

Lecture 8:







tejasgokhale.com

CMSC 475/675 Neural Networks

NETWORK

NEWORK

Some slides based on Ranjay Krishna, Phillip Isola, Lana Lazebnik,

Recap: AE, VAE





- Encodes and decodes the data
- Low-dimensional bottleneck

Variational Autoencoder



Generative Modeling vs. Representation Learning

Representation Learning: "Analysis"

- Mapping data to representations

Generative Modeling: "Synthesis"

- Mapping representations to data







Generative Adversarial Networks

Ian J. Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, Yoshua Bengio[‡] Département d'informatique et de recherche opérationnelle Université de Montréal Montréal, QC H3C 3J7

- **Problem:** We want to sample from a high-dimensional training distribution p(x)
 - But there is no direct way to do this ...
 - We don't know which z maps to which image (so we can't use autoencoders)
- We know how to sample from a random distribution (e.g. Gaussian)

• Can we map a random distribution **directly** to p(x)?



NIPS 2014

Output: Sample from training distribution

Generator Network

Ζ

Input: Random noise





Generative Adversarial Networks

Goal: Map all z to some realistic-looking x

Image synthesis from "noise"



 ${f Sa}\ G:$





x = G(z)

Sampler

$$: \mathcal{Z} \to \mathcal{X}$$
$$z \sim p(z)$$
$$x = G(z)$$

Image synthesis from "noise"



 ${f Sa}\ G:$

x

G

x = G(z)

Sampler

$$\mathcal{Z} \to \mathcal{X}$$
$$\sim p(z)$$
$$= G(z)$$



© aleju/cat-generator





A two-player game:

- *G* tries to generate fake images that can fool *D*.
- *D* tries to detect fake images.

nages that can fool D. ges.



Learning objective (GANs) $\min_{G} \max_{D} \mathbb{E}_{z}[\log(1 - D(G(z))]]$







Learning objective (GANs)





real image

Learning objective (GANs) $\min_{\substack{G \\ G}} \max_{D} \mathbb{E}_{z}[\log(1 - D(G(z))] + \mathbb{E}_{x}[\log D(x)]$



GAN Training Breakdown

• From the discriminator D's perspective: binary classification: real vs. fake. • Nothing special: similar to 1 vs. 7 or cat vs. dog $\max \mathbb{E}[\log(1 - D(\mathbf{M})] + \mathbb{E}[\log D(\mathbf{M})]$

GAN Training Breakdown

- From the discriminator D's perspective: o binary classification: real vs. fake. • Nothing special: similar to 1 vs. 7 or cat vs. dog
- From the generator G's perspective:
 - Optimizing a loss that depends on a classifier D

 $\min \mathbb{E}_{z}[\mathcal{L}_{D}(G(z))]$ GGAN loss for G

 $\max \mathbb{E}[\log(1 - D(\mathbb{N})] + \mathbb{E}[\log D(\mathbb{N})]$

 $\min_{G} \mathbb{E}_{(x,y)} || F(G(x)) - F(y) ||$ Perceptual Loss for G

GAN Training Breakdown



G tries to synthesize fake images that fool D D tries to identify the fakes

- Global optimum when G reproduces data distribution.

• Training: iterate between training D and G with backprop.

Discriminator network: try to distinguish between real and fake images **Generator network:** try to fool the discriminator by generating real-looking images







Discriminator network: Generator network:

try to distinguish between real and fake images try to fool the discriminator by generating real-looking images





Connection to Game Theory: Zero-Sum "Minimax" Game

• Each player trying to minimize the opponent's profits

• Each player trying to maximize their own profits



Discriminator network: Generator network:

try to distinguish between real and fake images try to fool the discriminator by generating real-looking images





Real Images (from training set)

Fake and real images copyright Emily Denton et al. 2015. Reproduced with permission.







Fake and real images copyright Emily Denton et al. 2015. Reproduced with permission.





Discriminator network: try t **Generator network**: try t

try to distinguish between real and fake images try to fool the discriminator by generating real-looking images

Train jointly in **minimax game**

Minimax objective function:



$\min_{\theta_g} \max_{\theta_d} \left[\mathbb{E}_{x \sim p_{data}} \log D_{\theta_d}(x) + \mathbb{E}_{z \sim p(z)} \log(1 - D_{\theta_d}(G_{\theta_g}(z))) \right]$



Discriminator network: **Generator network**:

try to distinguish between real and fake images try to fool the discriminator by generating real-looking images

Train jointly in **minimax game**

Minimax objective function:

Discriminator outputs likelihood in (0,1) of real image

 $\min_{\theta_g} \max_{\theta_d} \left[\mathbb{E}_{x \sim p_{data}} \log D_{\theta_d}(x) + \right]$ Discriminator output for real data x

- Discriminator (θ_d) wants to **maximize objective** s.t. D(x) is close to 1 (real) and D(G(z)) is close to 0 (fake)

- Generator (θ_{a}) wants to **minimize objective** s.t.

$$-\mathbb{E}_{z\sim p(z)}\log(1-D_{\theta_d}(G_{\theta_g}(z)))$$

Discriminator output for generated fake data G(z)

D(G(z)) is close to 1 (discriminator is fooled into thinking generated G(z) is real)



Minimax objective function:

$$\min_{\theta_g} \max_{\theta_d} \left[\mathbb{E}_{x \sim p_{data}} \log D_{\theta_d}(x) \right]$$

Alternate between:

1. Gradient ascent on discriminator

$$\max_{\theta_d} \left[\mathbb{E}_{x \sim p_{data}} \log D_{\theta_d}(x) + \mathbb{E}_{z \sim p(z)} \log(1 - D_{\theta_d}(G_{\theta_g}(z))) \right]$$

2. Gradient descent on generator

$$\min_{\theta_g} \mathbb{E}_{z \sim p(z)} \log(1 - D_{\theta_d}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g$$

In practice, optimizing this generator objective does not work well!

 $+ \mathbb{E}_{z \sim p(z)} \log(1 - D_{\theta_d}(G_{\theta_g}(z))) \Big|$

When sample is likely (z)) fake, want to learn from $\frac{3}{2}$ it to improve generator $\frac{3}{2}$ (move to the right on X $\frac{3}{2}$



Minimax objective function:

$$\min_{\theta_g} \max_{\theta_d} \left[\mathbb{E}_{x \sim p_{data}} \log D_{\theta_d}(x) \right]$$

Alternate between:

Gradient ascent on discriminator

$$\max_{\theta_d} \left[\mathbb{E}_{x \sim p_{data}} \log D_{\theta_d}(x) + \mathbb{E}_{z \sim p(z)} \log(1 - D_{\theta_d}(G_{\theta_g}(z))) \right]$$

Gradient descent on generator 2.

$$\min_{\theta_g} \mathbb{E}_{z \sim p(z)} \log(1 - D_{\theta_d}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g}(G_{\theta_g$$

In practice, optimizing this generator objective does not work well!

 $+ \mathbb{E}_{z \sim p(z)} \log(1 - D_{\theta_d}(G_{\theta_g}(z))) \Big|$

Gradient signal dominated by region where sample is already good

When sample is likely (z)) fake, want to learn from it to improve generator (move to the right on X axis).



But gradient in this region is relatively flat!





Minimax objective function:

$$\min_{\theta_g} \max_{\theta_d} \left[\mathbb{E}_{x \sim p_{data}} \log D_{\theta_d}(x) \right]$$

Alternate between:

1. Gradient ascent on discriminator

$$\max_{\theta_d} \left[\mathbb{E}_{x \sim p_{data}} \log D_{\theta_d}(x) + \mathbb{E} \right]$$

2. Instead: Gradient ascent on generator, different objective

$$\max_{\theta_g} \mathbb{E}_{z \sim p(z)} \log(D_{\theta_d}(G_{\theta_g}(z)))$$

Instead of minimizing likelihood of discriminator being correct, now maximize likelihood of discriminator being wrong. Same objective of fooling discriminator, but now higher gradient signal for bad samples => works much better! Standard in practice.

 $+ \mathbb{E}_{z \sim p(z)} \log(1 - D_{\theta_d}(G_{\theta_g}(z))) \Big|$

 $\mathbb{E}_{z \sim p(z)} \log(1 - D_{\theta_d}(G_{\theta_g}(z))) \Big|$



Putting it together: GAN training algorithm

for number of training iterations do for k steps do

- Sample minibatch of m noise samples $\{z^{(1)}, \ldots, z^{(m)}\}$ from noise prior $p_q(z)$. • Sample minibatch of m examples $\{x^{(1)}, \ldots, x^{(m)}\}$ from data generating distribution
- $p_{\text{data}}(\boldsymbol{x}).$
- Update the discriminator by ascending its stochastic gradient:

$$\nabla_{\theta_d} \frac{1}{m} \sum_{i=1}^m \left[\log D_{\theta_d}(x^{(i)}) + \log(1 - D_{\theta_d}(G_{\theta_g}(z^{(i)}))) \right]$$

end for

$$\nabla_{\theta_g} \frac{1}{m} \sum_{i=1}^m \log(D_{\theta_d})$$

end for

• Sample minibatch of m noise samples $\{z^{(1)}, \ldots, z^{(m)}\}$ from noise prior $p_q(z)$. • Update the generator by ascending its stochastic gradient (improved objective):

 $(G_{\theta_a}(z^{(i)})))$

Putting it together: GAN training algorithm

for number of training iterations do for k steps do

- Sample minibatch of m noise samples $\{z^{(1)}, \ldots, z^{(m)}\}$ from noise prior $p_g(z)$. • Sample minibatch of m examples $\{x^{(1)}, \ldots, x^{(m)}\}$ from data generating distribution
- $p_{\text{data}}(\boldsymbol{x}).$
- Update the discriminator by ascending its stochastic gradient:

$$\nabla_{\theta_d} \frac{1}{m} \sum_{i=1}^m \left[\log D_{\theta_d}(x^{(i)}) + \log(1 - D_{\theta_d}(G_{\theta_g}(z^{(i)}))) \right]$$

Followup work (e.g. Wasserstein GAN, BEGAN) alleviates this problem, better stability!

Some find k=1

others use k > 1,

more stable,

no best rule.

end for

$$abla_{ heta_g} rac{1}{m} \sum_{i=1}^m \log(D_{ heta_d})$$

end for

Arjovsky et al. "Wasserstein gan." arXiv preprint arXiv:1701.07875 (2017) Berthelot, et al. "Began: Boundary equilibrium generative adversarial networks." arXiv preprint arXiv:1703.10717 (2017)

• Sample minibatch of m noise samples $\{z^{(1)}, \ldots, z^{(m)}\}$ from noise prior $p_g(z)$. • Update the generator by ascending its stochastic gradient (improved objective):

 $(G_{\theta_a}(z^{(i)})))$

Discriminator network: Generator network:

try to distinguish between real and fake images try to fool the discriminator by generating real-looking images

Discriminator Network



Real or Fake



Real Images (from training set)

After training, use generator network to generate new images

Fake and real images copyright Emily Denton et al. 2015. Reproduced with permission.





Generative Adversarial Nets: Convolutional Architectures

Discriminator is a convolutional network

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.

Radford et al, "Unsupervised Representation Learning with Deep Convolutional Generative Adversarial Networks", ICLR 2016

- Generator is an upsampling network with fractionally-strided convolutions



Generative Adversarial Nets: Convolutional Architectures

Samples from the model look much better!

Radford et al, ICLR 2016





Generative Adversarial Nets: Convolutional Architectures

Interpolating between random points in laten space

Radford et al, ICLR 2016





Generative Adversarial Nets: Interpretable Vector Math

Neutral man Neutral woman Smiling woman

Samples from the model





Radford et al, ICLR 2016







Generative Adversarial Nets: Interpretable Vector Math

Neutral man Neutral woman Smiling woman

Samples from the model



Average Z vectors, do arithmetic









Radford et al, ICLR 2016









Generative Adversarial Nets: Interpretable Vector Math

Neutral woman Smiling woman

Samples from the model



Average Z vectors, do arithmetic









Neutral man

Radford et al, ICLR 2016









Smiling Man





Generative Adversarial Nets: Interpretable Vector Math Glasses man No glasses man

)= A





No glasses woman

Radford et al, ICLR 2016











GAN in PyTorch



https://pytorch.org/tutorials/beginner/dcgan_faces_tutorial.html

```
class Generator(nn.Module):
class Discriminator(nn.Module):
                                                                             def __init__(self, ngpu):
    def __init__(self, ngpu):
                                                                                 super(Generator, self).__init__()
        super(Discriminator, self).__init__()
                                                                                 self.ngpu = ngpu
        self.ngpu = ngpu
                                                                                 self.main = nn.Sequential(
        self.main = nn.Sequential(
                                                                                     # input is Z, going into a convolution
            # input is ``(nc) x 64 x 64``
                                                                                     nn.ConvTranspose2d( nz, ngf * 8, 4, 1, 0, bias=False),
            nn.Conv2d(nc, ndf, 4, 2, 1, bias=False),
                                                                                     nn.BatchNorm2d(ngf * 8),
            nn.LeakyReLU(0.2, inplace=True),
                                                                                     nn.ReLU(True),
            # state size. ``(ndf) x 32 x 32``
                                                                                     # state size. ``(ngf*8) x 4 x 4``
                                                                                     nn.ConvTranspose2d(ngf * 8, ngf * 4, 4, 2, 1, bias=False),
            nn.Conv2d(ndf, ndf \star 2, 4, 2, 1, bias=False),
                                                                                     nn.BatchNorm2d(ngf * 4),
            nn.BatchNorm2d(ndf \star 2),
                                                                                     nn.ReLU(True),
            nn.LeakyReLU(0.2, inplace=True),
                                                                                     # state size. ``(ngf*4) x 8 x 8``
            # state size. ``(ndf*2) x 16 x 16``
                                                                                     nn.ConvTranspose2d( ngf * 4, ngf * 2, 4, 2, 1, bias=False),
            nn.Conv2d(ndf * 2, ndf * 4, 4, 2, 1, bias=False),
                                                                                     nn.BatchNorm2d(ngf * 2),
            nn.BatchNorm2d(ndf * 4),
                                                                                     nn.ReLU(True),
            nn.LeakyReLU(0.2, inplace=True),
                                                                                     # state size. ``(ngf*2) x 16 x 16``
            # state size. ``(ndf*4) x 8 x 8``
                                                                                     nn.ConvTranspose2d( ngf * 2, ngf, 4, 2, 1, bias=False),
            nn.Conv2d(ndf \star 4, ndf \star 8, 4, 2, 1, bias=False),
                                                                                     nn.BatchNorm2d(ngf),
            nn.BatchNorm2d(ndf * 8),
                                                                                     nn.ReLU(True),
            nn.LeakyReLU(0.2, inplace=True),
                                                                                     # state size. ``(ngf) x 32 x 32``
            # state size. ``(ndf*8) x 4 x 4``
                                                                                     nn.ConvTranspose2d( ngf, nc, 4, 2, 1, bias=False),
            nn.Conv2d(ndf * 8, 1, 4, 1, 0, bias=False),
                                                                                     nn.Tanh()
            nn.Sigmoid()
                                                                                     # state size. ``(nc) x 64 x 64``
                                                                             def forward(self, input):
    def forward(self, input):
                                                                                 return self.main(input)
        return self.main(input)
```

https://github.com/soumith/ganhacks for tips and tricks for training GANs



Since then: Explosion of GANs

"The GAN Zoo"

- GAN Generative Adversarial Networks
- 3D-GAN Learning a Probabilistic Latent Space of Object Shapes via 3D Generative-Adversarial M
- acGAN Face Aging With Conditional Generative Adversarial Networks
- AC-GAN Conditional Image Synthesis With Auxiliary Classifier GANs
- AdaGAN AdaGAN: Boosting Generative Models
- AEGAN Learning Inverse Mapping by Autoencoder based Generative Adversarial Nets
- AffGAN Amortised MAP Inference for Image Super-resolution
- AL-CGAN Learning to Generate Images of Outdoor Scenes from Attributes and Semantic Layouts
- ALI Adversarially Learned Inference
- AM-GAN Generative Adversarial Nets with Labeled Data by Activation Maximization
- AnoGAN Unsupervised Anomaly Detection with Generative Adversarial Networks to Guide Market
- ArtGAN ArtGAN: Artwork Synthesis with Conditional Categorial GANs
- b-GAN b-GAN: Unified Framework of Generative Adversarial Networks
- Bayesian GAN Deep and Hierarchical Implicit Models
- **BEGAN BEGAN: Boundary Equilibrium Generative Adversarial Networks** •
- **BiGAN Adversarial Feature Learning** .
- BS-GAN Boundary-Seeking Generative Adversarial Networks
- CGAN Conditional Generative Adversarial Nets
- CaloGAN CaloGAN: Simulating 3D High Energy Particle Showers in Multi-Layer Electromagnetic with Generative Adversarial Networks
- CCGAN Semi-Supervised Learning with Context-Conditional Generative Adversarial Networks
- CatGAN Unsupervised and Semi-supervised Learning with Categorical Generative Adversarial Network
- CoGAN Coupled Generative Adversarial Networks

| | Context-RNN-GAN - Contextual RNN-GANs for Abstract Reasoning Diagram Generation |
|--------------|--------------------------------------------------------------------------------------------------------------------------------------------|
| And allowed | C-RNN-GAN - C-RNN-GAN: Continuous recurrent neural networks with adversarial training |
| Adeling | CS-GAN - Improving Neural Machine Translation with Conditional Sequence Generative Adversarial Nets |
| | CVAE-GAN - CVAE-GAN: Fine-Grained Image Generation through Asymmetric Training |
| | CycleGAN - Unpaired Image-to-Image Translation using Cycle-Consistent Adversarial Networks |
| | DTN - Unsupervised Cross-Domain Image Generation |
| 6 | DCGAN - Unsupervised Representation Learning with Deep Convolutional Generative Adversarial Networks |
| | DiscoGAN - Learning to Discover Cross-Domain Relations with Generative Adversarial Networks |
| | DR-GAN - Disentangled Representation Learning GAN for Pose-Invariant Face Recognition |
| | DualGAN - DualGAN: Unsupervised Dual Learning for Image-to-Image Translation |
| | EBGAN - Energy-based Generative Adversarial Network |
| | f-GAN - f-GAN: Training Generative Neural Samplers using Variational Divergence Minimization |
| r Discovery | FF-GAN - Towards Large-Pose Face Frontalization in the Wild |
| | GAWWN - Learning What and Where to Draw |
| | GeneGAN - GeneGAN: Learning Object Transfiguration and Attribute Subspace from Unpaired Data |
| | Geometric GAN - Geometric GAN |
| | GoGAN - Gang of GANs: Generative Adversarial Networks with Maximum Margin Ranking |
| | GP-GAN - GP-GAN: Towards Realistic High-Resolution Image Blending |
| | IAN - Neural Photo Editing with Introspective Adversarial Networks |
| | iGAN - Generative Visual Manipulation on the Natural Image Manifold |
| | IcGAN - Invertible Conditional GANs for image editing |
| Calorimeters | ID-CGAN - Image De-raining Using a Conditional Generative Adversarial Network |
| | Improved GAN - Improved Techniques for Training GANs |
| | InfoGAN - InfoGAN: Interpretable Representation Learning by Information Maximizing Generative Adversaria |
| etworks | LAGAN - Learning Particle Physics by Example: Location-Aware Generative Adversarial Networks for Physics Synthesis |
| | LAPGAN - Deep Generative Image Models using a Laplacian Pyramid of Adversarial Networks |
| | https://github.com/hindupuravinash/the-gan-zoo |







2017: Explosion of GANs Better training and generation



LSGAN, Zhu 2017.



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







Progressive GAN, Karras 2018.

Some challenges with GANs ...

Challenges with GANs

• Vanishing gradients:

 \circ the discriminator becomes too good and the generator gradient vanishes.

Non-Convergence:

o the generator and discriminator oscillate without reaching an equilibrium.

Mode Collapse:

 \circ the generator distribution collapses to a small set of examples.

• Mode Dropping:

 \circ the generator distribution doesnt fully cover the data distribution.

Challenges with GANs: Vanishing Gradients

- The minimax objective saturates when D_{θ_d} is close to perfect: $V(\theta_d, \theta_g) = \mathbb{E}_{p_{data}} \left[\log D_{\theta_d}(\mathbf{x}) \right] + \mathbb{E}_{p_z(\mathbf{z})} \left[\log \left(1 - D_{\theta_d}(G_{\theta_g}(\mathbf{z})) \right) \right].$
- A non-saturating heuristic objective for the generator is

$$J(G_{\theta_g}) = -\mathbb{E}_{p_z(z)} \left[\log \left(D_{\theta_d} (G_{\theta_g}) \right) \right]$$



https://arxiv.org/abs/1701.00160

(**z**)))].

Challenges with GANs: Vanishing Gradients

- The minimax objective saturates when D_{θ_d} is close to perfect: $V(\theta_d, \theta_g) = \mathbb{E}_{p_{\text{data}}} \left[\log D_{\theta_d}(\mathbf{x}) \right] + \mathbb{E}_{p_z(\mathbf{z})} \left[\log \left(1 - D_{\theta_d}(G_{\theta_g}(\mathbf{z})) \right) \right].$
- A non-saturating heuristic objective for the generator is

$$J(G_{\theta_g}) = -\mathbb{E}_{p_z(z)} \left[\log \left(D_{\theta_d} (G_{\theta_g}(z) \right) \right]$$



https://arxiv.org/abs/1701.00160



Potential Solutions:

Explore other training 1. objectives?

Discriminator Capacity: 2.

- make it small?
- train it less?
- slow learning rate?

Learning Schedule: 3.

try to balance training G and D



Problems: Nonconvergence

Deep Learning models (in general) involve a single player

- The player tries to maximize its reward (minimize its loss).
- Use SGD (with Backpropagation) to find the optimal parameters.
- SGD has convergence guarantees (under certain conditions).
- Problem: With non-convexity, we might converge to local optima.

GANs instead involve two (or more) players

- Discriminator is trying to maximize its reward.
- Generator is trying to minimize Discriminator's reward.

- SGD was not designed to find the Nash equilibrium of a game.
- **Problem:** We might not converge to the Nash equilibrium at all.

 $\min_{G} L(G)$

 $\min_{G} \max_{D} V(D,G)$

Challenges with GANs: Non-Convergence

- Simultaneous gradient descent is not guaranteed to converge for minimax objectives.
- Goodfellow et al. only showed convergence when updates are made in the function space.
- The parameterization of D and G results in highly non-convex objective.
- In practice, training tends to oscillate updates undo each other!

Challenges with GANs: Non-Convergence

- Simultaneous gradient descent is not guaranteed to converge for minimax objectives.
- Goodfellow et al. only showed convergence when updates are made in the function space.
- The parameterization of D and G results in highly non-convex objective.
- In practice, training tends to oscillate – updates undo each other!

Potential Solutions (HACKS) https://github.com/soumith/ganhacks

How to Train a GAN? Tips and tricks to make GANs work

While research in Generative Adversarial Networks (GANs) continues to improve the fundamental stability of these models, we use a bunch of tricks to train them and make them stable day to day.

Here are a summary of some of the tricks.

Here's a link to the authors of this document

If you find a trick that is particularly useful in practice, please open a Pull Request to add it to the document. If we find it to be reasonable and verified, we will merge it in.

1. Normalize the inputs

- normalize the images between -1 and 1
- Tanh as the last layer of the generator output

2: A modified loss function

In GAN papers, the loss function to optimize G is min (log 1-D), but in practice folks practically use max log D

- Goodfellow et. al (2014)

In practice, works well:

3: Use a spherical Z

Dont sample from a Uniform distribution

4: BatchNorm

- generated images.
- standard deviation)

5: Avoid Sparse Gradients: ReLU, MaxPool

- the stability of the GAN game suffers if you have sparse gradients
- LeakyReLU = good (in both G and D)
- For Upsampling, use: PixelShuffle, ConvTranspose2d + stride PixelShuffle: <u>https://arxiv.org/abs/1609.05158</u>

6: Use Soft and Noisy Labels

- 0.0 and 0.3 (for example). Salimans et. al. 2016

7: DCGAN / Hybrid Models

Use DCGAN when you can. It works!

• because the first formulation has vanishing gradients early on

Flip labels when training generator: real = fake, fake = real

• Construct different mini-batches for real and fake, i.e. each mini-batch needs to contain only all real images or all

• when batchnorm is not an option use instance normalization (for each sample, subtract mean and divide by

• For Downsampling, use: Average Pooling, Conv2d + stride

• Label Smoothing, i.e. if you have two target labels: Real=1 and Fake=0, then for each incoming sample, if it is real, then replace the label with a random number between 0.7 and 1.2, and if it is a fake sample, replace it with

• make the labels the noisy for the discriminator: occasionally flip the labels when training the discriminator

8: Use stability tricks from RL

Experience Replay

- Keep a replay buffer of past generations and occassionally show them
- Keep checkpoints from the past of G and D and occassionaly swap them out for a few iterations
- All stability tricks that work for deep deterministic policy gradients
- See Pfau & Vinyals (2016)

9: Use the ADAM Optimizer

- optim.Adam rules!
- See Radford et. al. 2015
- Use SGD for discriminator and ADAM for generator

10: Track failures early

- D loss goes to 0: failure mode
- check norms of gradients: if they are over 100 things are screwing up
- when things are working, D loss has low variance and goes down over time vs having huge variance and spiking
- if loss of generator steadily decreases, then it's fooling D with garbage (says martin)

11: Dont balance loss via statistics (unless you have a good reason to)

- Dont try to find a (number of G / number of D) schedule to uncollapse training
- It's hard and we've all tried it.
- If you do try it, have a principled approach to it, rather than intuition

For example

| - income of the second | | | |
|------------------------|--|--|--|
| e lossD > A: | | | |
| ain D | | | |
| e lossG > B: | | | |
| ain G | | | |
| | | | |

12: If you have labels, use them

if you have labels available, training the discriminator to also classify the samples: auxillary GANs

13: Add noise to inputs, decay over time

- Add some artificial noise to inputs to D (Arjovsky et. al., Huszar, 2016)
- http://www.inference.vc/instance-noise-a-trick-for-stabilising-gan-training/ <u>https://openreview.net/forum?id=Hk4_qw5xe</u>
- adding gaussian noise to every layer of generator (Zhao et. al. EBGAN) Improved GANs: OpenAl code also has it (commented out)

14: [notsure] Train discriminator more (sometimes)

- especially when you have noise
- hard to find a schedule of number of D iterations vs G iterations

15: [notsure] Batch Discrimination

Mixed results

- Use an Embedding layer
- Add as additional channels to images
- Keep embedding dimensionality low and upsample to match image channel size

17: Use Dropouts in G in both train and test phase

- Provide noise in the form of dropout (50%).
- Apply on several layers of our generator at both training and test time
- https://arxiv.org/pdf/1611.07004v1.pdf





Challenges with GANs: Mode Collapse

The generator maps all z values to the x that is mostly likely to fool the discriminator.





Some real examples



Reed, S., et al. Generating interpretable images with controllable structure. Technical report, 2016. 2, 2016.

Challenges with GANs: Mode Collapse

Possible Solutions:

There are a large variety of divergence measures for distributions:

• **f-Divergences:** (e.g. Jensen-Shannon, Kullback-Leibler)

$$D_f (P ||Q) = \int_{\chi} q(\mathbf{x}) f(\frac{p(\mathbf{x})}{q(\mathbf{x})})$$

• GANs [2], f-GANs [7], and more.

• Integral Probability Metrics: (e.g. Earth Movers Distance, Maximum Mean Discrepancy)

$$\gamma_F(P ||Q) = \sup_{f \in F} \left| \int f dP - \int f dP \right|$$

• Wasserstein GANs [1], Fisher GANs [6], Sobolev GANs [5] and more.

Wasserstein GAN

Wasserstein GANs

WGAN

- If our data are on a low-dimensional manifold of a high dimensional space the model's manifold and the true data manifold can have a negligible intersection in practice
- KL divergence is undefined or infinite
- The loss function and gradients may not be continuous and well behaved
- The Earth Mover's Distance is well defined: Minimum transportation cost for making one pile of dirt (pdf/pmf) look like the other

WGAN

$$J^{(D)}(\theta^{(D)},\theta^{(G)}) = -\left[\mathbb{E}_{x \sim p_{data}} D(x) - \mathbb{E}_{z} D(x) - \mathbb{E}_{z} D(G^{(D)},\theta^{(G)})\right] = -\mathbb{E}_{z} D(G(z))$$

- Importantly, the discriminator is trained for many steps before the generator is updated
- Gradient-clipping is used in the discriminator to ensure D(x) has the Lipschitz continuity required by the theory
- The authors argue that this solves many training issues, including mode collapse

(G(z))

Additional Sources:

There are lots of excellent references on GANs:

- Sebastian Nowozins presentation at MLSS 2018: https://github.com/nowozin/mlss2018-madrid-gan
- NIPS 2016 tutorial on GANs by Ian Goodfellow: https://arxiv.org/abs/1701.00160
- A nice explanation of Wasserstein GANs by Alex Irpan: https://www.alexirpan.com/2017/02/22/wasserstein-gan.html

Conditional GANs

MNIST digits generated conditioned on their class label.

| $[1, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0] \longrightarrow$ | 0 | 0 | 0 | Ø | Ô | 0 | Ç | 0 | Ô | 0 | 0 | ۵ | 0 | | Ò | Ø | 5 | Q | Õ | 0 |
|-----------------------------------------------------|----|------|---|------------|----------------|-------------|----------------------------------------------------------------------------------------|---|---|----------|-----|---|----|----|----------------|------|--------------|----|----|---|
| $[0, 1, 0, 0, 0, 0, 0, 0, 0, 0] \longrightarrow$ | X | | 1 | 1 | | | 1 | - | f | 1 | | 1 | | ſ | ţ | ł, | ł | 1. | 1 | Į |
| $[0, 0, 1, 0, 0, 0, 0, 0, 0, 0] \longrightarrow$ | 3 | 2 | 2 | 2. | 2 | 2 | $\left\{ \begin{array}{c} \mathbf{r} \\ \mathbf{r} \\ \mathbf{r} \end{array} \right\}$ | 1 | | Ĵ. | No. | | 2 | 2 | 0 | 3 | 1 | 3 | 3 | Ĵ |
| $[0, 0, 0, 1, 0, 0, 0, 0, 0, 0] \longrightarrow$ | | ţ, î | | | ננגע | Ĩ. | J | 3 | | tai | 3 | 3 | 3 | | (\mathbf{v}) | | مورد مورد | ~ | 3 | 3 |
| $[0, 0, 0, 0, 1, 0, 0, 0, 0, 0] \longrightarrow$ | ¥ | | ş | 51 | ¥. | Ŷ | ų, | 4 | 4 | <u>Ģ</u> | ş | 4 | 4 | ¥ | ý, | 4 | 4 | 4 | 4 | 4 |
| $[0, 0, 0, 0, 0, 1, 0, 0, 0, 0] \longrightarrow$ | \$ | Ĭ.) | Ş | Ş | 5 | 5 | ŝ | 5 | 5 | 434 | E) | 5 | Ę, | L. | цС). | 5 | 5 | ÷1 | 11 | 5 |
| $[0, 0, 0, 0, 0, 0, 0, 1, 0, 0, 0] \longrightarrow$ | ÷ | b | 4 | C . | v | 10 | 1 | ģ | (| €; | 6 | 6 | 6 | 5 | 6 | al a | ¢ | Ĵ | 3 | 6 |
| $[0, 0, 0, 0, 0, 0, 0, 0, 1, 0, 0] \longrightarrow$ | 7 | 7 | 7 | 7 | , , , , | <u>لي</u> (| | 5 | 7 | 7 | 7 | 7 | 7 | | ŗ, | 7 | 7 | 7 | 7 | |
| $[0, 0, 0, 0, 0, 0, 0, 0, 0, 1, 0] \longrightarrow$ | 8 | (L) | 2 | \$ | 2 | 2 | 7 | ŧ | 8 | | z | | 8 | 8 | 2 | | су L) | 8 | 5 | |
| $[0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 1] \longrightarrow$ | 9 | | 4 | 9 | ey | 9 | 4 | 7 | 9 | 4 | 5 | 9 | 4 | 9 | 9 | 9 | 4 | Ŷ | 9 | 4 |

Figure 2 in the original paper.

Mirza, Mehdi, and Simon Osindero. Conditional generative adversarial nets. arXiv preprint arXiv:1411.1784 (2014).

Conditional GANs

- Simple modification to the original GAN framework that conditions the model on additional information for better multi-modal learning.
- Lends to many practical applications of GANs when we have explicit supervision available.

Image Credit: Figure 2 in Odena, A., Olah, C. and Shlens, J., 2016. Conditional image synthesis with auxiliary classifier GANs. arXiv preprint arXiv:1610.09585.

Mirza, Mehdi, and Simon Osindero. "Conditional generative adversarial nets". arXiv preprint arXiv:1411.1784 (2014).

Conditional GAN (Mirza & Osindero, 2014)

Image-to-Image Translation

Labels to Street Scene

Link to an interactive demo of this paper

Isola, P., Zhu, J. Y., Zhou, T., & Efros, A. A. "Image-to-image translation with conditional adversarial networks". arXiv preprint arXiv:1611.07004. (2016).

Labels to Facade

BW to Color

Figure 1 in the original paper.

Image-to-Image Translation

- Architecture: *DCGAN*-based architecture
- Training is conditioned on the images from the source domain.
- Conditional GANs provide an effective way to handle many complex domains without worrying about designing structured loss functions explicitly.

Isola, P., Zhu, J.Y., Zhou, T., & Efros, A.A. "Image-to-image translation with conditional adversarial networks". arXiv preprint arXiv:1611.07004. (2016).

Positive examples

G tries to synthesize fake

images that fool **D**

Negative examples

D tries to identify the fakes

Figure 2 in the original paper.

Text-to-Image Synthesis

Motivation

Given a text description, generate images closely associated.

Uses a conditional GAN with the generator and discriminator being condition on "dense" text embedding.

this small bird has a pink breast and crown, and black primaries and secondaries.

the flower has petals that are bright pinkish purple with white stigma

this magnificent fellow is almost all black with a red crest, and white cheek patch.

this white and yellow flower have thin white petals and a round yellow stamen

Figure 1 in the original paper.

Reed, S., Akata, Z., Yan, X., Logeswaran, L., Schiele, B., & Lee, H. "Generative adversarial text to image synthesis". ICML (2016).

Text-to-Image Synthesis

Positive Example: Real Image, Right Text

Reed, S., Akata, Z., Yan, X., Logeswaran, L., Schiele, B., & Lee, H. "Generative adversarial text to image synthesis". ICML (2016).

This flower has small, round violet petals with a dark purple center

Figure 2 in the original paper.

Negative Examples: Real Image, Wrong Text Fake Image, Right Text

Face Aging with Conditional GANs

- input image.
- categories.

Antipov, G., Baccouche, M., & Dugelay, J. L. (2017). "Face Aging With Conditional Generative Adversarial Networks". arXiv preprint arXiv:1702.01983.

Differentiating Feature: Uses an *Identity Preservation Optimization* using an auxiliary network to get a better approximation of the latent code (z*) for an

Latent code is then conditioned on a discrete (one-hot) embedding of age

Figure 1 in the original paper.

Face Aging with Conditional GANs

Figure 3 in the original paper.

Antipov, G., Baccouche, M., & Dugelay, J. L. (2017). "Face Aging With Conditional Generative Adversarial Networks". arXiv preprint arXiv:1702.01983.

Conditional GANs

Conditional Model Collapse

- Scenario observed when the Conditional GAN starts *ignoring* either the code (c) or the noise variables (z).
- This limits the diversity of images generated.

A man in a orange jacket with sunglasses and a hat ski down a hill.

This guy is in black trunks and swimming underwater.

A tennis player in a blue polo shirt is looking down at the green court.

Mirza, Mehdi, and Simon Osindero. Conditional generative adversarial nets. arXiv preprint arXiv:1411.1784 (2014).

GANs got stuck ...

Later in this course: Diffusion Models