

TTIC 31230, Fundamentals of Deep Learning

David McAllester, April 2017

Generative Adversarial Networks (GANs)

The Generator and The Discriminator

A GAN consists of two networks: a generator $P_{\Theta}^{\text{gen}}(x)$ and a discriminator $P_{\Psi}^{\text{disc}}(y|x)$.

$$\Theta^* = \operatorname{argmax}_{\Theta} \min_{\Psi} \mathbb{E}_{(x,y) \sim (D \uplus P_{\Theta}^{\text{gen}})} \left[\log \frac{1}{P_{\Psi}^{\text{disc}}(y|x)} \right]$$

Here x is drawn from the data distribution D or the generator distribution P_{Θ}^{gen} with equal probability and $y = 1$ if x is drawn from D and -1 if x is drawn from P_{Θ}^{gen} .

The discriminator tries to determine which source x came from and the generator tries to fool the discriminator.

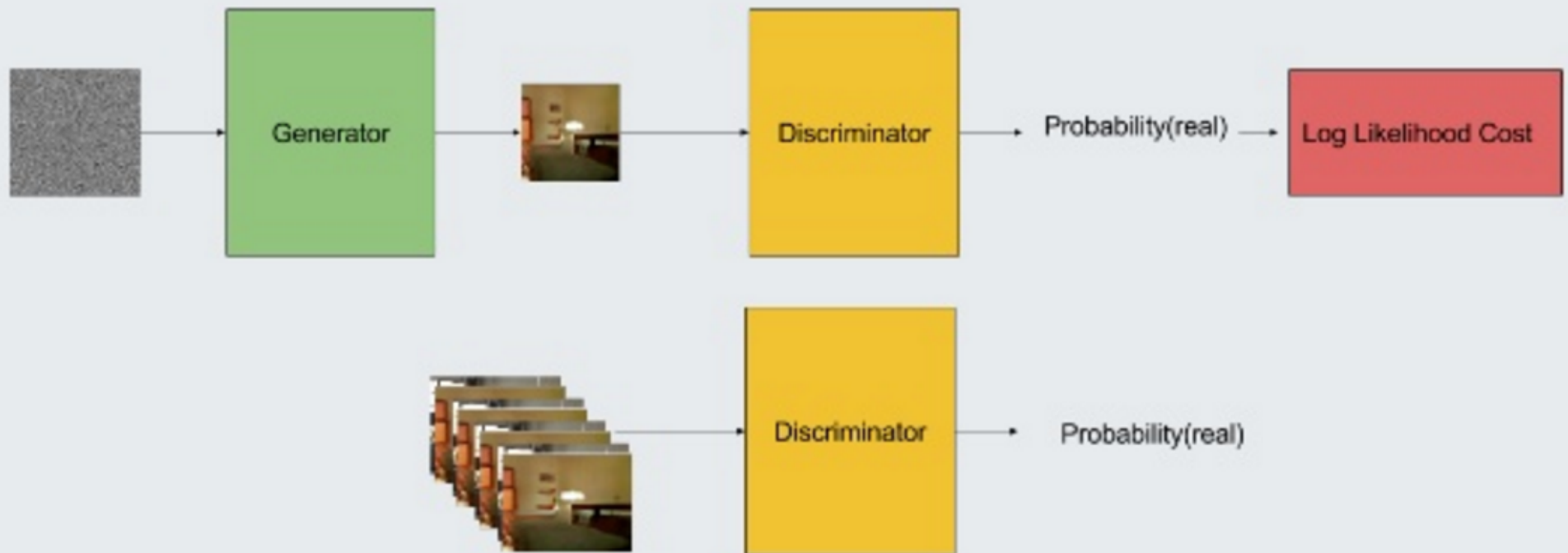
Consistency

If the discriminator is perfect, then the only way to fool it is to exactly copy the data distribution.

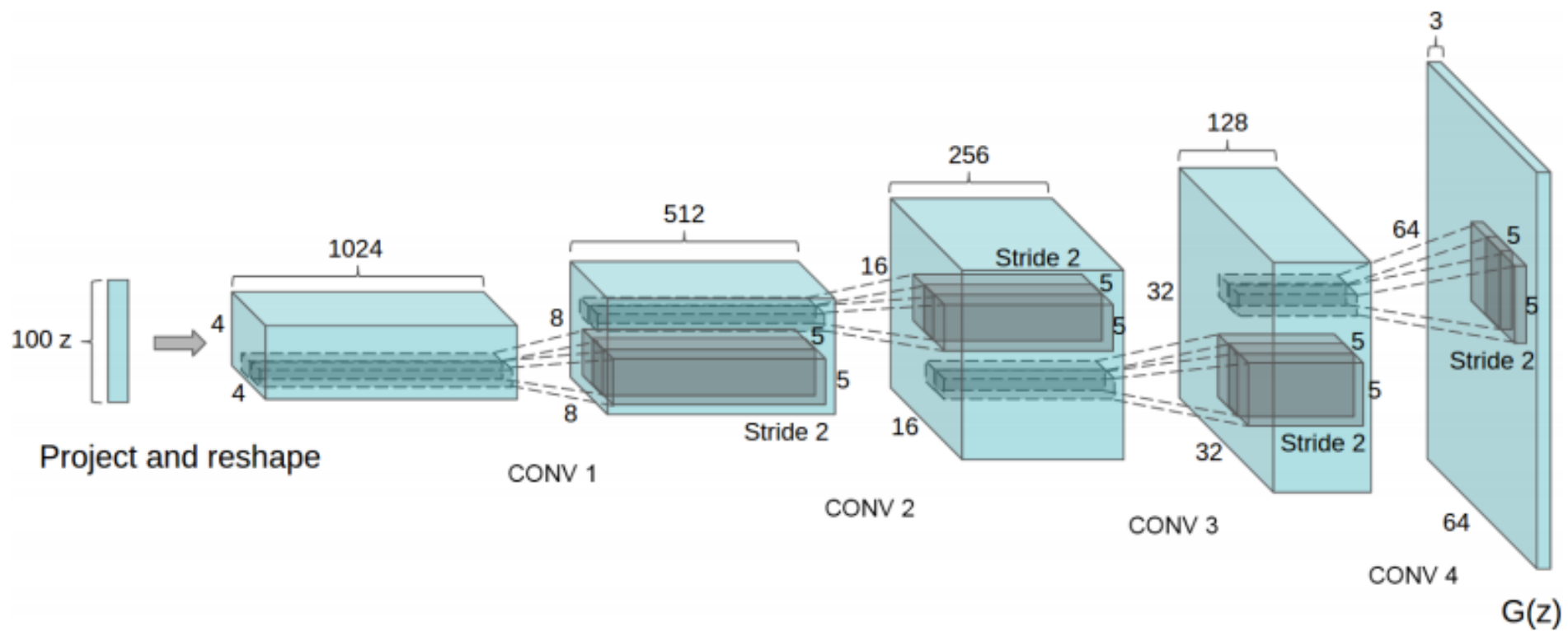
Consistency Theorem: If $P_{\Theta}^{\text{gen}}(x)$ and $P_{\Psi}^{\text{disc}}(y|x)$ are both universally expressive (any distribution can be represented) then $P_{\Theta^*}^{\text{gen}} = D$.

DC GANs, Radford, Metz and Chintala, ICLR 2016

Generative Adversarial Networks



The Generator



Generated Bedrooms



Interpolated Faces

[Ayan Chakrabarti]



Conditional Distribution Modeling

All distribution modeling methods apply to conditional distributions.

For conditional GANs we allow the generator to take x as an input and generate a conditional value c .

$$\Theta^* = \operatorname{argmax}_{\Theta} \min_{\Psi} \mathbb{E}_{x \sim D, (c,y) \sim (D(c|x) \uplus P_{\Theta}^{\text{gen}}(c|x))} \left[\log \frac{1}{P_{\Psi}^{\text{disc}}(y|c, x)} \right]$$

Here $y = 1$ if c is drawn from $D(c|x)$ and $y = -1$ if c is drawn from $P_{\Theta}^{\text{gen}}(c|x)$.

The Case of Imperfect Generation

$$\Theta^* = \operatorname{argmax}_{\Theta} \min_{\Psi} \mathbb{E}_{(x,y) \sim (D \uplus P_{\Theta}^{\text{gen}})} \left[\log \frac{1}{P_{\Psi}^{\text{disc}}(y|x)} \right]$$

$$\Psi^*(\Theta) = \operatorname{argmin}_{\Psi} \mathbb{E}_{(x,y) \sim (D \uplus P_{\Theta}^{\text{gen}})} \left[\log_2 \frac{1}{P(y|x)} \right]$$

$$P_{\Psi^*(\Theta)}^{\text{disc}}(y = 1|x) = \frac{P(x, y = 1)}{P(x)} = \frac{D(x)}{D(x) + P_{\Theta}^{\text{gen}}(x)}$$

$$\Theta^* = \operatorname{argmax}_{\Theta} \mathbb{E}_{(x,y) \sim (D \uplus P_{\Theta}^{\text{gen}})} \left[-\log_2 P_{\Psi^*(\Theta)}^{\text{disc}}(y|x) \right]$$

$$= \operatorname{argmax}_{\Theta} \frac{1}{2} \mathbb{E}_{(x,1) \sim D} \left[\log_2 \frac{D(x) + \pi(x|\Theta)}{D(x)} \right] \\ + \frac{1}{2} \mathbb{E}_{(x,-1) \sim P_{\Theta}^{\text{gen}}} \left[\log_2 \frac{D(x) + P_{\Theta}^{\text{gen}}(x)}{P_{\Theta}^{\text{gen}}(x)} \right]$$

$$= \operatorname{argmax}_{\Theta} 1 - \frac{1}{2} KL(D, A) - \frac{1}{2} KL(P_{\Theta}^{\text{gen}}, A)$$

$$A(x) = \frac{1}{2}(D(x) + P_{\Theta}^{\text{gen}}(x))$$

Jensen-Shannon Divergence (JSD)

We have arrived at the Jensen-Shannon divergence.

$$\Theta^* = \underset{\Theta}{\operatorname{argmin}} \quad \text{JSD}(D, P_{\Theta}^{\text{gen}})$$

$$\text{JSD}(P, Q) = \frac{1}{2}KL\left(P, \frac{P+Q}{2}\right) + \frac{1}{2}KL\left(Q, \frac{P+Q}{2}\right)$$

$$0 \leq \text{JSD}(P, Q) = \text{JSD}(Q, P) \leq 1$$

The Discriminator Tends to Win

If the discriminator “wins” the discriminator log loss goes to zero (becomes exponentially small) and there is no gradient to guide the generator.

In this case the learning stops and the generator is blocked from minimizing $\text{JSD}(D, P_{\Theta}^{\text{gen}})$.

The Standard Fix

The standard fix is to replace the loss

$$\ell = -\log P_{\Psi}^{\text{disc}}(y|x)$$

with

$$\tilde{\ell} = -y \log P_{\Psi}^{\text{disc}}(1|x)$$

These two loss functions agree when $y = 1$ (the case where x is drawn from D) but are very different when x is drawn from the generator ($y = -1$) and $P_{\Psi}^{\text{disc}}(1|x)$ is exponentially close to zero.

A Margin Interpretation of the Standard Fix

The standard fix can be interpreted in terms of the “margin” of binary classification.

For $y \in \{-1, 1\}$ we typically have $s_\Psi(1|x) = -s_\Psi(-1|x)$ and softmax over 1 and -1 gives

$$P_\Psi(y|x) = \frac{1}{1 + e^{-m}}$$

where the margin $m = 2ys_\Psi(x)$.

The margin is large when the prediction is confidently correct.

A Margin Interpretation of the Standard Fix

In the standard fix we (essentially) take the loss to be the margin of the discriminator.

The generator wants to reduce the discriminator's margin.

The direction of the update is the same but the step is much larger under margin-loss for generated inputs and large discriminator margins.

END