

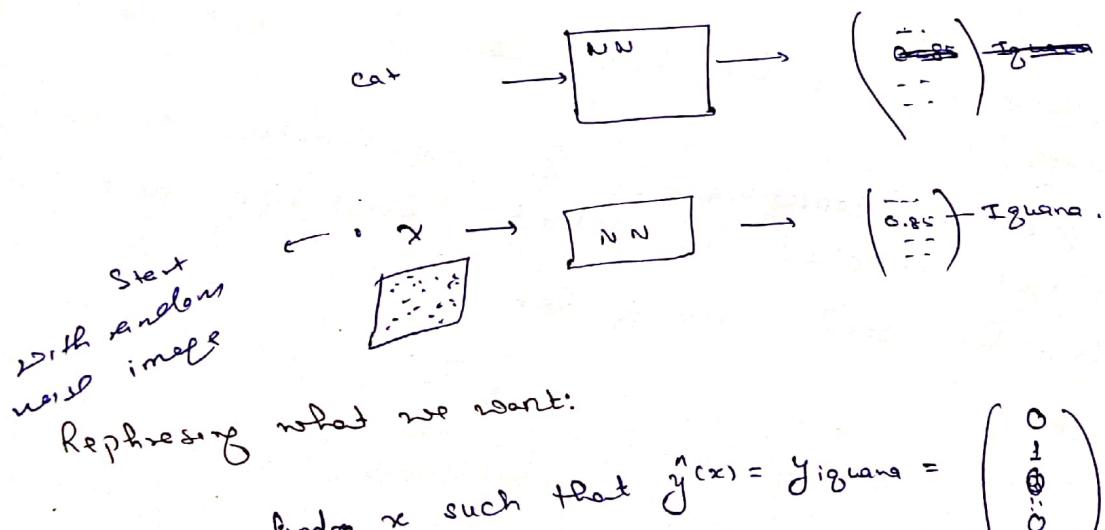
- Attacking NN with Adversarial Examples
- Generative Adversarial Examples.

- Attacking a network with adversarial
- Defenses against adversarial examples
- Why are neural networks vulnerable to Adversarial examples.

Paper - Intriguing properties of Neural networks 2013 See adj.  
Explaining & Harnessing Adversarial Examples: 2015 Ian.

### (1) Attacking a network with adversarial:

Goal: Given a network pretrained on Imagenet, find an input image that is not iguana but will be classified as an iguana.

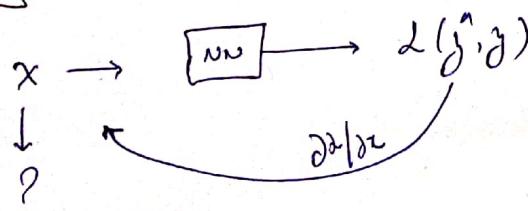


Loss function:

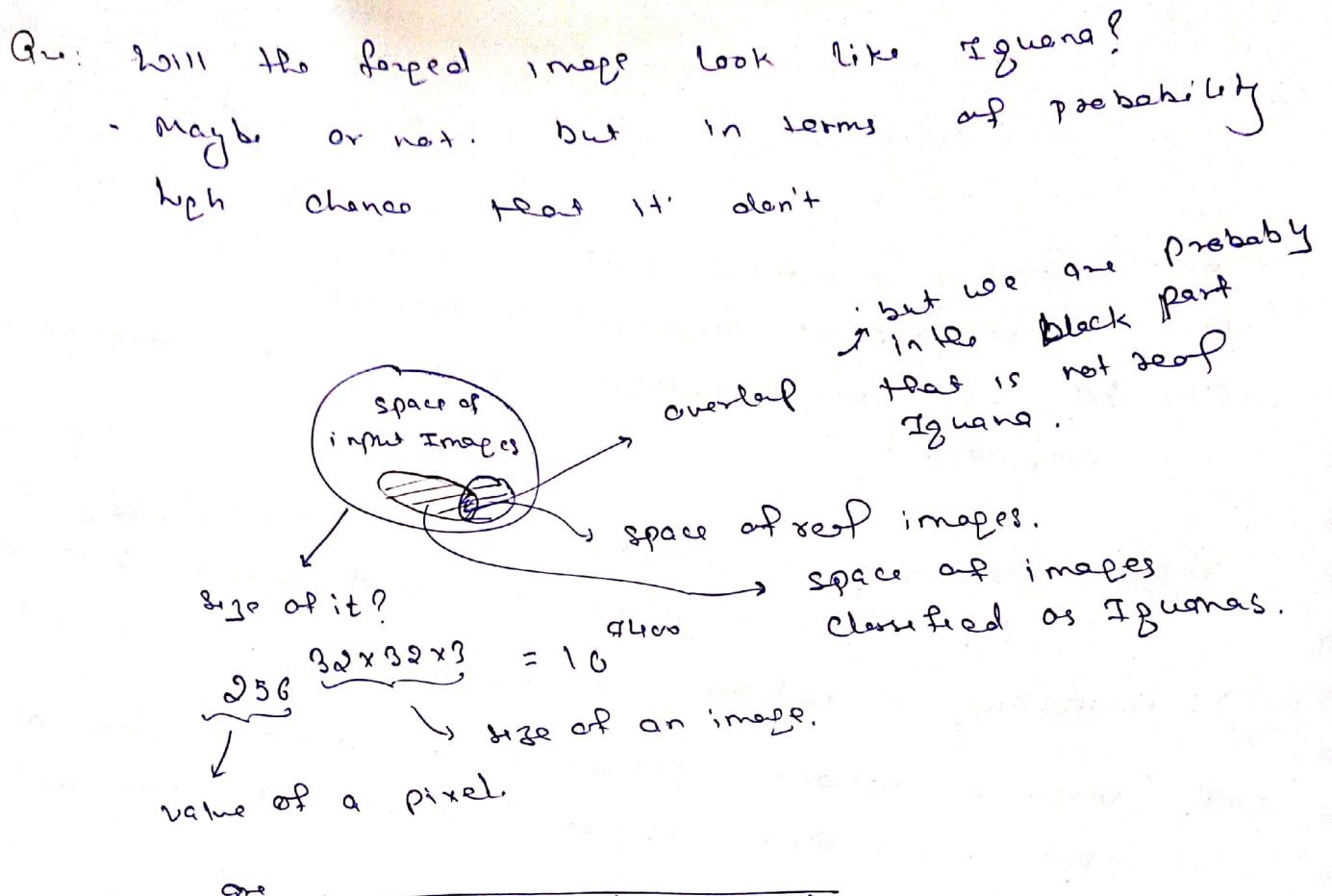
$$L(\hat{y}, y) = \frac{1}{2} \| \hat{y}(w, b, x) - y_{\text{iguana}} \|^2$$

$\hookrightarrow$  loss.

Optimize  $L(x)$ :

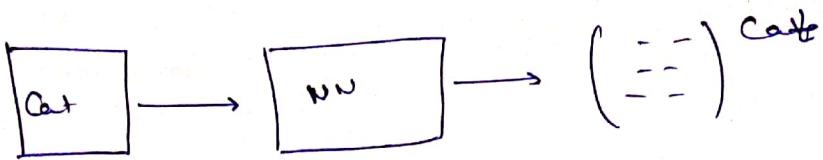


$$z = x - \alpha \frac{\partial L}{\partial x}$$



- ~~the car to be in green block part~~
- This is dangerous when human sees cat best network see Iguana.
- Harder problems / consequences - Face verification, violent Content not detected as violent, Step sign net Step sign for self-driving car.

Goal: Given a network pretrained on Imagenet, find an input image that is cat but will be classified as an iguana.



Rephrasing:

Find  $x$  such that  $\hat{y}(x) = y_{\text{iguana}} = \begin{pmatrix} 0 \\ 1 \\ 0 \\ 0 \\ 0 \end{pmatrix}$

And  $x \approx x_{\text{cat}}$ .

Loss:

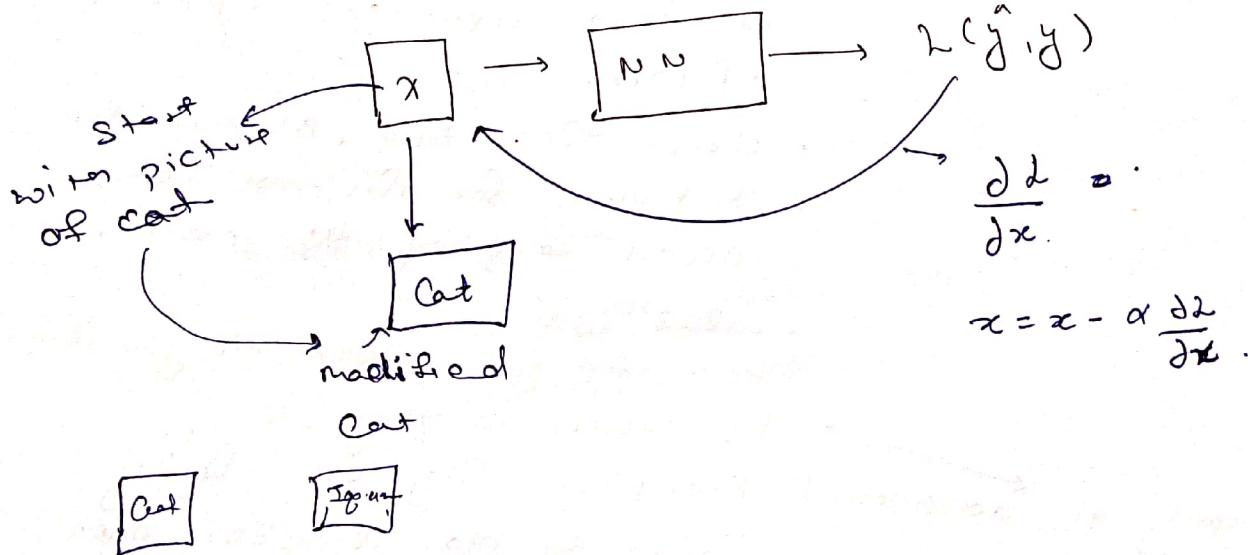
$$L(\hat{y}, y) = \frac{1}{2} \|\hat{y} - y\|_2^2 + \lambda \|x - x_{\text{cat}}\|_2^2$$

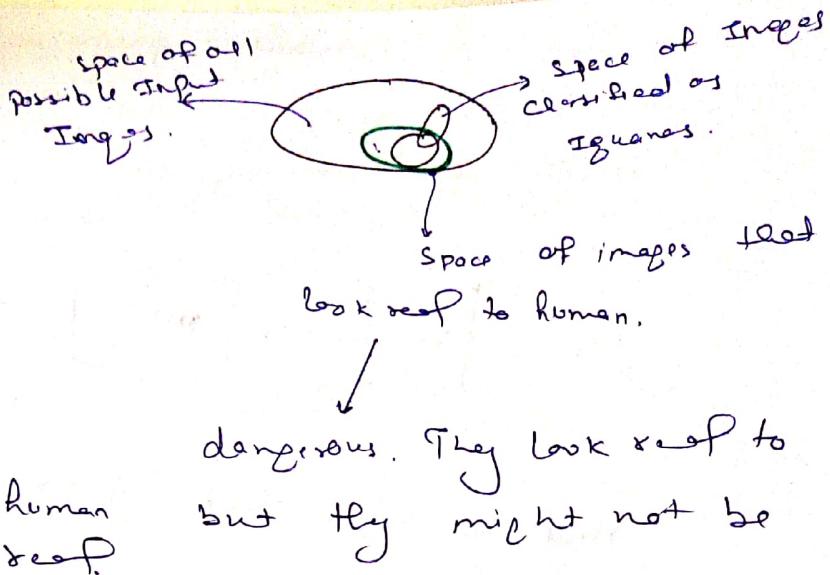
$\downarrow$

$$\hat{y} = \hat{y}(w, b, x)$$

$\uparrow$  variable

Optimize the Image:





## ③ Defenses against Adversarial Examples:

### Type of attacks

- non-targeted attacks.
- Targeted Attacks.

~~knowledge~~ of attacker.

white box <sup>access to</sup>

Black box.

<sup>even model is encrypted.</sup>

How would you attack when you have a black box?

- we cannot backpropagate but we can have numerical approximation of loss.
- change the input, observe the output. For this, we have access to query the model.
- related to it
- How will you attack if you don't have access to query?

On property of adversarial examples. They are highly transferable. What I am going to do is, say another model is animal classifier. I will build my own animal classifier, build a adversarial image on it. It is highly likely that it is an adversarial example for another model also.

Open Research Topic

Why this happen?  
How to defend?

Paper:

Solution 1:

- Create a SafetyNet. This will allow only those that are safe & not adversarial.
- Maybe we try to fool this also & set the constraint there is that we have to fool both at the same time.

Paper: SafetyNet - Detecting & Rejecting Adversarial Examples Robustly.

Solution 2:

- Generate adversarial examples & train on it.  
- very high complexity.

Solution 3:

Adversarial training.  $L_{\text{new}} = L(w, b, x, y) + \gamma L(w, b, \underline{x}_{\text{adv}}, \underline{y})$

↓  
train on adversarial examples at the same time.

↳ Complexity - ?

- Forward propagate  $x$  through the network to compute the first term.
- Generate  $\underline{x}$  adversarial with the optimization process & forward propagate it to calculate the second term, then back propagate over the weights of network.

super costly.

Another:

Adversarial legit pairing:

$$L_{\text{new}} = L(w, b, x, y) + \gamma \| f(x; w, b) - f(\underline{x}_{\text{adv}}; w, b) \|^2$$

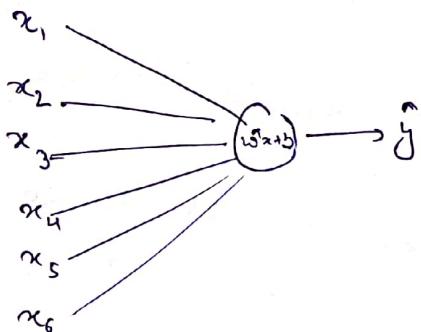
Paper: Adversarial legit pairing.

(C) Why are neural networks vulnerable to adversarial examples?

Why do adversarial examples exist?

Tan & his team argued that there are linear parts that are the cause of existence of adversarial examples. But, linear parts are what makes the training fast. like ReLU, also general, we focus on area that is linear by linear regression:

forward prop:  $\hat{y} = \mathbf{w}^T \mathbf{x} + b$ .



Network has been trained & converged to  $\approx$

$$\mathbf{w} = [1, 3, -1, 2, 2, 3]$$

$$b = 0$$

$$\hat{y} \quad \mathbf{x} = \begin{bmatrix} 1 \\ -1 \\ 2 \\ 0 \\ 3 \\ -2 \end{bmatrix} \Rightarrow \hat{y} = \mathbf{w}^T \mathbf{x} + b \\ = 1 \cdot 1 + 3 \cdot (-1) + (-1) \cdot 2 + 2 \cdot 0 + 2 \cdot 3 + 3 \cdot (-2) \\ = -4$$

Ques: How to change  $\mathbf{x}$  into  $\mathbf{x}'$  such that  $\hat{y} \rightarrow$  changes radically but  $\mathbf{x}' \approx \mathbf{x}$ .?

$\frac{\partial \hat{y}}{\partial \mathbf{x}}$  ~ impact on  $\hat{y}$  on small changes of  $\mathbf{x}$ .

$$\frac{\partial \hat{y}}{\partial \mathbf{x}} = \mathbf{w}^T$$

shape same as  $\mathbf{x}$ , size  $\mathbf{w}^T$ , not  $\mathbf{w}$ .

$$x^* = x + \epsilon w^\top$$

value of perturbation.

$$\hat{y}^* = w x^* + b = w x + \epsilon \underbrace{w w^\top}_{\|w\|^2}$$

push  $\hat{y}$  to larger value due to  $\hookrightarrow$  always positive.

+  $\epsilon \|w\|^2$ , if  $-\epsilon \|w\|^2$  then push  $\hat{y}$  to lower value.

$$x^* = x + \epsilon w^\top$$

$\downarrow 0.2$

$$\begin{bmatrix} 1 & +0.2 \\ -1 & +0.6 \\ 2 & -0.2 \\ 0 & +0.4 \\ 3 & +0.4 \\ -2 & +0.6 \end{bmatrix} = \begin{bmatrix} 1.2 \\ -0.4 \\ 1.8 \\ 0.4 \\ 3.4 \\ -1.4 \end{bmatrix}$$

$$\hat{y}^* = w x^* = 1.2 + (-1.2) - 1.8 + 0.8 + 6.8 - 4.2 \\ = \cancel{-0.6} \quad 1.6$$

insights: If  $x^*$  is not similar to  $x$ .

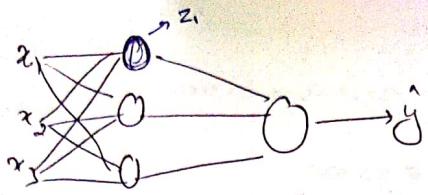
→ If  $w$  is large, then  $x^*$  is not similar to  $x$ .  
 $\hookrightarrow$  take  $\text{sign}(w)$

→ As  $x$  grows in dimension the impact of  $\epsilon \text{sign}(w)$  increases. (Because  $\sum w_i x_i$  grows)  
 $\hookrightarrow$  So we get it! For higher dimension is very high. Thus Adversarial examples are more prone to happen.

Fast gradient sign Method

$$x^* = x + \epsilon \text{sign}(\nabla f(w, x, y))$$

- formula to generate adversarial example



$$\frac{\partial \mathcal{L}}{\partial x} = \underbrace{\sum \frac{\partial \mathcal{L}}{\partial z_i^{[1]}}}_{\text{we want this high. Because to train}} \cdot \frac{\partial z_i^{[1]}}{\partial x}$$

↓  
for neuron corresponding to it, we need it

The networks similarly high.

$$\frac{\partial \mathcal{L}}{\partial w_i^{[1]}} = \frac{\partial \mathcal{L}}{\partial z_i^{[1]}} \cdot \frac{\partial z_i^{[1]}}{\partial w_i^{[1]}}$$

high.

Similarly to train  $w_i$ , we would need  $\frac{\partial \mathcal{L}}{\partial z_i^{[1]}}$  to be high.

The networks that have - high gradients & thus operating in the linear regime are more vulnerable to adversarial attacks. Because of above observation.

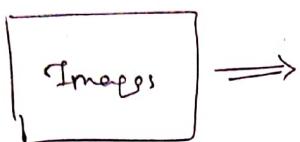
### GANS

Do neural networks actually understand the data?

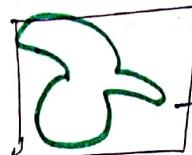
- A. Motivation
- B. G/D Game
- C. Training GAN
- D. Nice Results
- E. Evaluating GAN.

- (A) Motivation: Endow computers with an understanding of world.
- # Goal: Collect a lot of data, use it to train a model to generate similar data from scratch.
- Intuition: no. of parameters of model << amount of data

Probability distribution:

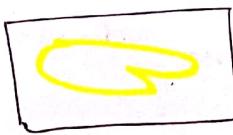
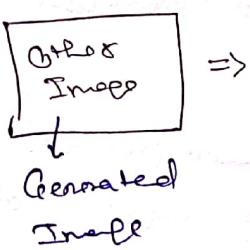
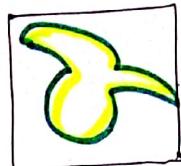


real data distribution



space of real world Images.  
⇒ Goal

Matching distributions.

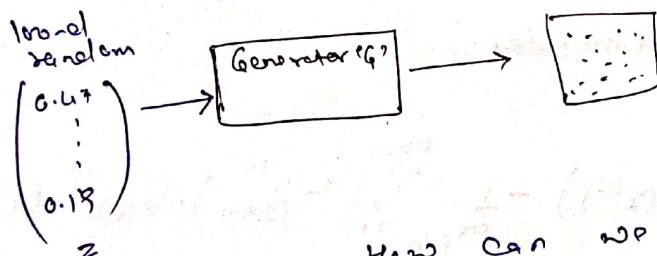


generated probability distribution.

we have so much amount of data. No. of parameters of the model is smaller than the amount of data.

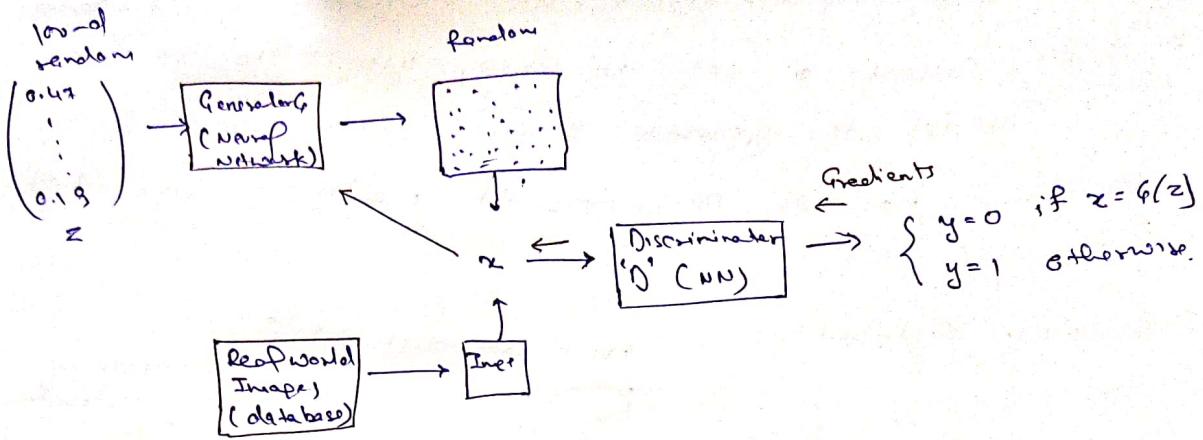
This is the intuition why generative networks exist. Because there is too much data, & there are not enough parameters to mimic this data.

Because network has enough features to understand the salient features, but doesn't have enough to overfit everything.



# Real world images

How can we train G to generate images from true data distributions?



$$\text{Gradient} \leftarrow \begin{cases} y=0 & \text{if } x=G(z) \\ y=1 & \text{otherwise.} \end{cases}$$

- + we take any image from real world database as one.
- + we have any image from real world database as one.
- If  $x$  is  $G(z)$  - means we can back propagate all the way back to the input of the discriminator.
- If  $x$  is from real database - then there is no relation between  $z$  &  $x$  & hence  $\partial L_D / \partial z = 0$ .

- + Run algorithm such as adam for both simultaneously
- + There are many methods that are tried to train GAN properly.

Loss  $L^D$ : Cost of Discriminator.

$$L^D = -\frac{1}{m} \sum_{i=1}^{m_{\text{real}}} y_{\text{real}}^{(i)} \cdot \log(D(x^{(i)})) + \frac{1}{m} \sum_{i=1}^{m_{\text{gen}}} (1-y_{\text{gen}}^{(i)}) \cdot \log(1-D(G(z^{(i)})))$$

Cross entropy 1:  
D should correctly label  
real data as 1

Cross entropy 2:  
D should correctly label  
generated data as 0.

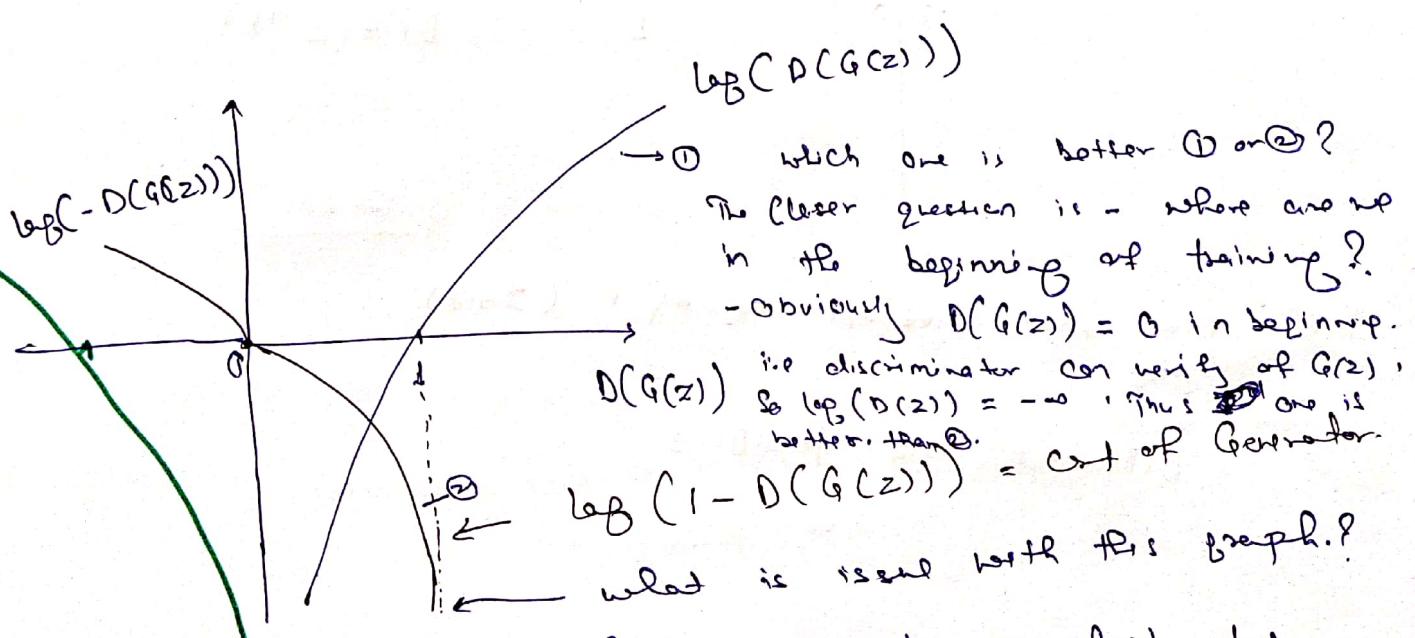
\*  $y_{gen}^{(i)} = 1$  always.  $1 - y_{gen}^{(i)} = 1$  also. So we can remove them.

$$JG = -JD = \frac{1}{m_{gen}} \sum_{i=1}^{m_{gen}} \log (1 - D(G(z^{(i)})))$$

Goal is to fool the Discriminator. If  $G$  managed to fool  $D$ , &  $D$  is very good then  $G$  is very good.

It's a game. If  $D$  is bad, we cannot say  $G$  is good. We want  $D$  to be up to be very good &  $G$  to go up at the same time. Until the equilibrium is reached at a certain point where  $D$  is always output one-half, like random probabilities. Because it cannot ~~sample~~ distinguish the samples coming from  $G$  versus ~~real~~ samples.

For generated images we want  $D$  to classify as one.



Thus:

But

$$J(G) = -J(D) = -\frac{1}{m} \sum_{i=1}^{m_{\text{gen}}} \log(1 - D(G(z^{(i)}))) \quad \textcircled{1}$$

Convert it to  $\textcircled{1}$  form of  $\log(D(G))$ .

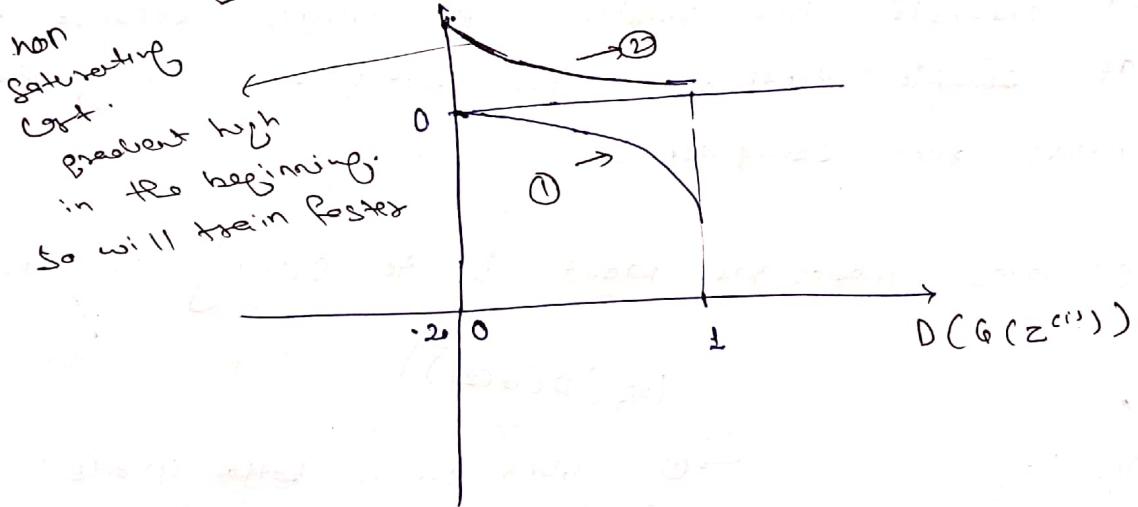
$$\min \log(1-x)$$

$$\equiv \max \log(x)$$

$$\equiv \min -\log(x)$$

Now  $J_G = \frac{1}{m} \sum_{i=1}^{m_{\text{gen}}} -\log(D(G(z^{(i)})))$

$$= \underbrace{-\frac{1}{m} \sum_{i=1}^{m_{\text{gen}}} \log(D(G(z^{(i)})))}_{\min} \quad \textcircled{2}$$



Paper: Are GANs created equal? (2018)

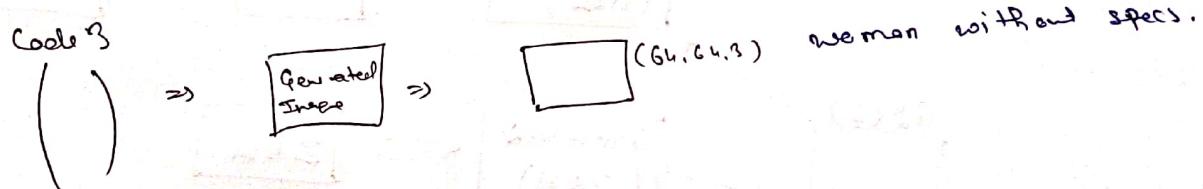
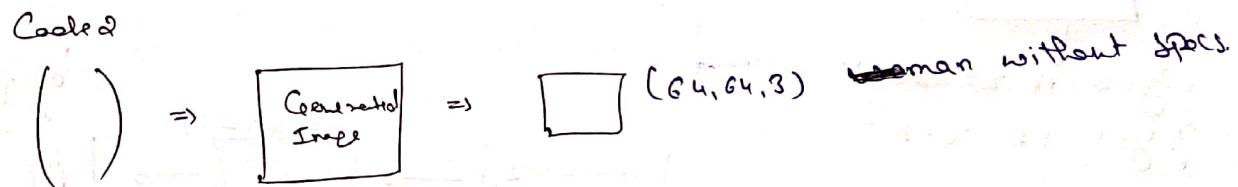
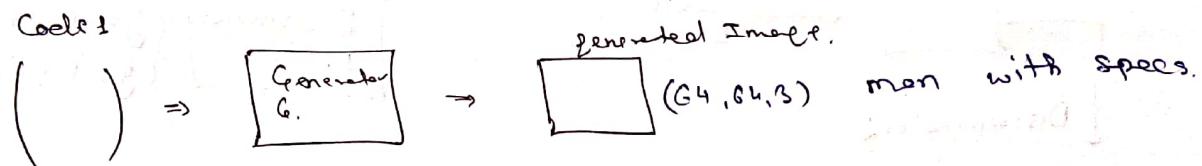
↓  
Different loss functions.

## ③ Training GAN:

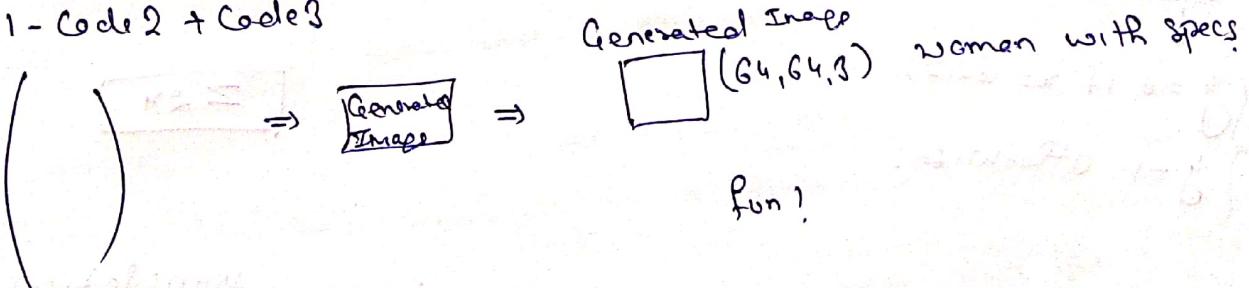
- Modification to the cost function
- keep  $\theta$  up to date with respect to  $G$ . ( $\theta$  update for  $D$  /  $\theta$  update for  $G$ )
- Use virtual BatchNorm.
- (not presented till now) One sided label smoothing.

## ④ Nice Results

1. Operation on codes:



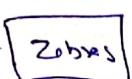
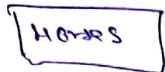
Code 1 - Code 2 + Code 3



Paper: Unsupervised Representation learning with Deep generating adversarial networks.

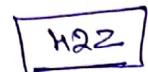
Goal: Convert horses to zebras on images, vice-versa.

Data:

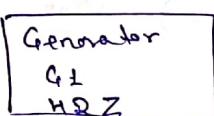


- unpaired images.

Architecture: Cycle GAN.



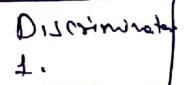
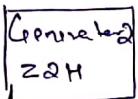
Horse



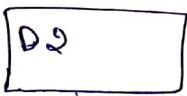
$G_1(H)$



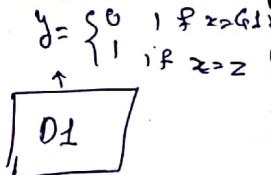
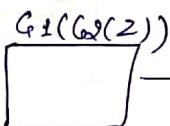
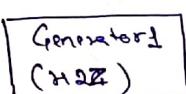
$G_2(G_1(H))$



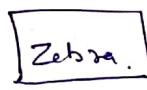
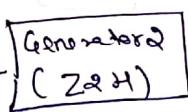
$$\begin{cases} y = 0 \text{ if } x = G_2(H) \\ y = 1 \text{ if } x = h \end{cases}$$



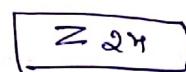
$$\begin{cases} y = 0 \text{ if } x = G_2(Z) \\ y = 1 \text{ if } x = H \end{cases}$$



$D_1$



$$\begin{cases} y = 0 \text{ if } x = G_2(Z) \\ y = 1 \text{ otherwise. } \end{cases}$$



Paper: Unpaired image to image translation  
Using cycle-consistent Adversarial Networks

Loss function:

$$J^{D_1} = -\frac{1}{m_{real}} \sum_{i=1}^{m_{real}} \log(D_1(h^{(i)})) - \frac{1}{m_{gen}} \sum_{i=1}^{m_{gen}} \log(1 - D_1(G_1(h^{(i)})))$$

$$J^{G_1} = -\frac{1}{m_{gen}} \sum_{i=1}^{m_{gen}} \log(D_1(G_1(h^{(i)})))$$

$$J^{D_2} = -\frac{1}{m_{real}} \sum_{i=1}^{m_{real}} \log(D_2(z^{(i)})) - \frac{1}{m_{gen}} \sum_{i=1}^{m_{gen}} \log(1 - D_2(G_2(z^{(i)})))$$

$$J^{G_2} = -\frac{1}{m_{gen}} \sum_{i=1}^{m_{gen}} \log(D_2(G_2(z^{(i)})))$$

$$J = J^{D_1} + J^{G_1} + J^{D_2} + J^{G_2} \rightarrow J^{\text{Cycle}}$$

CycleGAN:

Papers:-

- face2Ramen.
- FaceRamen using CycleGAN.
- Unpaired Image-to-Image Translation using Cycle Consistent Adversarial Networks.
- Pix2pix - generate image from drawing.  
train specifically for a domain.
- Photo-realistic single-Image super resolution using GAN.  
↳ Higher resolution image from lower resolution.
- Privacy preserving generative model support clinical data sharing.
- Learning beyond Human expertise with generative models for ~~learning~~ Dental restorations.
- Humanized Cifar cracking using discrete GAN.

How to evaluate GAN?

- Human annotators.

- Inception Score (IS)

$$IS(G) = \exp\left(-E_{x \sim P_D} [KL(P(y|x) || P(y))]\right)$$

↓ ↓  
measure of image quality → measure of image diversity.

Papers: Are GANs Created Equal?

: Improved Techniques for training GANs.

Also Frechet Inception Distance (FID)