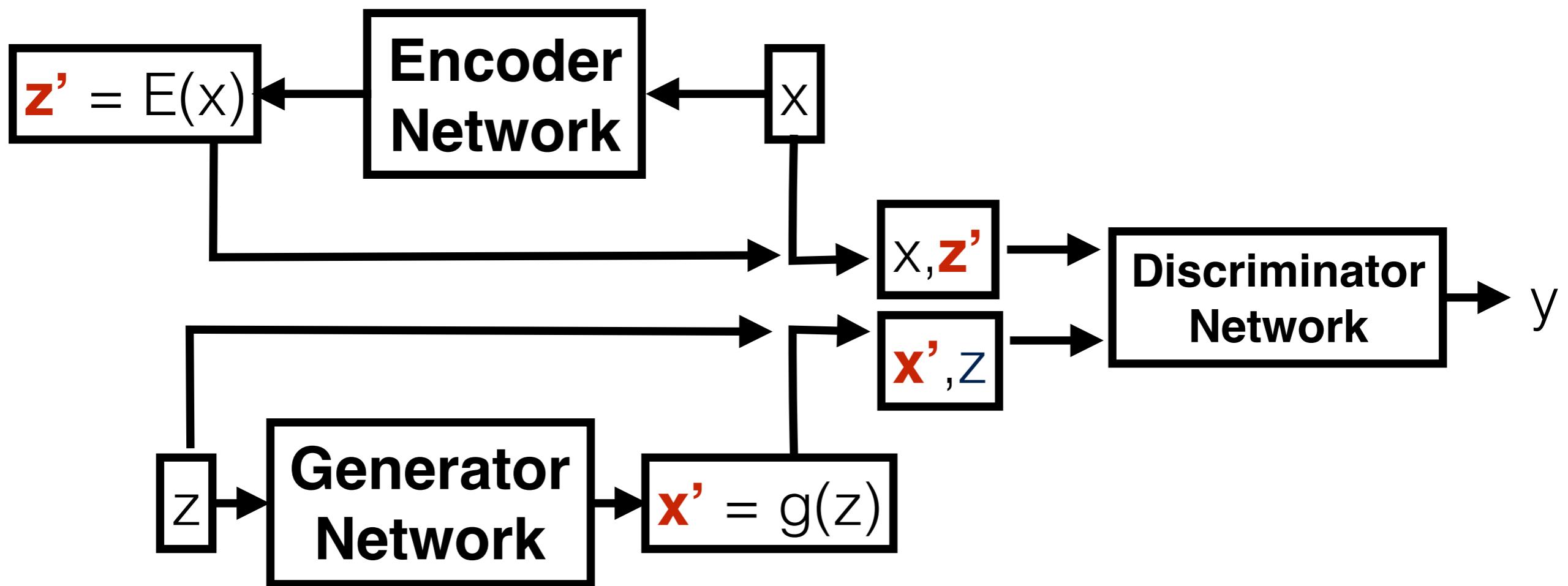


how about give  $x$ , generate  $z'$









## Cost function

$$C(w, \theta, \lambda) = - \left( \frac{1}{N^+} \sum_i \log\{D_w(x_i, E_\lambda(x_i))\} + \frac{1}{N^-} \sum_j \log\{1 - D_w(g_\theta(z_j), z_j)\} \right)$$
$$\max_{\theta, \lambda} \min_w C(w, \theta, \lambda)$$

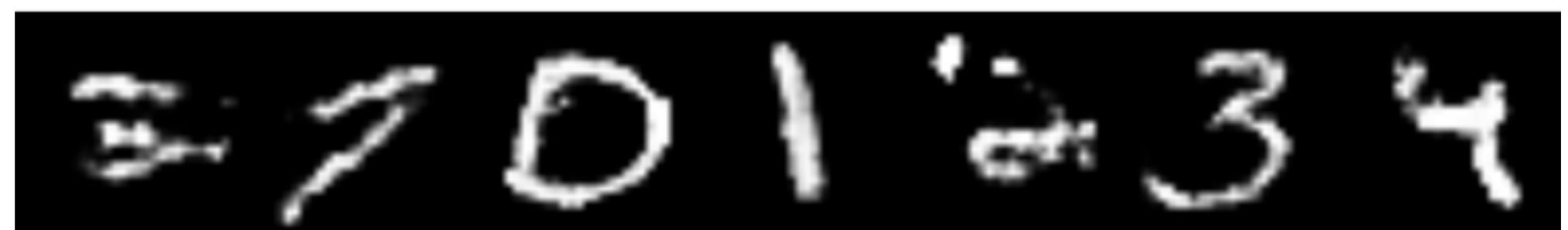
Adjust generator and encoder such that all fake data looks like real ones

$$g_\theta(z_j) = x' \rightarrow x \quad E_\lambda(x_i) = z' \rightarrow z$$

By this procedure,  $g(E_\lambda(x_i)) \rightarrow x_i$

$$E_\lambda \approx g_\theta^{-1}$$

$x$	 0 1 2 3 4 5 6 7 8 9 0 1 2 3 4 5 6 7 8 9
$g(E(x))$	 0 1 2 3 7 5 1 7 3 7 0 1 6 3 4 4 6 7 8 7

$x$	 0 9 0 1 2 5 4
$g(E(x))$	 0 7 0 1 6 3 4

X



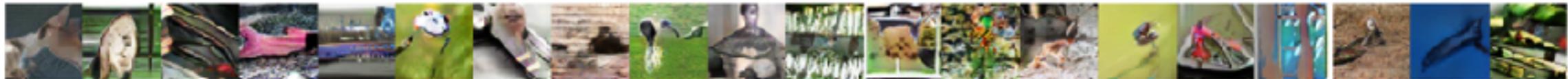
$g(E(x))$



X



$g(E(x))$



X



$g(E(x))$



X



$g(E(x))$



Assignment #3  
we will announce this weekend

