Topics:

Generative Adversarial Networks

## CS 4644-DL / 7643-A ZSOLT KIRA

#### • Assignment 3

- Due March 8th 11:59pm EST
- Projects
  - Project proposal due March 14th

 Next Meta office hours 03/14 3pm ET on machine translation

W9: Mar 3	Generative Models (Part I): Generative Adversarial Networks	Generative Adversarial Networks ;
W9: Mar 5	Project Planning Session	
W10: Mar 10	Generative Models (Part II): Diffusion Models PS3/HW3 due Mar 8th 11:59pm (grace period Mar 10th), PS4/HW4 out (due Mar 30th)	
W10: Mar 12	Variational Autoencoders (VAEs) Project Check-in Due March 14th 11:59pm (grace period until March 16th)	Tutorial on Variational Autoencoders
W11: Mar 17	Spring Break	
W11: Mar 19	Spring Break	



#### Symmetry in Encoder/Decoder



#### **U-Net**

You can have skip connections to bypass bottleneck!



Ronneberger, et al., "U-Net: Convolutional Networks for Biomedical Image Segmentation", 2015



Single-shot detectors use an idea of **grids** as anchors, with different scales and aspect ratios around them

 Various tricks used to increase the resolution (decrease subsampling ratio)





Liu, et al., "SSD: Single Shot MultiBox Detector", 2015



## DEtector TRansformer - DETR overview



Slides by R. Q. FEITOSA



# Generative Models: Introduction











#### Traditional unsupervised learning methods:



Similar in deep learning, but from neural network/learning perspective





#### **Discriminative vs. Generative Models**

- Discriminative models model P(y|x)
  - Example: Model this via neural network, SVM, etc.
- Generative models model P(x)

Goodfellow, NeurIPS 2016 Tutorial: Generative Adversarial Networks





#### **Discriminative vs. Generative Models**

- Discriminative models model P(y|x)
  - Example: Model this via neural network, SVM, etc.
- Generative models model P(x)
- We can parameterize our model as  $P(x, \theta)$  and use maximum likelihood to optimize the parameters given an unlabeled dataset:  $\theta^* = \arg \max \prod_{m=1}^{m} p_{model}(x^{(i)}; \theta)$

$$P^* = \arg \max_{\boldsymbol{\theta}} \prod_{i=1}^m p_{\text{model}} \left( \boldsymbol{x}^{(i)}; \boldsymbol{\theta} \right)$$
$$= \arg \max_{\boldsymbol{\theta}} \log \prod_{i=1}^m p_{\text{model}} \left( \boldsymbol{x}^{(i)}; \boldsymbol{\theta} \right)$$
$$= \arg \max_{\boldsymbol{\theta}} \sum_{i=1}^m \log p_{\text{model}} \left( \boldsymbol{x}^{(i)}; \boldsymbol{\theta} \right).$$

- They are called generative because they can often generate samples
  - Example: Multivariate Gaussian with estimated parameters  $\mu, \sigma$

Goodfellow, NeurIPS 2016 Tutorial: Generative Adversarial Netvorks

**Generative Models** 





Goodfellow, NeurIPS 2016 Tutorial: Generative Adversarial Networks

**Generative Models** 



# PixelRNN & PixelCNN





Goodfellow, NeurIPS 2016 Tutorial: Generative Adversarial Networks





#### We can use chain rule to decompose the joint distribution

- Factorizes joint distribution into a product of conditional distributions
  - Similar to Bayesian Network (factorizing a joint distribution)
  - Similar to language models!

p

Same as language modeling!  

$$(x) = \prod_{i=1}^{n^2} p(x_i | x_1, \dots, x_{i-1})$$

$$p(s) = \prod_i p(W_i | W_{i-1}, \dots, W_1)$$

$$next$$
word
history
word

- Requires some ordering of variables (edges in a probabilistic graphical model)
- We can estimate this conditional distribution as a neural network

Oord et al., Pixel Recurrent Neural Networks



 Language modeling involves estimating a probability distribution over sequences of words.

$$p(\mathbf{s}) = p(w_1, w_2, \dots, w_n) = \prod_{\substack{i \\ wor}} p(w_i \mid w_{i-1}, \dots, w_1)$$

RNNs are a family of neural architectures for modeling sequences.









$$p(x) = \prod_{i=1}^{n^2} p(x_i | x_1, \dots, x_{i-1})$$
$$p(x) = p(x_1) \prod_{i=2}^{n^2} p(x_i | x_1, \dots, x_{i-1})$$

1. Choose ordering (upper left, top to bottom, left to right.

Separate out pixel 1

Oord et al., Pixel Recurrent Neural Networks



Georg



 $p(x) = p(x_1)p(x_2|x_1)p(x_3|x_1)\prod_{i=1}^{n^2} p(x_i|x_1, \dots, x_{i-1})$ 

- Model this as RNN with parameters
- Training:
  - We can train similar to language models:
  - Maximum likelihood approach

Downsides?

- Downsides:
  - Slow sequential generation process
  - Only considers few context pixels

Oord et al., Pixel Recurrent Neural Networks



**Factorized Models for Images** 





- Idea: Represent conditional distribution as a convolution layer!
  - Because of spatial locality in images
- Considers larger context (receptive field)
- Practically can be implemented by applying a mask, zeroing out "future" pixels
- Faster training but still slow generation
  - Limited to smaller images

Oord et al., Conditional Image Generation with PixelCNN Decoders





## occluded

## completions

## original



Oord et al., Conditional Image Generation with PixelCNN Decoders

Example Results: Image Completion (PixelRNN)





#### Geyser



Hartebeest



Grey whale



Tiger

#### Can we update this to modern times?

Oord et al., Conditional Image Generation with PixelCNN Decoders

Example Images (PixelCNN)





Chameleon: Mixed-Modal Early-Fusion Foundation Models

Multi/Mixed-Modal Large Language Models



Generative Adversarial Networks (GANs)





Goodfellow, NeurIPS 2016 Tutorial: Generative Adversarial Networks





• Implicit generative models do not actually learn an explicit model for p(x)

- Instead, learn to generate samples from p(x)
  - Learn good feature representations
  - Perform data augmentation
  - Learn world models (a simulator!) for reinforcement learning

How?

Decode architecture

- What architecture lets us generate images? How do we generate a different image every time?
- Learn to sample from a neural network output





• We would like to sample from p(x) using a neural network

#### Idea:

- Sample from a simple distribution (Gaussian)
- Transform the sample to p(x)



0



Input can be a vector with (independent) Gaussian random numbers

We can use a CNN to generate images!



#### How do we train this (loss)?





Instead, learn to generate samples from p(x)

How?

- Adversarial training that uses one network's predictions to train the other (dynamic loss function!)
- Lots of tricks to make the optimization more stable





- Goal: We would like to generate *realistic* images. How can we drive the network to learn how to do this?
- Idea: Have another network try to distinguish a real image from a generated (fake) image
  - Why? Signal can be used to determine how well it's doing at generation
    - Can be seen as a dynamic (adversarial) loss!













 $N(\mu, \sigma)$ 

**Neural Network** 

p(x)

### **Adversarial Networks**





#### Question: What loss functions can we use (for each network)?





Since we have two networks competing, this is a mini-max two player game

- Ties to game theory
- Not clear what (even local) Nash equilibria are for this game

Goodfellow, NeurIPS 2016 Tutorial: Generative Adversarial Networks





Since we have two networks competing, this is a mini-max two player game

- Ties to game theory
- Not clear what (even local) Nash equilibria are for this game
- The full mini-max objective is:

$$\begin{split} \min_{G} \max_{D} V(D,G) &= \mathbb{E}_{\boldsymbol{x} \sim p_{\mathsf{data}}(\boldsymbol{x})} [\log D(\boldsymbol{x})] + \mathbb{E}_{\boldsymbol{z} \sim p_{\boldsymbol{z}}(\boldsymbol{z})} [\log(1 - D(G(\boldsymbol{z})))] \\ & \text{Sample from real} \end{split}$$

Goodfellow, NeurIPS 2016 Tutorial: Generative Adversarial Networks





$$\min_{G} \max_{D} V(D,G) = \mathbb{E}_{\boldsymbol{x} \sim p_{\text{data}}(\boldsymbol{x})} [\log D(\boldsymbol{x})] + \mathbb{E}_{\boldsymbol{z} \sim p_{\boldsymbol{z}}(\boldsymbol{z})} [\log(1 - D(G(\boldsymbol{z})))]$$

where D(x) is the discriminator outputs probability ([0,1]) of real image
x is a real image and G(z) is a generated image

#### The discriminator wants to maximize this:

- D(x) is pushed up (to 1) because x is a real image
- 1 D(G(z)) is also pushed up to 1 (so that D(G(z)) is pushed down to 0)
- In other words, discriminator wants to classify real images as real (1) and fake images as fake (0)

## **Discriminator Perspective**



$$\min_{G} \max_{D} V(D,G) = \mathbb{E}_{\boldsymbol{x} \sim p_{\text{data}}(\boldsymbol{x})} [\log D(\boldsymbol{x})] + \mathbb{E}_{\boldsymbol{z} \sim p_{\boldsymbol{z}}(\boldsymbol{z})} [\log(1 - D(G(\boldsymbol{z})))]$$

where D(x) is the discriminator outputs probability ([0,1]) of real image
x is a real image and G(z) is a generated image

#### The generator wants to minimize this:

- First term: G(..) doesn't appear in it!
- 1 D(G(z)) is pushed down to 0 (so that D(G(z)) is pushed up to 1)
- This means that the generator is fooling the discriminator, i.e. succeeding at generating images that the discriminator can't discriminate from real

### **Generator Perspective**



Since we have two networks competing, this is a mini-max two player game

- Ties to game theory
- Not clear what (even local) Nash equilibria are for this game
- The full mini-max objective is:

Sample from fake

$$\min_{G} \max_{D} V(D,G) = \mathbb{E}_{\boldsymbol{x} \sim p_{\mathsf{data}}(\boldsymbol{x})} [\log D(\boldsymbol{x})] + \mathbb{E}_{\boldsymbol{z} \sim p_{\boldsymbol{z}}(\boldsymbol{z})} [\log(1 - D(G(\boldsymbol{z})))]$$

**Generator** *minimizes* 

How well discriminator does (0 for fake)

• where D(x) is the discriminator outputs probability ([0,1]) of real image

• x is a **real image** and G(z) is a **generated** image

Goodfellow, NeurIPS 2016 Tutorial: Generative Adversarial Networks

Mini-max Two Player Game

Since we have two networks competing, this is a mini-max two player game

Ties to game theory

Not clear what (even local) Nash equilibria are for this game

The full mini-max objective is: Sample from real Sample from fake  $\min_{G} \max_{D} V(D,G) = \mathbb{E}_{\boldsymbol{x} \sim p_{\mathsf{data}}(\boldsymbol{x})} [\log D(\boldsymbol{x})] + \mathbb{E}_{\boldsymbol{z} \sim p_{\boldsymbol{z}}(\boldsymbol{z})} [\log(1 - D(G(\boldsymbol{z})))]$ **Discriminator** *maximizes* How well discriminator How well discriminator does (1 for real) does (0 for fake) where D(x) is the discriminator outputs probability ([0,1]) of real image • x is a real image and G(z) is a generated image Goodfellow, NeurIPS 2016 Tutorial: Generative Adversarial Networks

Mini-max Two Player Game





**Generative Adversarial Networks (GANs)** 



The generator part of the objective does not have good gradient properties

 $\min_{G} \max_{D} V(D,G) = \mathbb{E}_{\boldsymbol{x} \sim p_{\mathsf{data}}(\boldsymbol{x})} [\log D(\boldsymbol{x})] + \mathbb{E}_{\boldsymbol{z} \sim p_{\boldsymbol{z}}(\boldsymbol{z})} [\log(1 - D(G(\boldsymbol{z})))]$ 

- High gradient when D(G(z)) is high (that is, discriminator is wrong)
- We want it to improve when samples are bad (discriminator is right)



Algorithm 1 Minibatch stochastic gradient descent training of generative adversarial nets. The number of steps to apply to the discriminator, k, is a hyperparameter. We used k = 1, the least expensive option, in our experiments.

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  $p_{\text{data}}(x)$ .
- Update the discriminator by ascending its stochastic gradient:

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

end for

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

$$\nabla_{\theta_g} \frac{1}{m} \sum_{i=1}^m \log\left(1 - D\left(G\left(\boldsymbol{z}^{(i)}\right)\right)\right).$$

#### end for

The gradient-based updates can use any standard gradient-based learning rule. We used momentum in our experiments.

Goodfellow, NeurIPS 2016 Generative Adversarial Nets







**Generative Adversarial Networks (GANs)** 







- Low-resolution images but look decent!
- Last column are nearest neighbor matches in dataset



c)

**Early Results** 



d)



GANs are very difficult to train due to the mini-max objective

#### Advancements include:

- More stable architectures
- Regularization methods to improve optimization
- Progressive growing/training and scaling

Goodfellow, NeurIPS 2016 Generative Adversarial Nets





Architecture guidelines for stable Deep Convolutional GANs

DCGAN

- 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



- Training GANs is difficult due to:
  - Minimax objective For example, what if generator learns to memorize training data (no variety) or only generates part of the distribution?
  - Mode collapse Capturing only some modes of distribution
- Several theoretically-motivated regularization methods
  - Simple example: Add noise to real samples!

$$\lambda \cdot \mathbb{E}_{x \sim P_{real}, \delta \sim N_d(0, cI)} \left[ \left\| \nabla_{\mathbf{x}} D_{\theta}(x + \delta) \right\| - k \right]^2$$

Kodali et al., On Convergence and Stability of GANs (also known as How to Train your DRAGAN)





#### Generative Adversarial Nets: Convolutional Architectures

DI9

Samples from the model look much better!

Radford et al, ICLR 2016



Slide Credit: Fei-Fei Li, Justin Johnson, Serena Yeung, CS 231n

#### **Generative Adversarial Nets: Convolutional Architectures**

Interpolating between random points in latent space

Radford et al, ICLR 2016







Brock et al., Large Scale GAN Training for High Fidelity Natural Image Synthesis

**Example Generated Images - BigGAN** 





(a) 128×128

(b) 256×256



1	EN.
10	D
10	•/

Figure 4: Samples from our model with truncation threshold 0.5 (a-c) and an example of class leakage in a partially trained model (d).

Brock et al., Large Scale GAN Training for High Fidelity Natural Image Synthesis







https://www.youtube.com/watch?v=PCBTZh41Ris





 Generative Adversarial Networks (GANs) can produce amazing images!

#### Several drawbacks

- High-fidelity generation heavy to train
- Training can be unstable
- No explicit model for distribution
- Larger number of extensions:
  - GANs conditioned on labels or other information
  - Adversarial losses for other applications





## **Comparison of Methods**



Gradually add Gaussian noise and then reverse