# Introduction

**Supervised Learning**

- Train Input: $\{X, Y\}$
- Learning output: $f: X \rightarrow Y, P(y|x)$
- e.g. classification
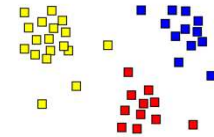
Sheep
Dog
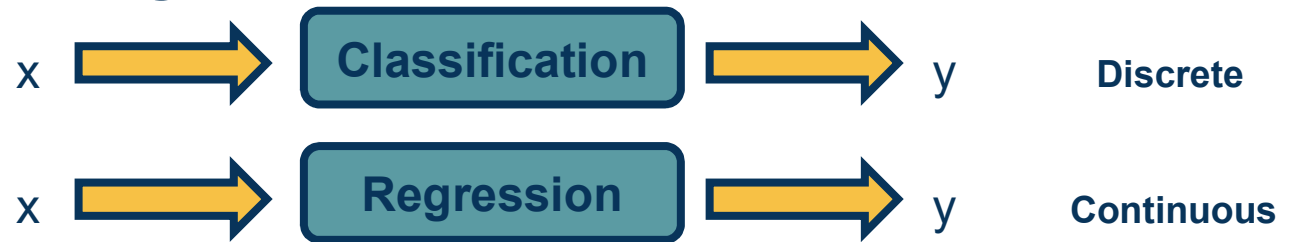**Cat**
Lion
Giraffe

**Less Labels**

**Unsupervised Learning**

- Input: $\{X\}$
- Learning output: $P(x)$
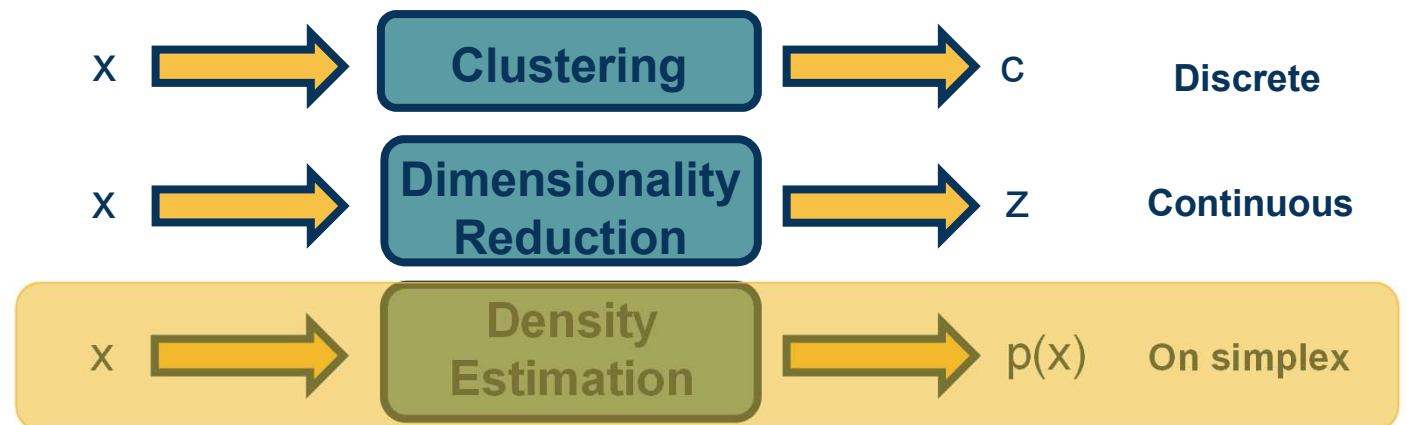- Example: Clustering, density estimation, etc.

# Supervised Learning

$x$ → **Classification** → $y$     **Discrete**

$x$ → **Regression** → $y$     **Continuous**

# Unsupervised Learning

$x$ → **Clustering** → $c$     **Discrete**

$x$ → **Dimensionality Reduction** → $z$     **Continuous**

$x$ → **Density Estimation** → $p(x)$     **On simplex**

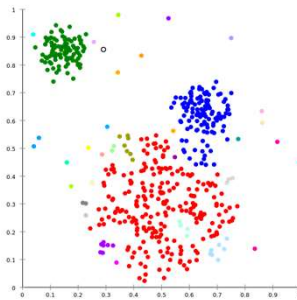Georgia Tech

**Traditional unsupervised learning methods:**



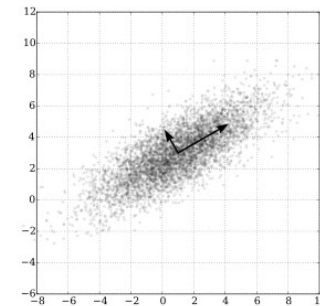Density estimation

**Modeling $P(x)$**

Deep Generative Models

Clustering

**Comparing/ Grouping**

Metric learning & clustering

Principal Component Analysis

**Representation Learning**

Almost all deep learning!

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

**What to Learn?**

## Discriminative vs. Generative Models

- Discriminative models model $P(y|x)$

    - Example: Model this via neural network, SVM, etc.

- Generative models model $P(x)$
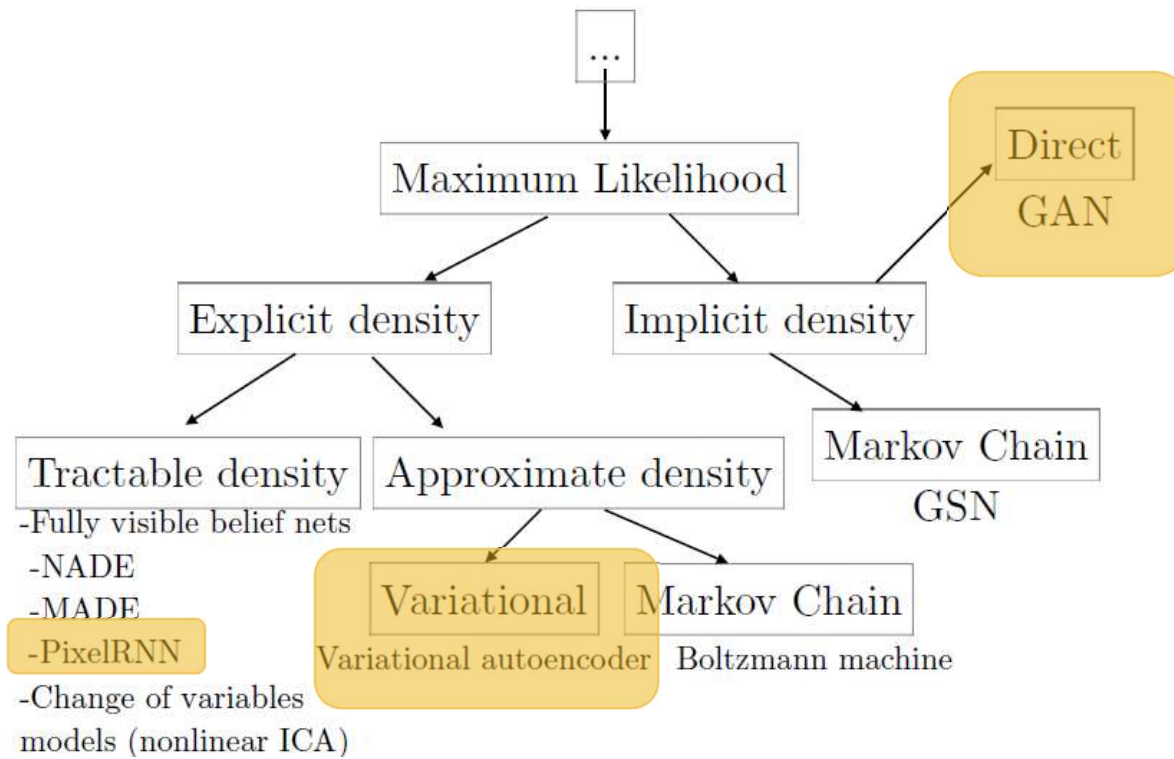
**Generative Models**

Georgia Tech

## 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:

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

- They are called generative because they can often generate *samples*

  - Example: Multivariate Gaussian with estimated parameters $\mu, \sigma$

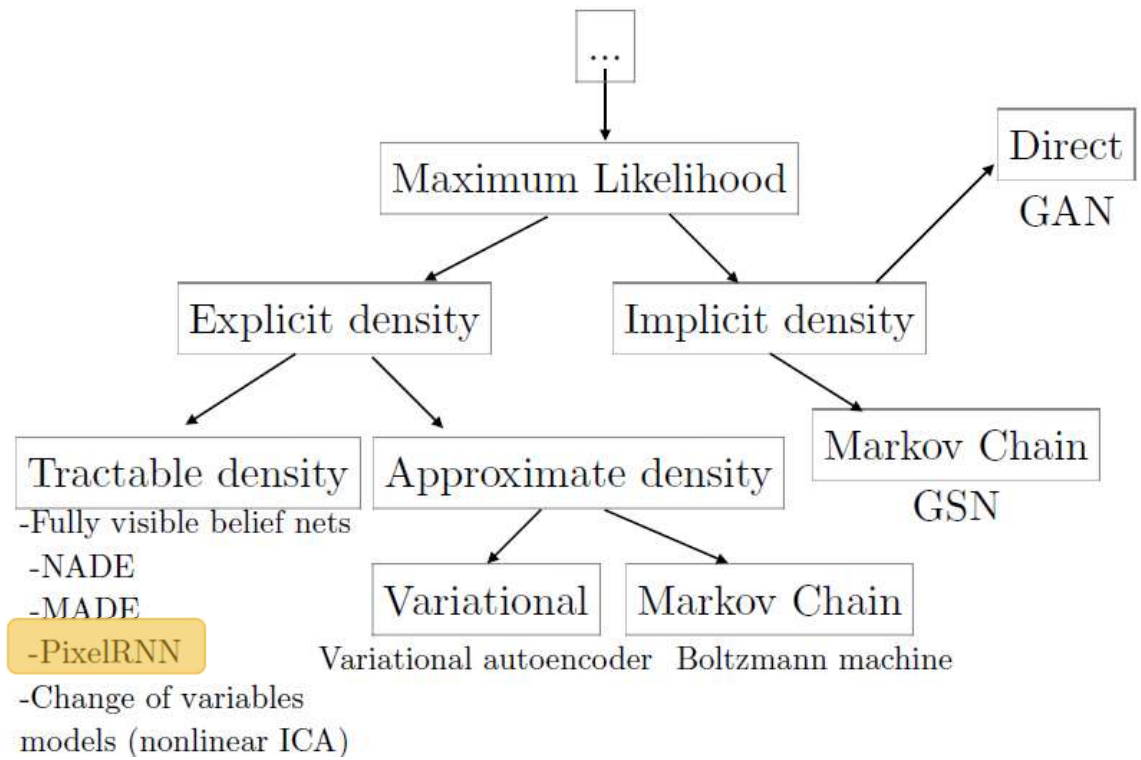*Goodfellow, NeurIPS 2016 Tutorial: Generative Adversarial Networks*

**Generative Models**

Georgia Tech

**Generative Models**

# PixelRNN  &
# PixelCNN

*Goodfellow, NeurIPS 2016 Tutorial: Generative Adversarial Networks*

**Generative Models**

**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!

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

- 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*

**Factorizing P(x)**

$$p(s) = p(w_1, w_2, \ldots, w_n)$$

$$= p(w_1) \, p(w_2 \mid w_1) \, p(w_3 \mid w_1, w_2) \cdots p(w_n \mid w_{n-1}, \ldots, w_1)$$

$$= \prod_i p(\underset{\text{next word}}{w_i} \mid \underset{\text{history}}{w_{i-1}, \ldots, w_1})$$

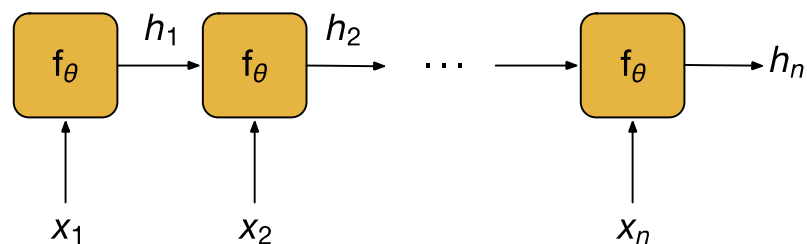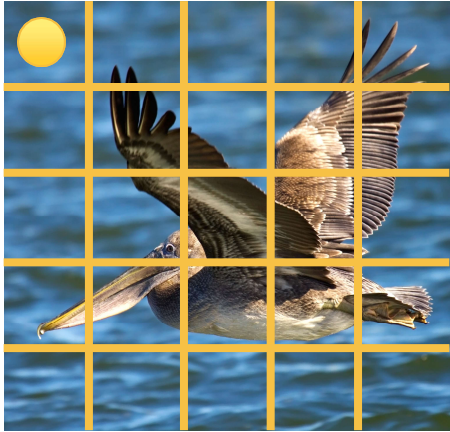Georgia Tech

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

$$p(\mathbf{s}) = p(w_1, w_2, \ldots, w_n) = \prod_i p(\underbrace{w_i}_{\substack{\text{next} \\ \text{wor} \\ \text{d}}} \mid \underbrace{w_{i-1}, \ldots, w_1}_{\text{history}})$$
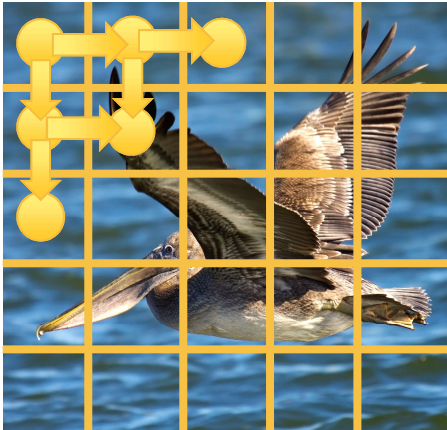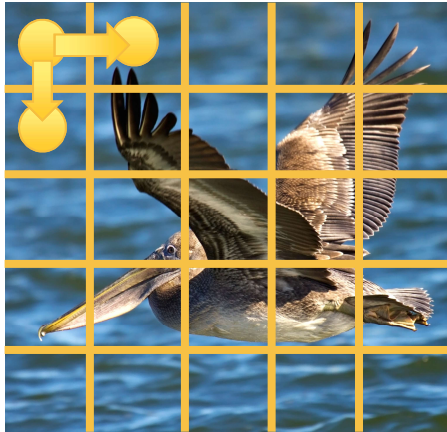
- RNNs are a family of neural architectures for modeling sequences.

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

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

*Oord et al., Pixel Recurrent Neural Networks*

**Factorized Models for Images**
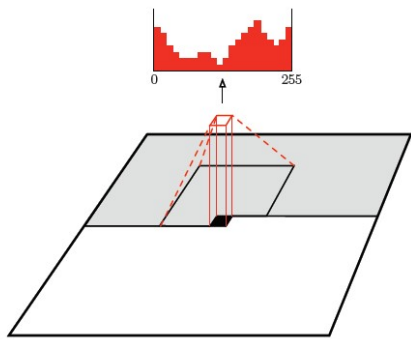
$$p(x) = p(x_1)p(x_2|x_1)p(x_3|x_1) \prod_{i=1}^{n^2} p(x_i|x_1, \ldots, x_{i-1})$$

- Training:
  - We can train similar to language models: Teacher/student forcing
  - Maximum likelihood approach

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

*Oord et al., Pixel Recurrent Neural Networks*

**Factorized Models for Images**

Georgia Tech

- **Idea:** Represent conditional distribution as a convolution layer!

- 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*

Georgia Tech

occluded            completions            original

*Oord et al., Conditional Image Generation with PixelCNN Decoders*

**Example Results: Image Completion (PixelRNN)**
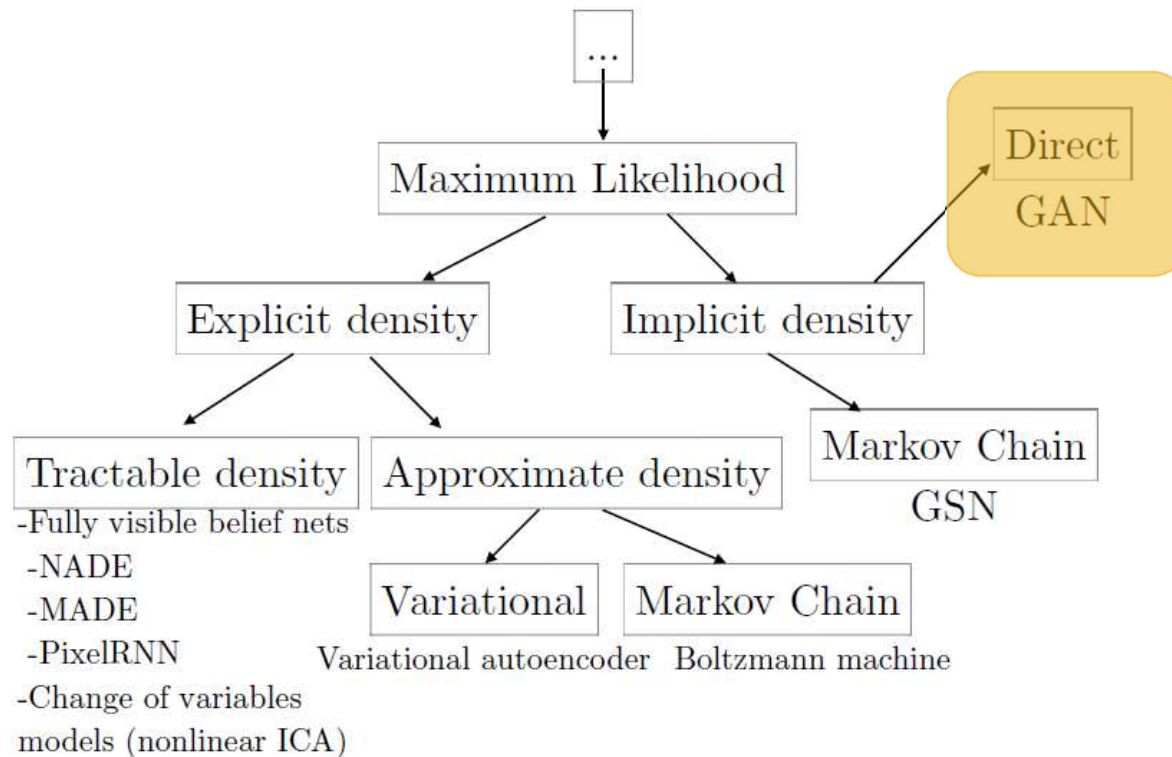
Geyser

Hartebeest

Grey whale

Tiger

*Oord et al., Conditional Image Generation with PixelCNN Decoders*

**Example Images (PixelCNN)**

Georgia Tech

Goodfellow, NeurIPS 2016 Tutorial: Generative Adversarial Networks

**Generative Models**

- *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?
  - **Learn to sample** from a neural network output
  - **Adversarial training** that uses one network's predictions to train the other (dynamic loss function!)
  - **Lots of tricks** to make the optimization more stable

**Implicit Models**

Georgia Tech

- We would like to *sample* from $p(x)$ using a neural network
- **Idea:**
  - Sample from a simple distribution (Gaussian)
  - Transform the sample to $p(x)$

Samples

Samples

$$N(\mu, \sigma)$$

**Neural Network**

$$p(x)$$

Georgia Tech

- Input can be a vector with (independent) Gaussian random numbers
- We can use a CNN to generate images!

**Generator**

**Vector of Random Numbers**

$N(\mu, \sigma)$  **Neural Network**  $p(x)$

- **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

**Generator**    **Discriminator**

**Vector of Random Numbers**



Real or Fake?

$N(\mu, \sigma)$    **Neural Network**    $p(x)$

**Generator:** Update weights to improve realism of generated images
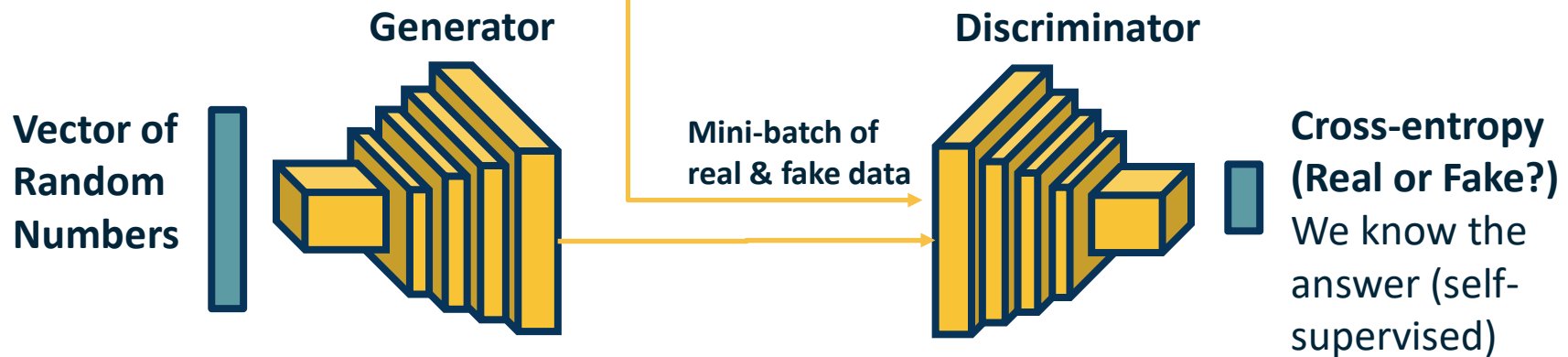
**Discriminator:** Update weights to better discriminate

**Vector of Random Numbers**

**Generator**

Mini-batch of real & fake data

**Discriminator**

**Cross-entropy (Real or Fake?)** We know the answer (self-supervised)

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

**Generative Adversarial Networks (GANs)**

Georgia Tech

- 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*

**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:

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

*Goodfellow, NeurIPS 2016 Tutorial: Generative Adversarial Networks*

**Mini-max Two Player Game**

Georgia Tech

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

- 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 $\text{D}(\text{G}(\text{z}))$ 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:
  - $1 - D(G(z))$ is pushed down to 0 (so that $\text{D}(\text{G}(z))$ 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**

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

**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*
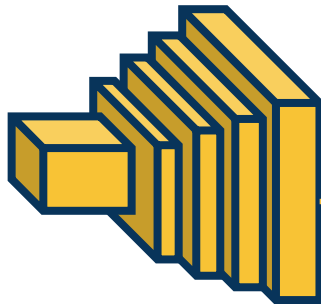
**Mini-max Two Player Game**

Georgia Tech

- 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**

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

**Discriminator *maximizes***

**How well discriminator does (1 for real)**   **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*
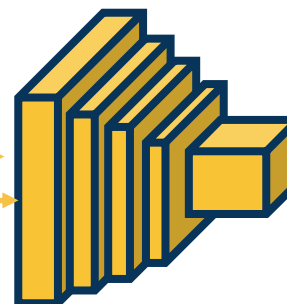
**Mini-max Two Player Game**

Georgia Tech

**Generator**

**Discriminator**

**Vector of Random Numbers**

**Mini-batch of real & fake data**

**Cross-entropy (Real or Fake?)** We know the answer (self-supervised)

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

**Generator Loss**

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

**Discriminator Loss**

**Generative Adversarial Networks (GANs)**

Georgia Tech

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

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

- 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)



- Alternative objective, *maximize*:

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

*Plot from CS231n, Fei-Fei Li, Justin Johnson, Serena Yeung*

**Converting to Max-Max Game**

Georgia Tech

**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:

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

    **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:

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

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

---

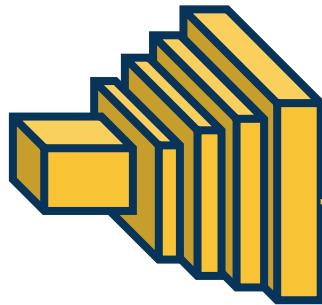*Goodfellow, NeurIPS 2016 Generative Adversarial Nets*
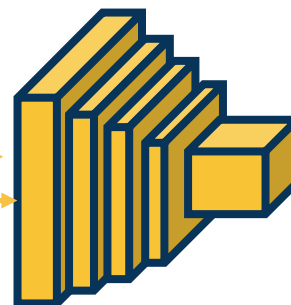
**Final Algorithm**

At the end, we have:
- An *implicit* generative model!
- Features from discriminator

**Vector of Random Numbers**

**Generator**

**Mini-batch of real & fake data**

**Discriminator**

**Cross-entropy (Real or Fake?)** We know the answer (self-supervised)

## Generative Adversarial Networks (GANs)

a)

b)

c)

d)

- Low-resolution images but look decent!

- Last column are nearest neighbor matches in dataset

**Early Results**

Georgia Tech

- 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*

Georgia Tech

## 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*

**DCGAN**

- 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!

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

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

**Regularization**

Georgia Tech

# 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
latent space



Radford et al,
ICLR 2016

Georgia
Tech

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

**Example Generated Images - BigGAN**

(a) 128×128    (b) 256×256    (c) 512×512    (d)

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).

**Failure Examples - BigGAN**

Georgia Tech

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

**Video Generation**

- A few other examples:
  - Deep nostalgia: https://www.myheritage.com/deep-nostalgia
  - High-resolution outputs: https://compvis.github.io/taming-transformers/

# GANs

Don't work with an explicit density function
Take game-theoretic approach: learn to generate from training distribution through 2-player game

Pros:
- Beautiful, state-of-the-art samples!

Cons:
- Trickier / more unstable to train
- Can't solve inference queries such as p(x), p(z|x)

Active areas of research:
- Better loss functions, more stable training (Wasserstein GAN, LSGAN, many others)
- Conditional GANs, GANs for all kinds of applications

Georgia Tech

# Mode Collapse

- Optimization of GANs is tricky
  - Not guaranteed to find Nash equilibrium

- Large number of methods to combat:
  - Use history of discriminators
  - Regularization
  - Different divergence measures

# Application: Data Augmentation

Georgia Tech

# Application: Domain Adaptation

- **Idea:** Train a model on *source* data and adapt to *target* data using unlabeled examples from target

# Approach



| | MNIST → USPS | USPS → MNIST | SVHN → MNIST |
|---|---|---|---|
| **Method** | | | |
| Source only | $0.752 \pm 0.016$ | $0.571 \pm 0.017$ | $0.601 \pm 0.011$ |
| Gradient reversal | $0.771 \pm 0.018$ | $0.730 \pm 0.020$ | 0.739 [16] |
| Domain confusion | $0.791 \pm 0.005$ | $0.665 \pm 0.033$ | $0.681 \pm 0.003$ |
| CoGAN | $0.912 \pm 0.008$ | $0.891 \pm 0.008$ | did not converge |
| ADDA (Ours) | $0.894 \pm 0.002$ | $0.901 \pm 0.008$ | $0.760 \pm 0.018$ |

Table 2: Experimental results on unsupervised adaptation among MNIST, USPS, and SVHN.

Georgia Tech

# Aside: Other ways to Align

- 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

Georgia Tech