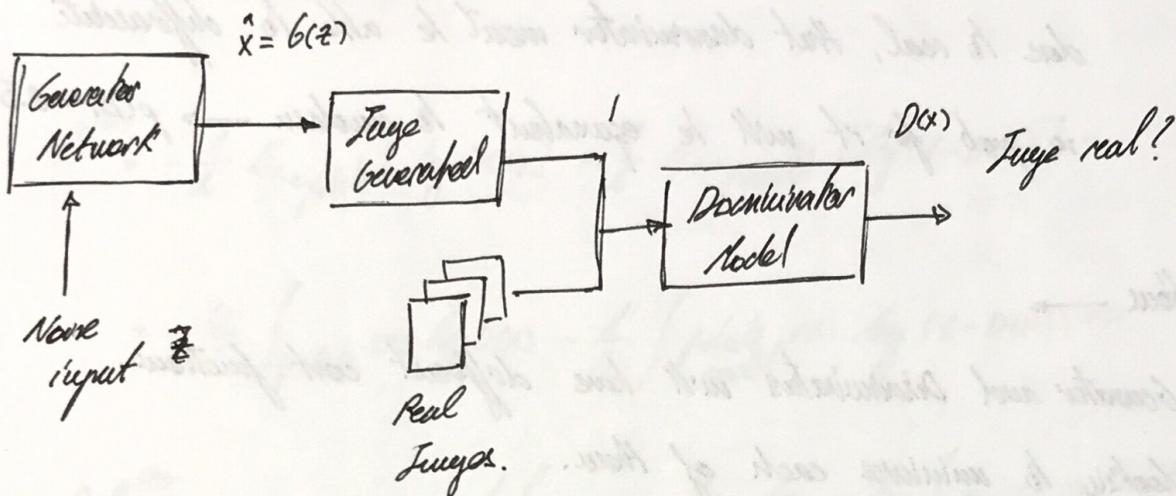


## Generative Adversarial Networks

\* Model →



\* Generator is a differentiable network so it can be trained

through gradient descent.

\* Gain intuition →

1. Generator →

- \* Generative model that maps input noise  $\mathbb{Z} = G(\mathbb{Z})$ .
- \* Its goal is to produce samples  $\hat{x}$ , from the distribution of training data  $p(x)$ .

\* Trained to maximize that the generated model probability will be close to the training data.

2. Discriminator →

- \* Differentiate between generated and real samples.

$$* D(x) \in (0, 1)$$

3. Generator and discriminator will be trained to fight against each other.

4. bad  $\rightarrow$  Ideal case the generator will create images so close to real, that discriminator won't be able to differentiate

so prob for it will be equivalent to random  $\rightarrow p(D) = 0.5$

\* last function  $\rightarrow$

- Generator and Discriminators will have different cost functions, looking to minimize each of them.

- Based on Game Theory, Equilibrium  $\rightarrow$

\* There will be a point in which the G will create images so 0.5 of them are taken as real, and D will discriminate properly 0.5  $\rightarrow$  balanced cost for 2 of them.

- Discriminator loss function  $\rightarrow$

$$L^{(D)}(\theta_D, \theta_G) = -\frac{1}{2} E_{x \sim \text{real}} [\log D(x)] + -\frac{1}{2} E_{x \sim \text{model}} [\log (1 - D(\hat{x}))]$$

$$= -\frac{1}{2} E_{x \sim \text{real}} [\log D(x)] - \frac{1}{2} E_{x \sim \text{model}} [\log (1 - D(G(x)))]$$

\* Cross-entropy for binary classification.

\* Structure  $\rightarrow D(x) = 1$  real data;  $D(\hat{x}) = 0$  for generated.

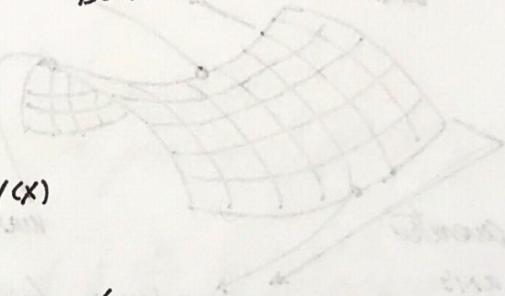
- Optimal minimization for Discontinuator  $\rightarrow$ 
  - \* Assuming continuous distribution, easier to make the profit.
  - \* Note:  $p_{data}, p_{model}$  are non-zero everywhere, assumption.

$$L^{(D)} = -\frac{1}{2} \mathbb{E}_{x \sim p_{data}} [\log D(x)] - \frac{1}{2} \mathbb{E}_{x \sim p_{model}} [\log (1-D(x))] =$$

$$= -\frac{1}{2} \int p_{data}(x) \log D(x) - \frac{1}{2} \int p_{model}(x) \log (1-D(x)) \Rightarrow$$

$$\frac{\partial L^{(D)}}{\partial D(x)} = 0 \rightarrow 0 = -\frac{1}{2} \int p_{data}(x) \frac{1}{D(x)} dx + \frac{1}{2} \int p_{model}(x) \frac{1}{1-D(x)} dx$$

$$D(x) = \frac{p_{data}(x)}{p_{data}(x) + p_{model}(x)}$$



$$p_{data}(x) + p_{model}(x)$$

- \* If D, G have enough capacity  $\rightarrow p_{model}(x) = p_{data}(x)$

Nash equilibrium  $\rightarrow D(x) = 0.5$

- Zero-sum game  $\rightarrow$

- \* Discontinuator  $\rightarrow$

- Wants to maximize the objective such that  $D(x) = 1$  and

$$D(G(z)) \approx 0.$$

- \* Generator  $\rightarrow$

- Wants to minimize the objective such that  $D(G(z)) \approx 1$ .

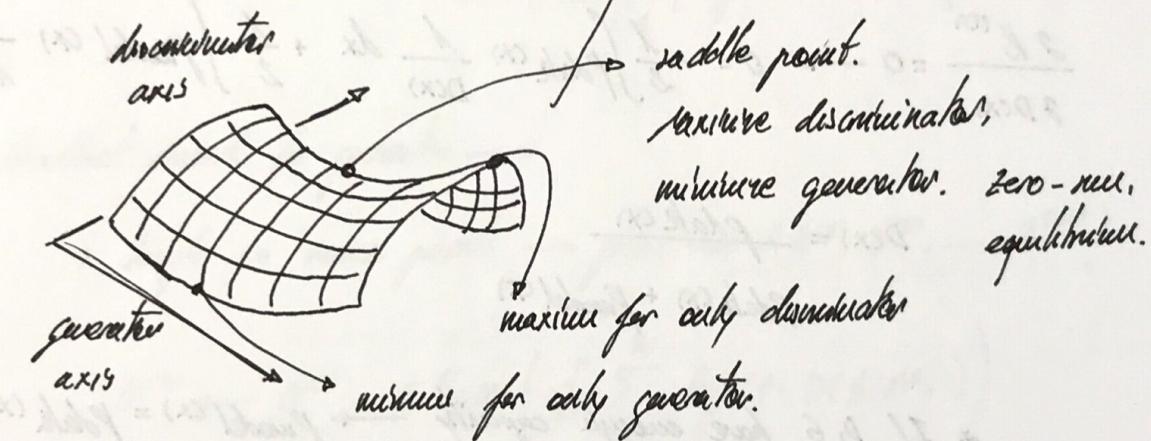
\* In this scenario, we can think of the sum of losses as 0 →

$$L^{(D)} + L^{(G)} = 0.$$

\* Find a way to minimize the maximum loss →

$$\min_{\theta^{(G)}} \max_{\theta^{(D)}} \left[ -\frac{1}{2} \mathbb{E}_{x \sim \text{update}(x)} [\log(D(x))] + \frac{1}{2} \mathbb{E}_z [\log(1 - D(G(z)))] \right]$$

Note: you change from loss.



- Optimization →

\* Gradient ascent for discriminator →

$$\theta^{(D)} \leftarrow \arg \max_{\theta^{(D)}} \left[ \frac{1}{2} \mathbb{E}_{x \sim \text{update}} [\log D(x)] + \frac{1}{2} \mathbb{E}_z [\log(1 - D(G(z)))] \right]$$

\* Gradient descent generator →

$$\theta^{(G)} \leftarrow \arg \min_{\theta^{(G)}} \underbrace{\frac{1}{2} \mathbb{E}_z [\log(1 - D(G(z)))]}_{\text{not } \mathbb{E}_{x \sim \text{update}} [\log(D(x))]}$$

This term depends on values for  $\theta^{(G)}$ ,  
not  $\mathbb{E}_{x \sim \text{update}} [\log(D(x))]$

- Tricky algorithm  $\rightarrow$

-  $k$  gradient steps for generator  $\rightarrow R = \text{hyperparameter}$ .

- $m =$   
mini-batch  
input data
- \* Sample  $m$  noise points  $\rightarrow p(z) \rightarrow \{z^{(1)}, z^{(2)}, z^{(3)}, \dots, z^{(m)}\}$
  - \* Sample  $m$  data points  $\rightarrow p_{\text{data}}(x) \rightarrow \{x^{(1)}, x^{(2)}, \dots, x^{(m)}\}$
  - \* Gradient step update  $\rightarrow$  ascent

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

- Gradient descent for generator  $\rightarrow$

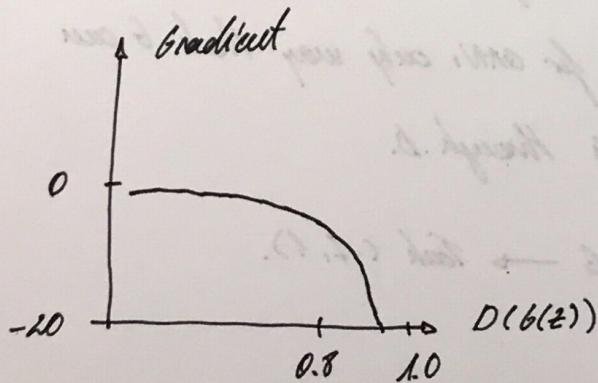
\* Sample  $m$  noise points  $\rightarrow p(z) \rightarrow \{z^{(1)}, z^{(2)}, \dots, z^{(m)}\}$

$$\theta^{(G)} \leftarrow \theta^{(G)} - \epsilon \nabla_{\theta^{(G)}} \left( \frac{1}{m} \sum_{i=1}^m \log(1 - D(G(z^{(i)}))) \right)$$

- Ian Goodfellow  $\rightarrow$  famous practice  $R=1$ .

- Gradient problem on  $\log(1 - D(G(z))) \rightarrow$

$$\frac{\partial}{\partial D(G(z))} (\log(1 - D(G(z)))) = -\frac{1}{1 - D(G(z))}$$

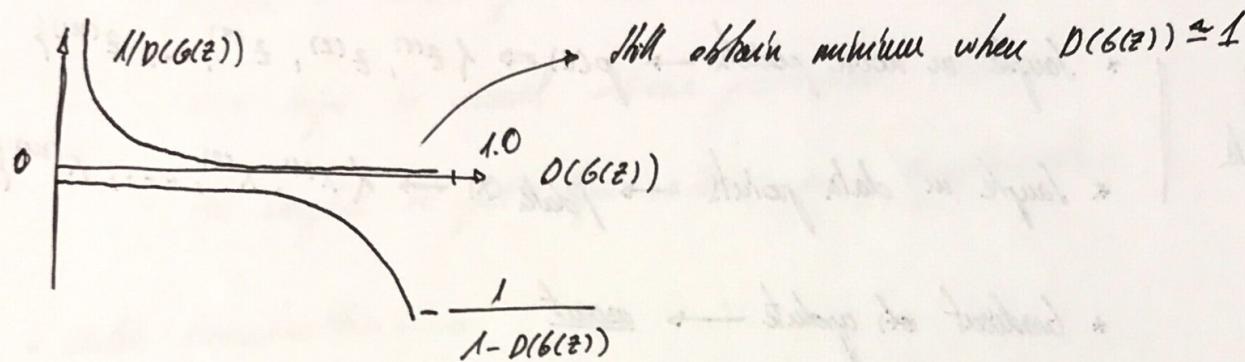


\*  $D(G(z)) = 0 \rightarrow \nabla \text{ almost } 0.$

\* At the beginning denominator differentiates well, the update is really small.

\* Instead alternative cost function for gradient  $\rightarrow$

$$\text{gradient} - \mathbb{E}_{\mathcal{D}} [\log(D(G(z)))] \rightarrow \text{grad. max} - 1/D(G(z))$$



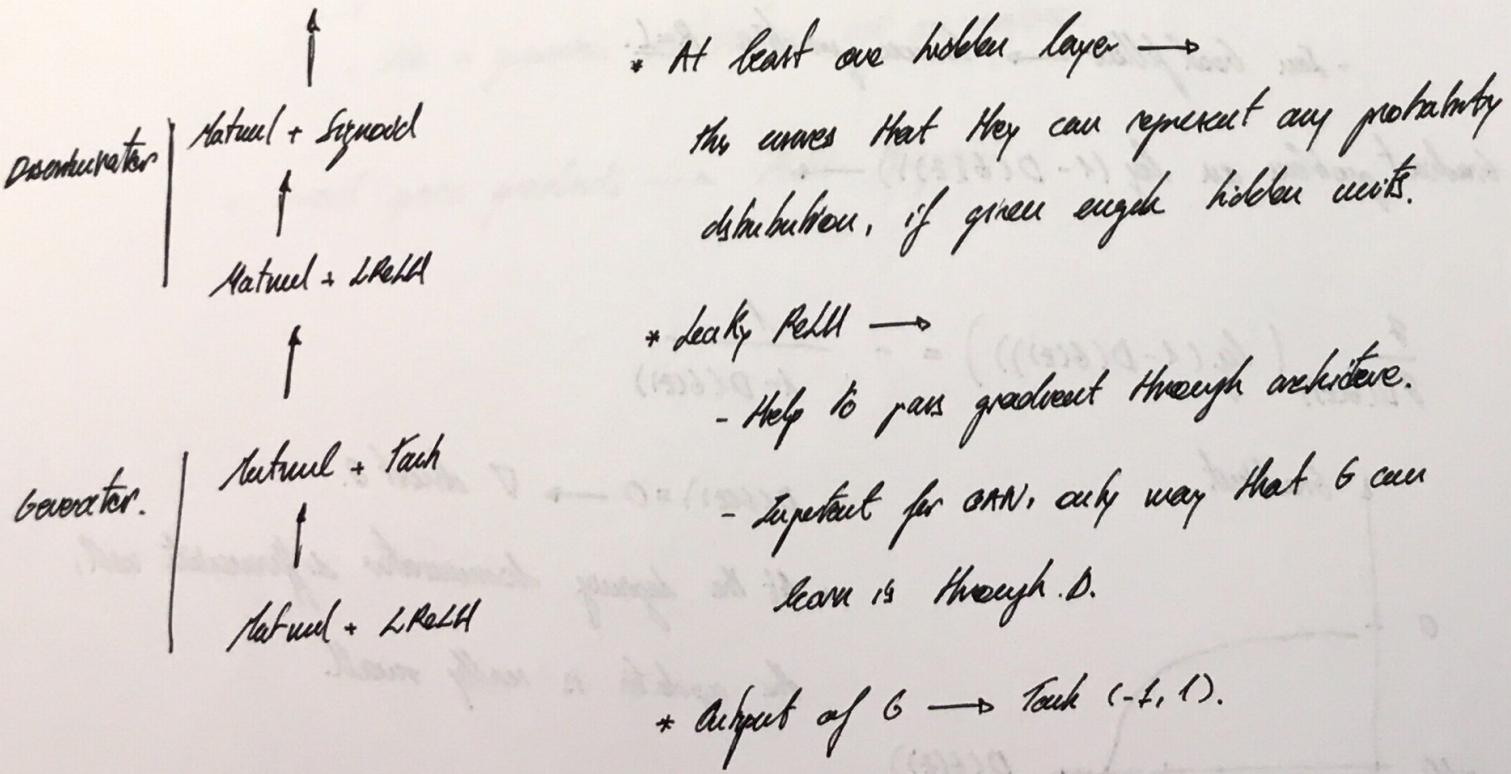
\* Algorithm changes on generator update  $\rightarrow$

$$\theta^{(G)} \leftarrow \theta^{(G)} + \epsilon \nabla_{\theta^{(G)}} \left[ \frac{1}{m} \sum_{i=1}^m \log D(G(z^{(i)})) \right]$$

\* No longer a zero-loss.

\* In practice, this is achieved by flipping labels.

- Back GAN  $\rightarrow$



- Training GANs →

\* One-sided label smoothing →

- provide D labels for real data as  $1-\kappa$ , fake 0.

-  $1-\kappa$  helps to avoid zero predictions. Prevent

the weights to get high values → regularization.

\* Batch Normalization →

- Apply every layer, except output generator and  
input of discriminator.

- Batch norm →

1. Read data for mini-batch

2. Generated data for mini-batch.

\* Noise input →

- Use a gaussian distribution, not an uniform.

\* Annoy space gradients → ReLU, Maxpool.