Artificial faces synthesized by StyleGAN (Nvidia)

# COMP541

# DEEP LEARNING
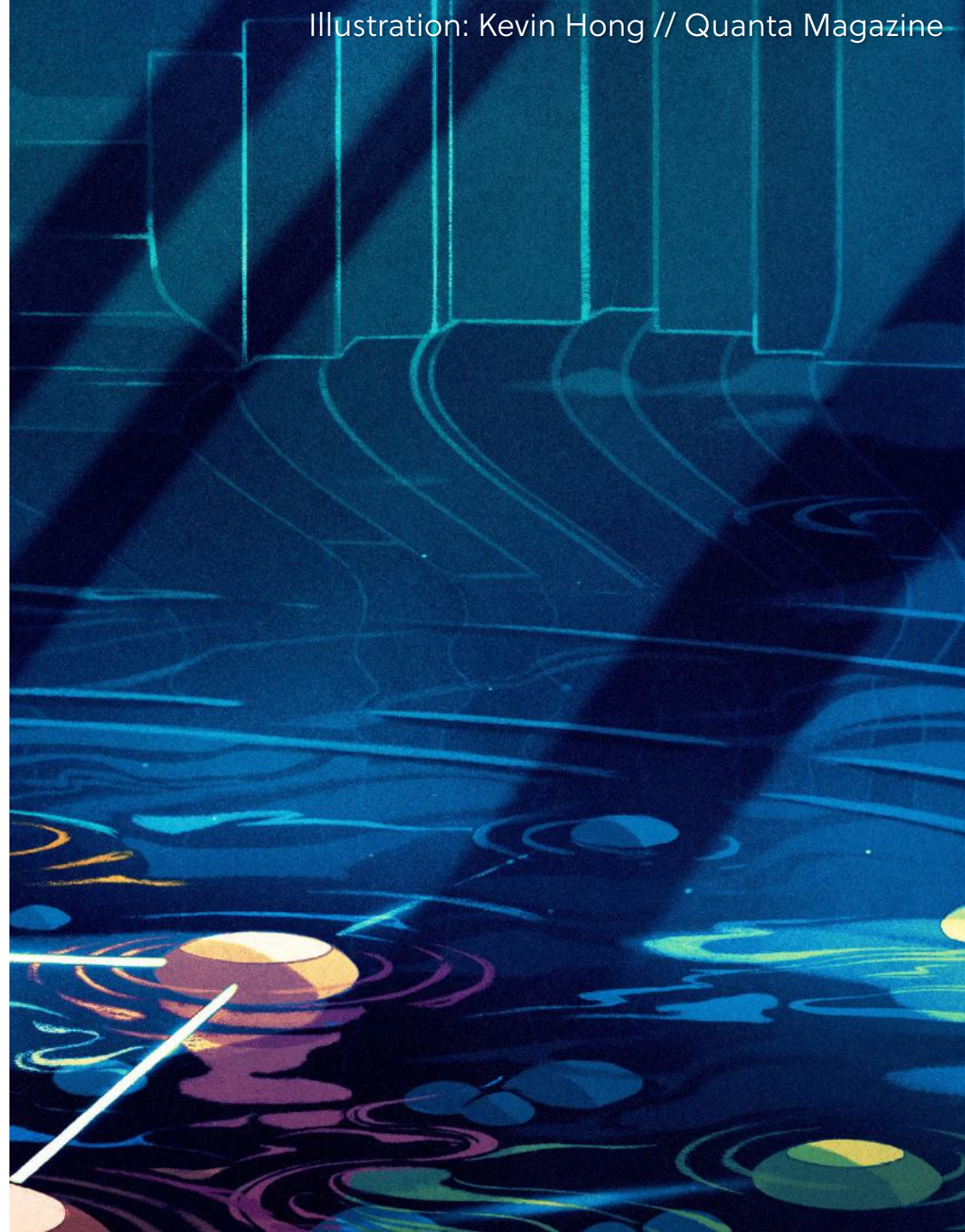
## Lecture #10 –Generative Adversarial Networks

Aykut Erdem // Koç University // Fall 2023

KOÇ UNIVERSITY

# Previously on COMP541

- graph structured data
- graph neural nets (GNNs)
- GNNs for "classical" network problems

# Lecture overview

- supervised vs unsupervised learning

- generative modeling

- basic foundations
  - sparse coding
  - autoencoders

- generative adversarial networks (GANs)

Disclaimer: Some of the material and slides for this lecture were borrowed from

—Justin Johnson's EECS 498/598 class

—Ruslan Salakhutdinov's talk titled "Unsupervised Learning: Learning Deep Generative Models"

—Ian Goodfellow's tutorial on "Generative Adversarial Networks"

—Aaron Courville's IFT6135 class

# Supervised vs Unsupervised Learning

**Supervised Learning**

**Data:** (x, y)

x is data, y is label

**Goal:** Learn a function to map x → y

**Examples:** classification, regression, object detection, semantic segmentation, image captioning, sentiment analysis, etc.

Classification



Cat

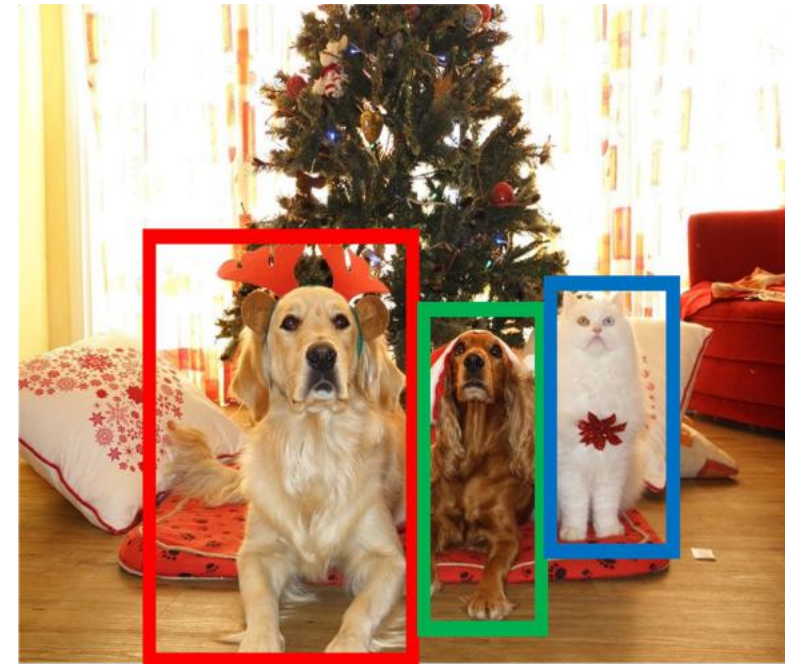# Supervised vs Unsupervised Learning

**Supervised Learning**

**Data:** (x, y)

x is data, y is label

**Goal:** Learn a function to map x → y

**Examples:** classification, regression, object detection, semantic segmentation, image captioning, sentiment analysis, etc.

Object Detection



**DOG, DOG, CAT**

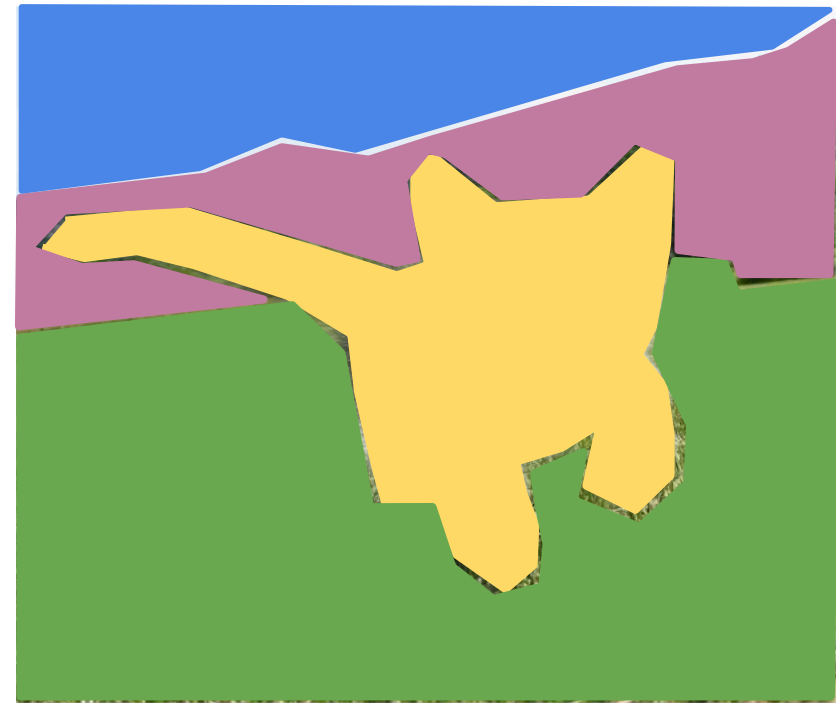# Supervised vs Unsupervised Learning

## Supervised Learning

**Data:** (x, y)

x is data, y is label

**Goal:** Learn a function to map x → y

**Examples:** classification, regression, object detection, semantic segmentation, image captioning, sentiment analysis, etc.

### Semantic Segmentation



**GRASS**, **CAT**, **TREE**, **SKY**

# Supervised vs Unsupervised Learning

**Supervised Learning**

**Data:** ⟨x, y⟩

x is data, y is label

**Goal:** Learn a function to map x → y

**Examples:** classification, regression, object detection, semantic segmentation, image captioning, sentiment analysis, etc.

Image captioning



*A cat sitting on a suitcase on the floor*

# Supervised vs Unsupervised Learning

## Supervised Learning

**Data:** (x, y)

x is data, y is label

**Goal:** Learn a function to map x → y

**Examples:** classification, regression, object detection, semantic segmentation, image captioning, sentiment analysis, etc.

## Sentiment Analysis

"This Movie is amazing.
It has a great plot and
talented actors, and
the supporting cast is
really good as well."

👍

# Supervised vs Unsupervised Learning

## Supervised Learning

**Data:** (x, y)

x is data, y is label

**Goal:** Learn a function to map x → y

**Examples:** classification, regression, object detection, semantic segmentation, image captioning, sentiment analysis, etc.
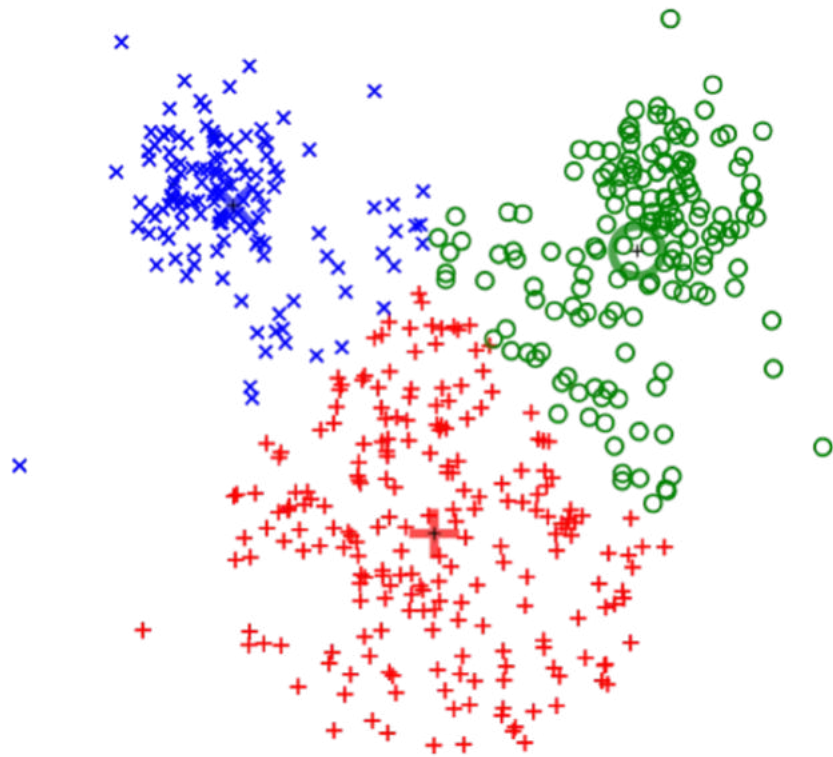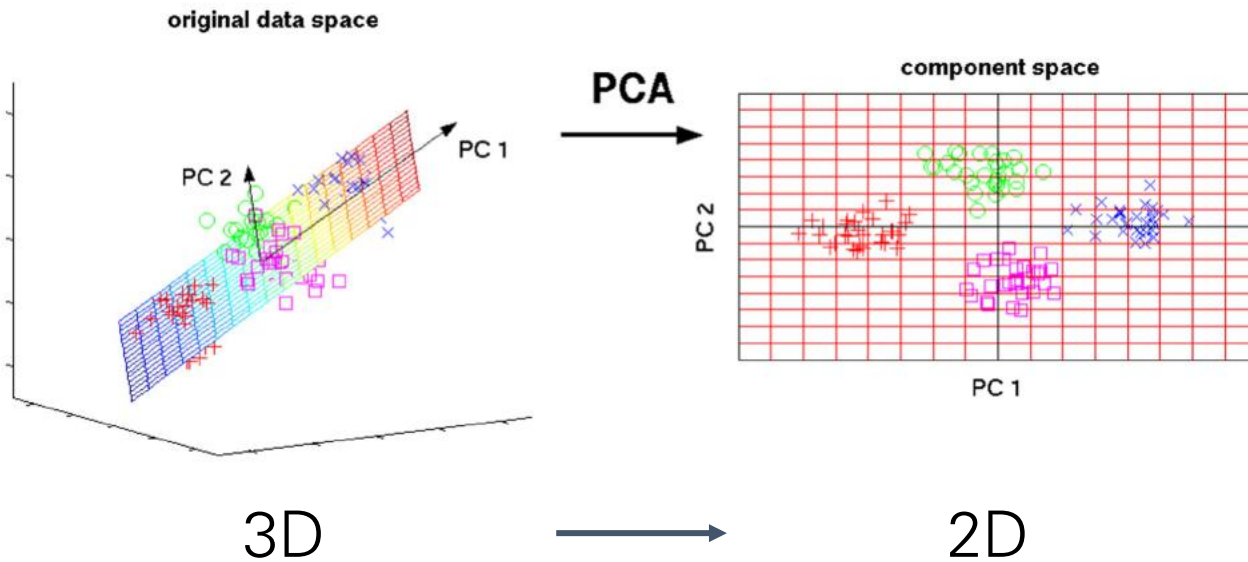
## Unsupervised Learning

**Data:** x

Just data, no labels!

**Goal:** Learn some underlying hidden structure of the data

**Examples:** clustering, dimensionality reduction, feature learning, density estimation, etc.

# Supervised vs Unsupervised Learning

Clustering
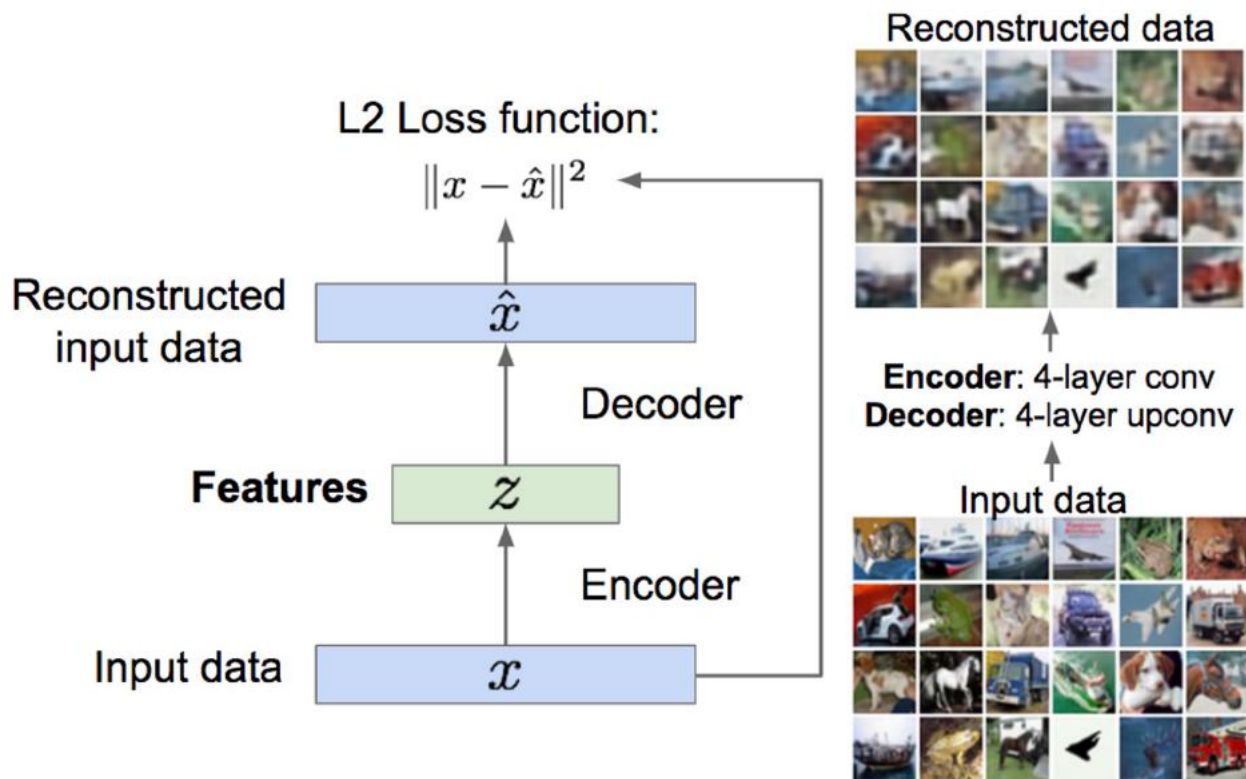(e.g. K-Means)



**Unsupervised Learning**

**Data:** x

Just data, no labels!

**Goal:** Learn some underlying hidden structure of the data

**Examples:** clustering, dimensionality reduction, feature learning, density estimation, etc.

# Supervised vs Unsupervised Learning

### Dimensionality Reduction
### (e.g. Principal Components Analysis)


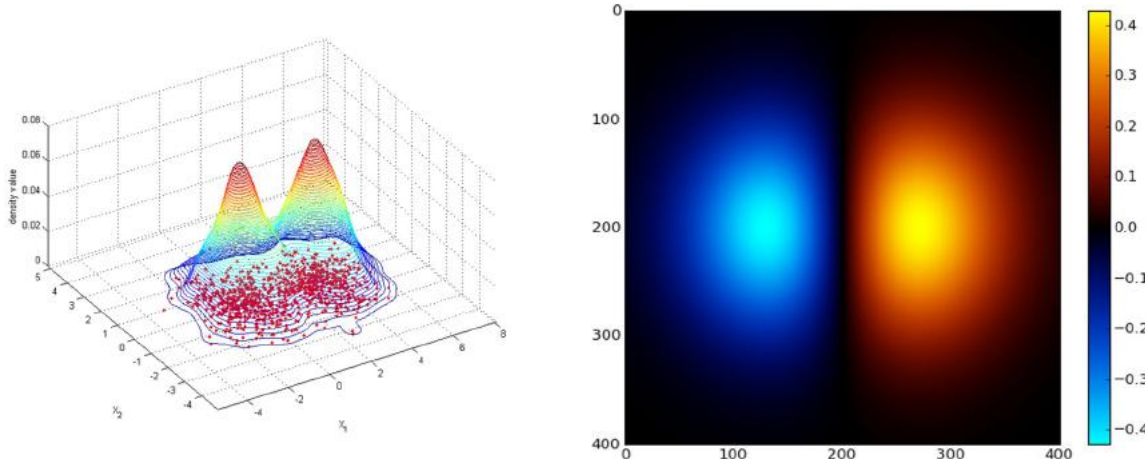
3D ⟶ 2D

**Unsupervised Learning**

**Data:** x

Just data, no labels!

**Goal:** Learn some underlying hidden structure of the data

**Examples:** clustering, dimensionality reduction, feature learning, density estimation, etc.

# Supervised vs Unsupervised Learning

### Feature Learning
### (e.g. autoencoders)

L2 Loss function:

$$\|x - \hat{x}\|^2$$

Reconstructed input data — $\hat{x}$

Decoder

**Features** — $z$

Encoder

Input data — $x$

Reconstructed data

**Encoder:** 4-layer conv
**Decoder:** 4-layer upconv

Input data

**Unsupervised Learning**

**Data:** x

Just data, no labels!

**Goal:** Learn some underlying hidden structure of the data

**Examples:** clustering, dimensionality reduction, feature learning, density estimation, etc.

# Supervised vs Unsupervised Learning

### Density Estimation



## Unsupervised Learning

**Data:** x

Just data, no labels!

**Goal:** Learn some underlying hidden structure of the data

**Examples:** clustering, dimensionality reduction, feature learning, density estimation, etc.

# Supervised vs Unsupervised Learning

## Supervised Learning

**Data:** (x, y)

x is data, y is label

**Goal:** Learn a function to map x → y

**Examples:** classification, regression, object detection, semantic segmentation, image captioning, sentiment analysis, etc.

## Unsupervised Learning

**Data:** x

Just data, no labels!

**Goal:** Learn some underlying hidden structure of the data

**Examples:** clustering, dimensionality reduction, feature learning, density estimation, etc.

# Discriminative vs Generative Models

**Discriminative Model:**
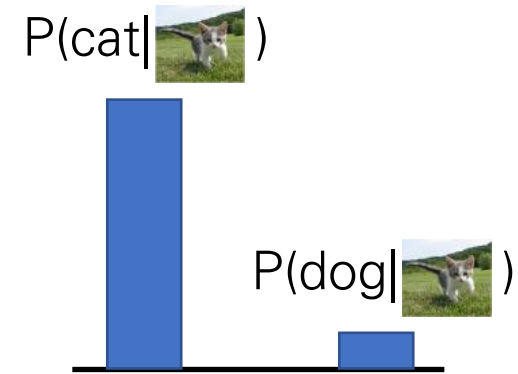Learn a probability
distribution p(y|x)

**Generative Model**:
Learn a probability
distribution p(x)

**Conditional
Generative Model:**
Learn p(x|y)

Data: x



Label: y

Cat

# Discriminative vs Generative Models

**Discriminative Model:**
Learn a probability distribution p(y|x)

**Generative Model**:
Learn a probability distribution p(x)

**Conditional Generative Model:**
Learn p(x|y)

Data: x

Label: y
Cat

**Probability Recap:**

**Density Function**
p(x) assigns a positive number to each possible x; higher numbers mean x is more likely
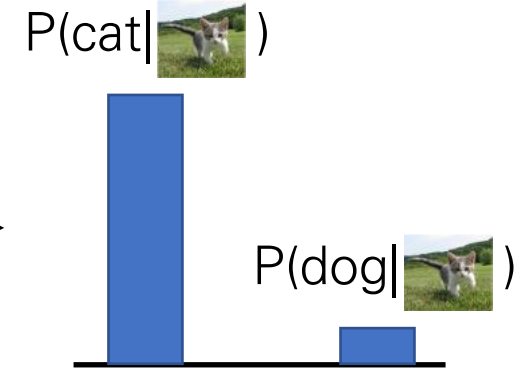
Density functions are **normalized**:

$$\int_X p(x)dx = 1$$

Different values of x **compete** for density
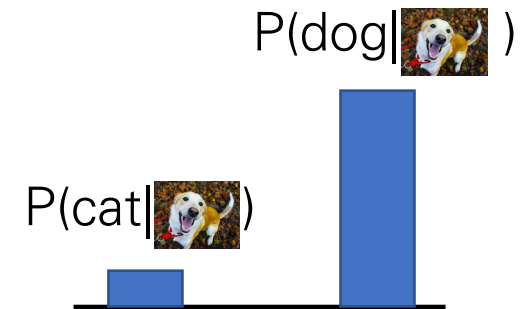
# Discriminative vs Generative Models

**Discriminative Model:**
Learn a probability distribution p(y|x)

**Generative Model:**
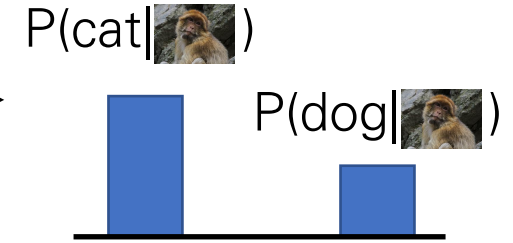Learn a probability distribution p(x)

**Conditional Generative Model:**
Learn p(x|y)

Data: x



P(cat| )

P(dog| )

**Density Function**
p(x) assigns a positive number to each possible x; higher numbers mean x is more likely

Density functions are **normalized**:

$$\int_X p(x)dx = 1$$

Different values of x **compete** for density

# Discriminative vs Generative Models

Discriminative Model:
Learn a probability
distribution p(y|x)

Generative Model:
Learn a probability
distribution p(x)

Conditional
Generative Model:
Learn p(x|y)

P(cat|  )

P(dog|  )

P(dog|  )

P(cat|  )

Discriminative model: the possible labels for each input "compete" for probability mass. But no competition between **images**
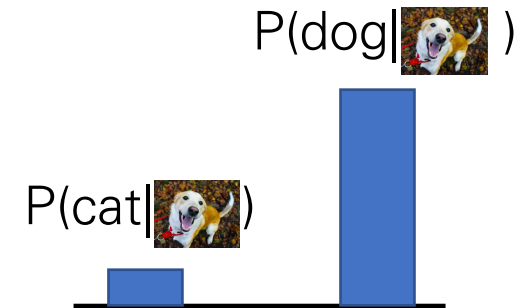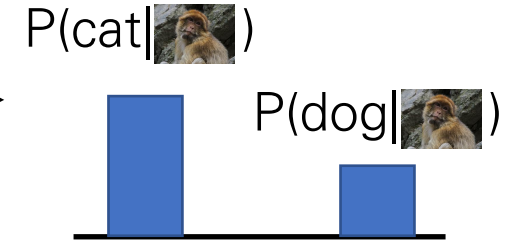
# Discriminative vs Generative Models

Discriminative Model:
Learn a probability
distribution p(y|x)

Generative Model:
Learn a probability
distribution p(x)

Conditional
Generative Model:
Learn p(x|y)



P(cat| )

P(dog| )

P(dog| )

P(cat| )

Discriminative model: No way for the model to handle unreasonable inputs; it must give label distributions for all images
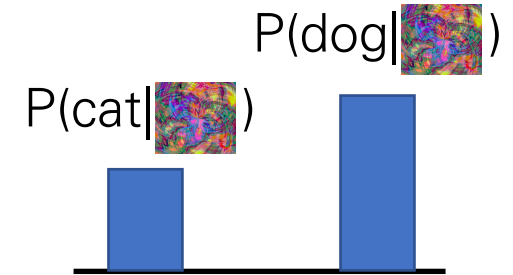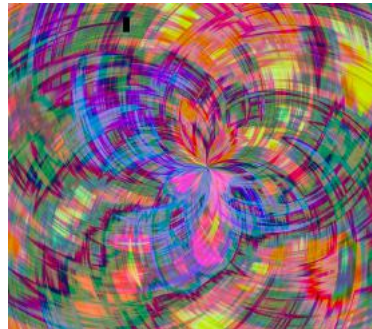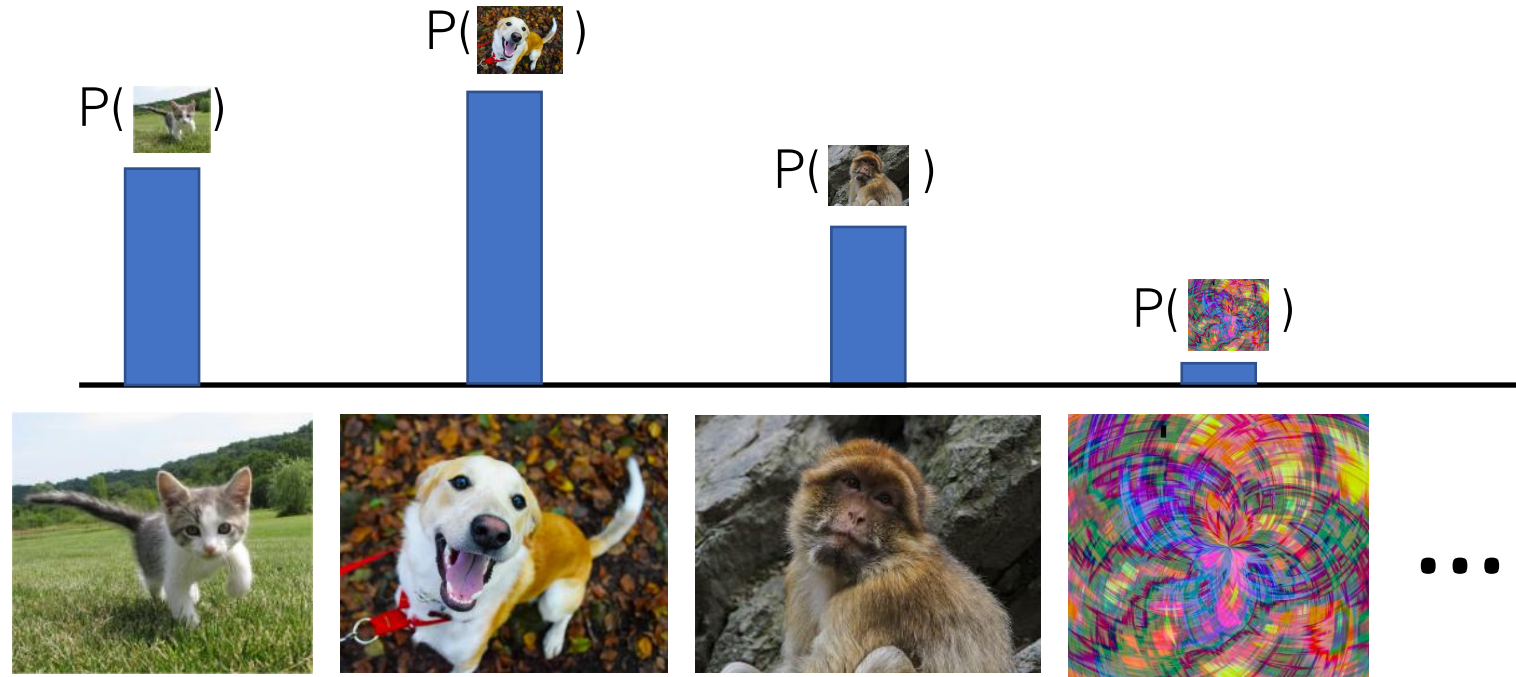
# Discriminative vs Generative Models

Discriminative Model:
Learn a probability
distribution p(y|x)

Generative Model:
Learn a probability
distribution p(x)

Conditional
Generative Model:
Learn p(x|y)



P(cat| )

P(dog| )

P(dog| )

P(cat| )

Discriminative model: No way for the model to handle unreasonable inputs; it must give label distributions for all images

# Discriminative vs Generative Models

**Discriminative Model:**
Learn a probability distribution p(y|x)

**Generative Model:**
Learn a probability distribution p(x)

**Conditional Generative Model:**
Learn p(x|y)



Generative model: All possible images compete with each other for probability mass

# Discriminative vs Generative Models

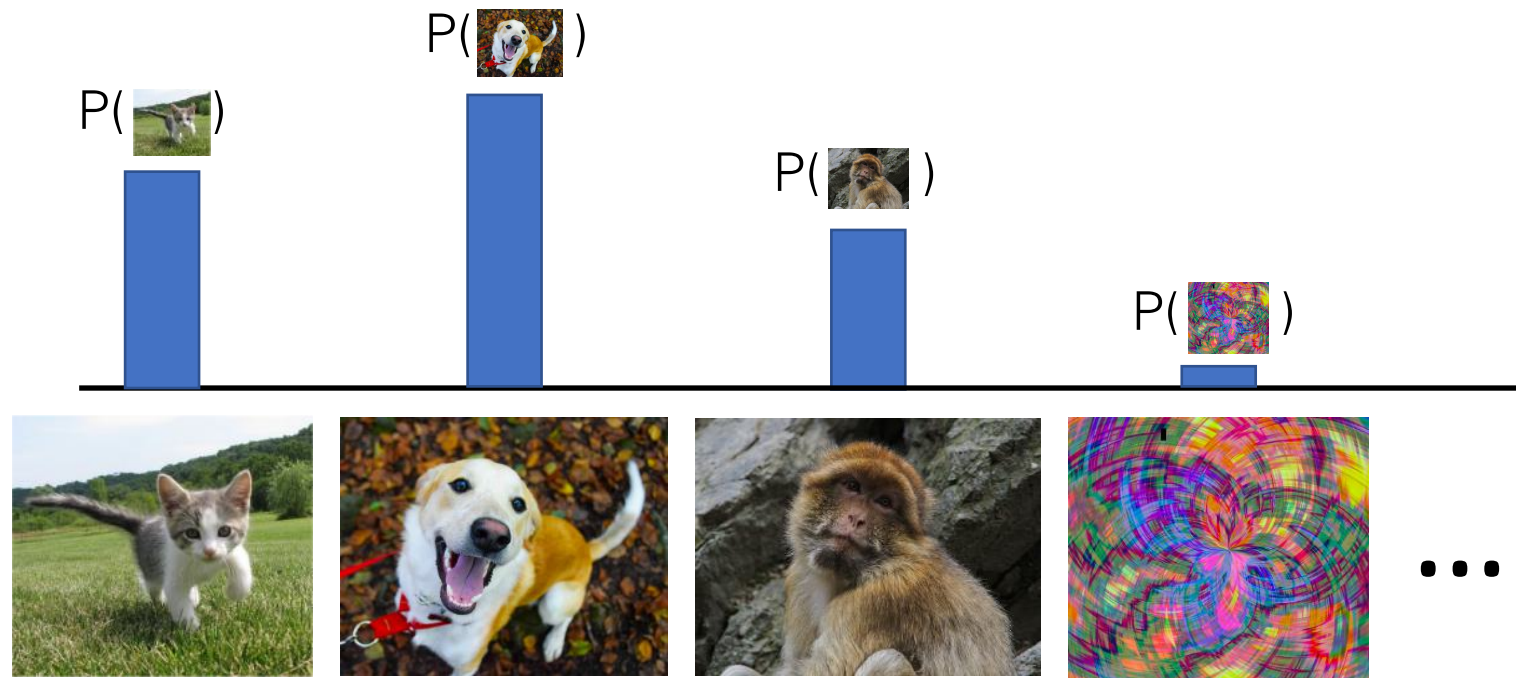**Discriminative Model:**
Learn a probability distribution p(y|x)

**Generative Model:**
Learn a probability distribution p(x)

**Conditional Generative Model:**
Learn p(x|y)



Generative model: All possible images compete with each other for probability mass

Requires deep image understanding! Is a dog more likely to sit or stand? How about 3-legged dog vs 3-armed monkey?

# Discriminative vs Generative Models

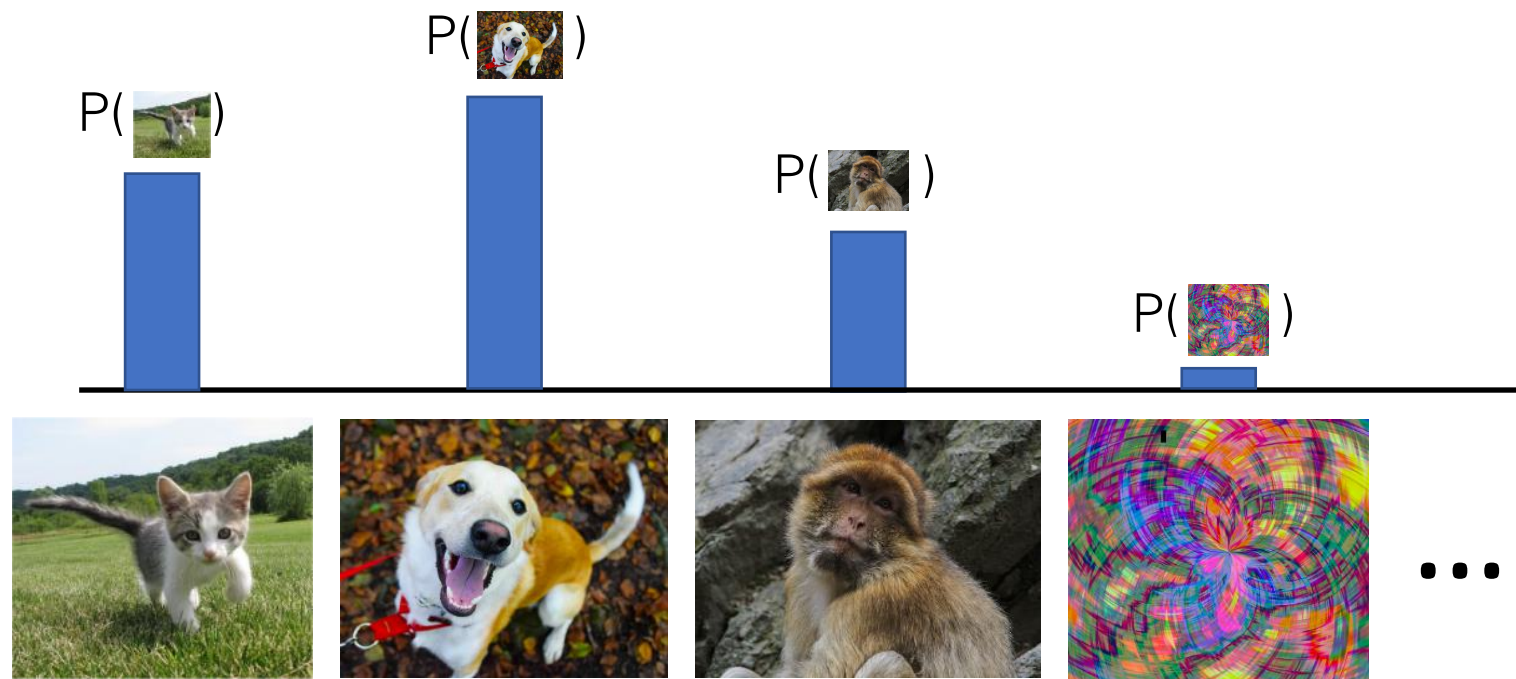**Discriminative Model:**
Learn a probability distribution p(y|x)

**Generative Model:**
Learn a probability distribution p(x)

**Conditional Generative Model:**
Learn p(x|y)



Generative model: All possible images compete with each other for probability mass

Model can "reject" unreasonable inputs by assigning them small values
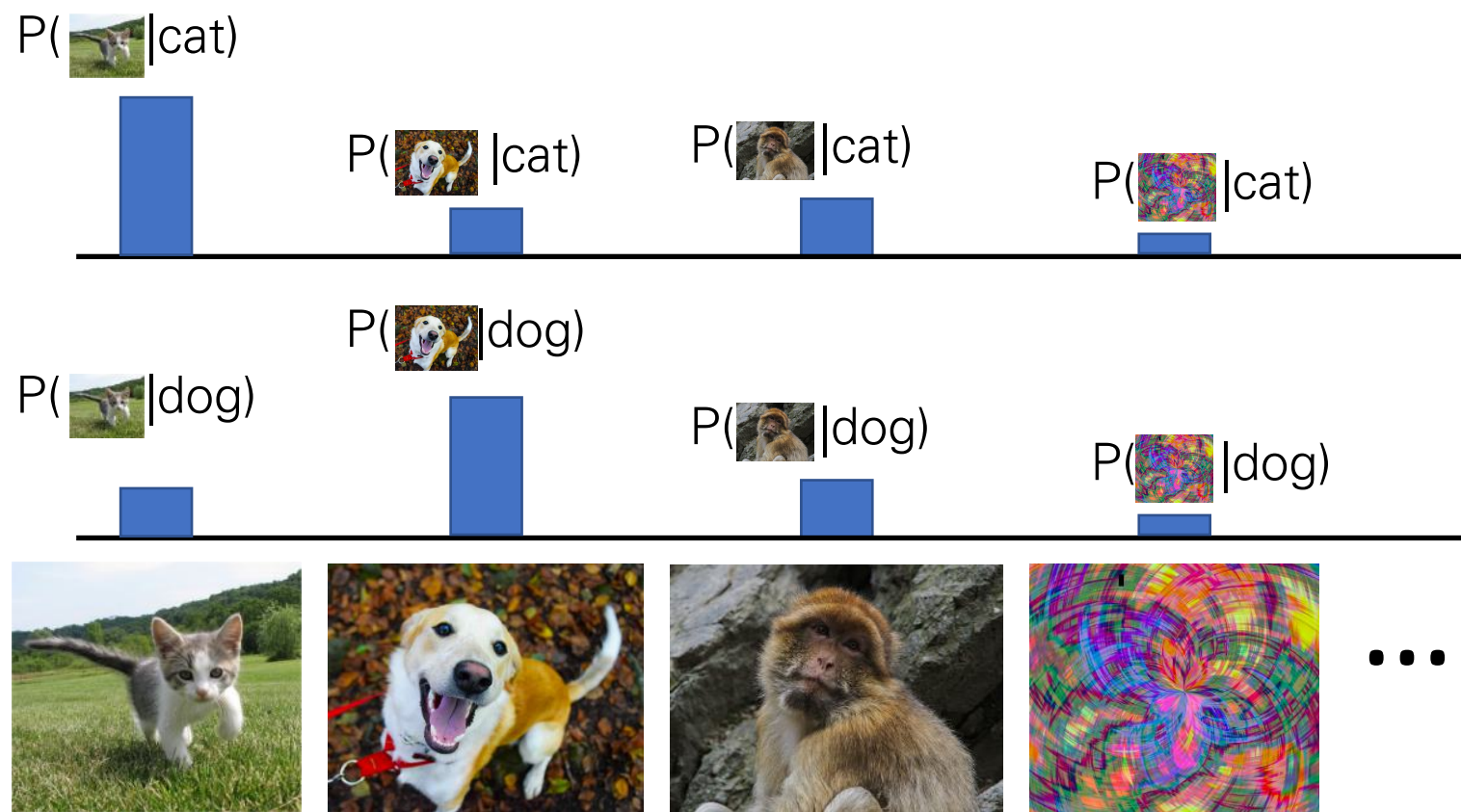
# Discriminative vs Generative Models

**Discriminative Model:**
Learn a probability distribution p(y|x)

**Generative Model:**
Learn a probability distribution p(x)

**Conditional Generative Model:**
Learn p(x|y)



Conditional Generative Model: Each possible label induces a competition among all images

# Discriminative vs Generative Models

**Discriminative Model:**
Learn a probability
distribution p(y|x)

**Generative Model**:
Learn a probability
distribution p(x)

**Conditional
Generative Model:**
Learn p(x|y)

Recall **Bayes' Rule:**

$$P(x \mid y) = \frac{P(y \mid x)}{P(y)} P(x)$$

# Discriminative vs Generative Models

**Discriminative Model:**
Learn a probability
distribution p(y|x)

**Generative Model**:
Learn a probability
distribution p(x)

**Conditional
Generative Model:**
Learn p(x|y)

Recall **Bayes' Rule:**

Discriminative Model

(Unconditional)
Generative Model

$$P(x \mid y) = \frac{P(y \mid x)}{P(y)} P(x)$$

Conditional
Generative Model

Prior over labels

We can build a conditional generative
model from other components!

# What can we do with a discriminative model?

**Discriminative Model:**
Learn a probability
distribution $p(y|x)$

➡ Assign labels to data
Feature learning (with labels)

**Generative Model**:
Learn a probability
distribution $p(x)$

**Conditional
Generative Model:**
Learn $p(x|y)$

# What can we do with a discriminative model?

**Discriminative Model:**
Learn a probability distribution p(y|x)

⟶ Assign labels to data
Feature learning (with labels)

**Generative Model:**
Learn a probability distribution p(x)

⟶ Detect outliers
Feature learning (without labels)
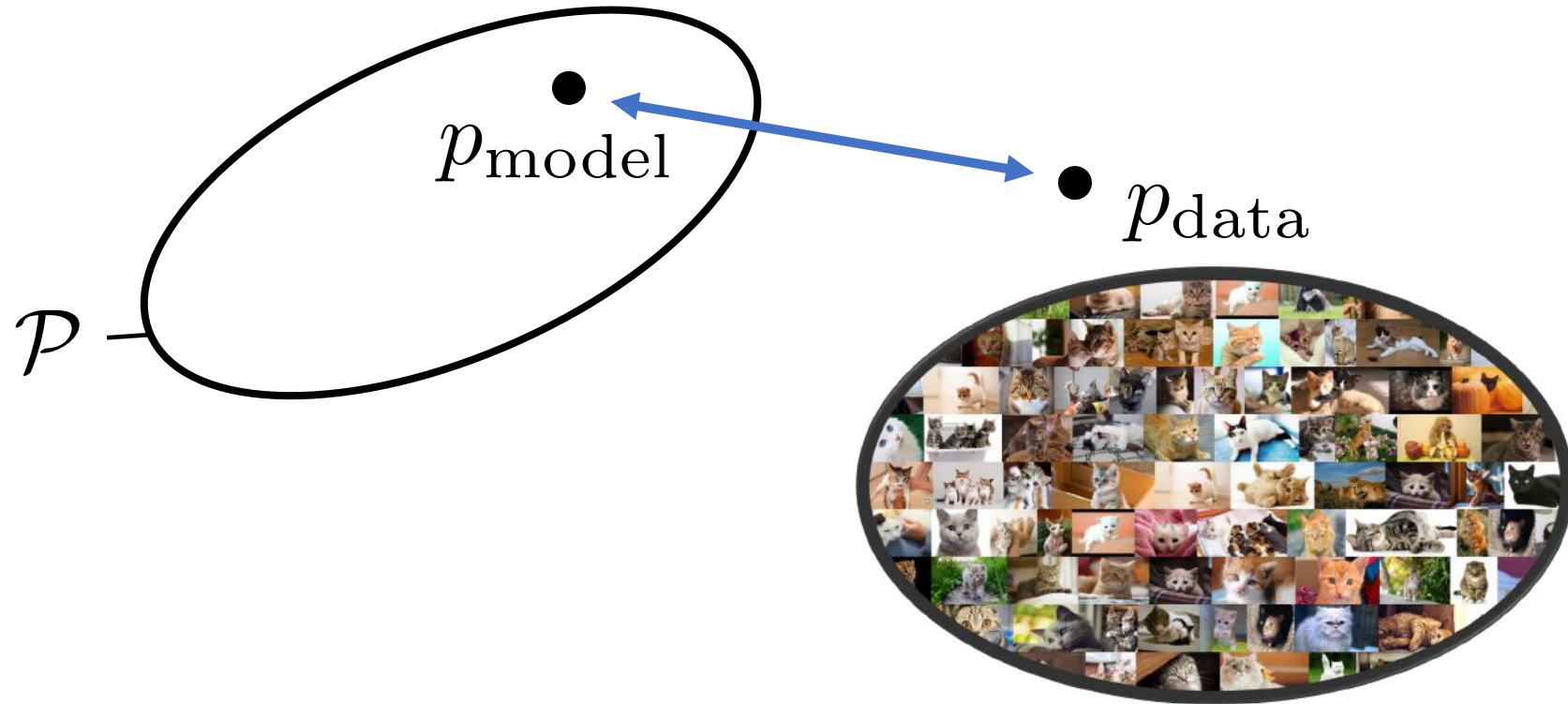Sample to **generate** new data

**Conditional Generative Model:**
Learn p(x|y)

# What can we do with a discriminative model?

**Discriminative Model:**
Learn a probability distribution $p(y|x)$

⟶

Assign labels to data
Feature learning (with labels)

**Generative Model:**
Learn a probability distribution $p(x)$

⟶

Detect outliers
Feature learning (without labels)
Sample to **generate** new data

**Conditional Generative Model:**
Learn $p(x|y)$

⟶

Assign labels, while rejecting outliers!
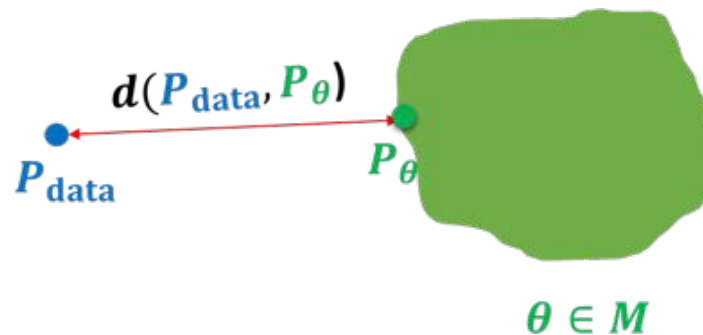Generate new data conditioned on input labels

# Generative Modeling



$p_{\mathrm{model}}$

$p_{\mathrm{data}}$

$\mathcal{P}$

- **Goal:** Learn some underlying hidden structure of the training samples to generate novel samples from same data distribution

# Learning a generative model

- We are given a training set of examples, e.g., images of dogs
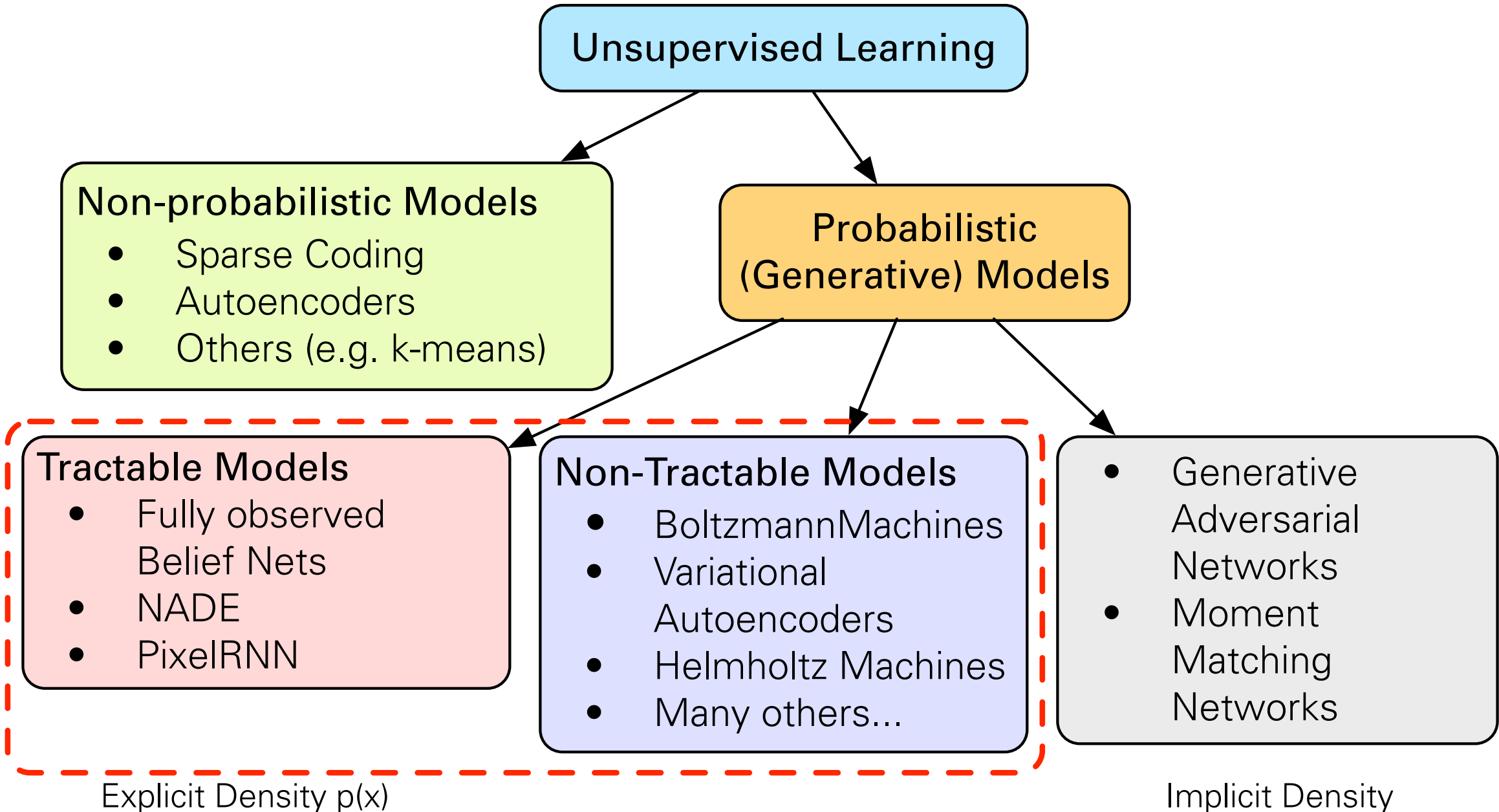


$$x_i \sim P_{data}$$
$$i = 1, 2, \ldots, n$$

$$d(P_{data}, P_\theta)$$

$P_{data}$

$P_\theta$

$\theta \in M$

Model family

- We want to learn a probability distribution p(x) over images x s.t.
  - **Generation**: If we sample $x_{new} \sim$ p(x), $x_{new}$ should look like a dog (sampling)
  - **Density estimation:** p(x) should be high if x looks like a dog, and low otherwise (anomaly detection)
  - **Unsupervised representation learning**: We should be able to learn what these images have in common, e.g., ears, tail, etc. (features)

# Why Unsupervised Learning?

- Given high-dimensional data $X = (x_1, \ldots, x_n)$ we want to find a low-dimensional model characterizing the population.

- Recent progress mostly in supervised DL

- Real challenges for unsupervised DL

- Potential benefits:
    - **Exploit tons of unlabeled data**
    - Answer new questions about the variables observed
    - Regularizer – transfer learning – domain adaptation
    - Easier optimization (divide and conquer)
    - Joint (structured) outputs

Unsupervised Learning

Non-probabilistic Models
- Sparse Coding
- Autoencoders
- Others (e.g. k-means)

Probabilistic (Generative) Models

Tractable Models
- Fully observed Belief Nets
- NADE
- PixelRNN

Non-Tractable Models
- BoltzmannMachines
- Variational Autoencoders
- Helmholtz Machines
- Many others...

- Generative Adversarial Networks
- Moment Matching Networks

Explicit Density p(x)

Implicit Density

# Unsupervised Learning

- Basic Building Blocks:
    - Sparse Coding
    - Autoencoders

- Autoregressive Generative Models

- Generative Adversarial Networks

- Variational Autoencoders

- Normalizing Flow Models

# Sparse Coding

- Sparse coding (Olshausen & Field, 1996). Originally developed to explain early visual processing in the brain (edge detection).

- **Objective:** Given a set of input data vectors $\{\mathbf{x}_1, \mathbf{x}_2, ..., \mathbf{x}_N\}$, learn a dictionary of bases, such that:

$$\mathbf{x}_n = \sum_{k=1}^{K} a_{nk} \phi_k$$

Sparse: mostly zeros

- Each data vector is represented as a sparse linear combination of bases.

# Sparse Coding

Natural Images



Learned bases: "Edges"



New example



$x \qquad = 0.8 * \qquad \phi_{36} \qquad + 0.3 * \qquad \phi_{42} \qquad + 0.5 * \qquad \phi_{65}$

[0.0, 0.0, ... **0.8**, ..., **0.3**, ..., **0.5**, ...] = coefficients (feature representation)

# Sparse Coding: Training

- Input image patches: $\mathbf{x}_1, \mathbf{x}_2, \ldots, \mathbf{x}_N \in \mathbb{R}^D$
- Learn dictionary of bases: $\phi_1, \phi_2, \ldots, \phi_K \in \mathbb{R}^D$

$$\min_{\mathbf{a}, \phi} \sum_{n=1}^{N} \left\| \mathbf{x}_n - \sum_{k=1}^{K} a_{nk} \phi_k \right\|_2^2 + \lambda \sum_{n=1}^{N} \sum_{k=1}^{K} |a_{nk}|$$

Reconstruction error    Sparsity penalty

- Alternating Optimization:
  1. Fix dictionary of bases and solve for activations **a** (a standard Lasso problem).
  2. Fix activations **a**, optimize the dictionary of bases (convex QP problem).

# Sparse Coding: Testing Time

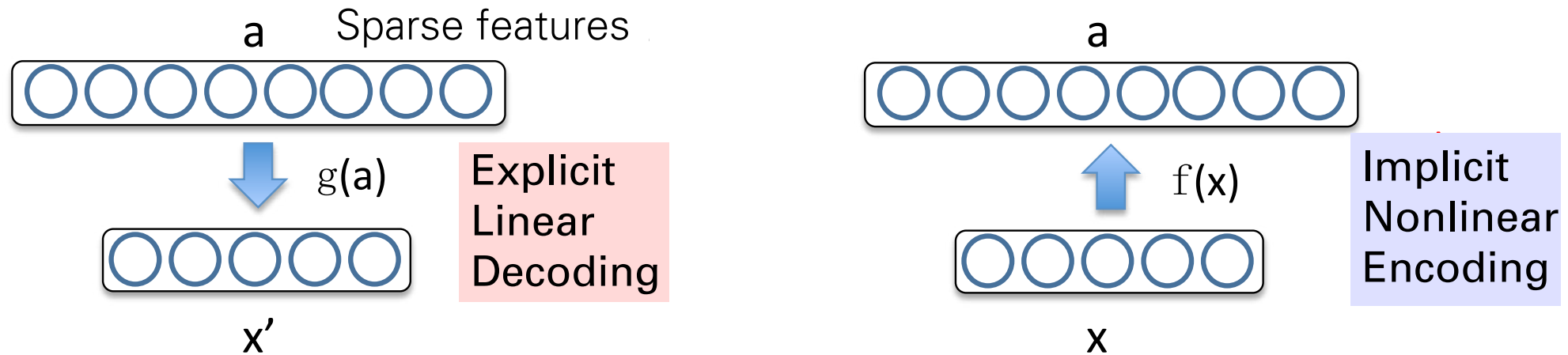- **Input:** a new image patch x* , and K learned bases $\phi_1, \phi_2, ..., \phi_K, ..., \phi_K$
- **Output:** sparse representation **a** of an image patch x*.

$$\min_{\mathbf{a}} \left\| \mathbf{x}^* - \sum_{k=1}^{K} a_k \phi_k \right\|_2^2 + \lambda \sum_{k=1}^{K} |a_k| \sum_{:=1}^{K} |a_k|$$



$\phi_{36}$    $\phi_{36}$    $\phi_{42}$    $\phi_{42}$    $\phi_{65}$    $\phi_{65}$

# Sparse Coding: Testing Time

- **Input:** a new image patch x* , and K learned bases $\phi_1, \phi_2, ..., \phi_K, ..., \phi_K$
- **Output:** sparse representation **a** of an image patch x*.

$$\min_{\mathbf{a}} \left\| \mathbf{x}^* - \sum_{k=1}^{K} a_k \phi_k \right\|_2^2 + \lambda \sum_{k=1}^{K} |a_k| \sum_{:=1}^{K} |a_k|$$



$$x^* = \phi_{36} \cdot \phi_{36} + \phi_{42} \cdot \phi_{42} + \phi_{65} \cdot \phi_{65}$$

[0.0, 0.0, ... **0.8**, ..., **0.3**, ..., **0.5**, ...] = coefficients (feature representation)

# Image Classification

- Evaluated on Caltech101 object category dataset.



Input Image    Learned bases    Features (coefficients)    Classification Algorithm (SVM)

**9K images, 101 classes**

| Algorithm | Accuracy |
|---|---|
| Baseline (Fei-Fei et al., 2004) | 16% |
| PCA | 37% |
| **Sparse Coding** | **47%** |

(Lee, Battle, Raina, Ng, NIPS 2007)

# Modeling Image Patches

- Natural image patches:
  - small **image regions** extracted from an image of nature (forest, grass, …)



Image

# Relationship to V1

- When trained on natural image patches
  - the dictionary columns (''atoms'') look like **edge detectors**

  -  each atom is tuned to a particular **position**, **orientation** and **spatial frequency**

  - V1 neurons in the mammalian brain have a similar behavior



Emergence of simple-cell receptive field properties by learning a sparse code of natural images. Olshausen and Field, 1996

# Relationship to V1

- Suggests that the brain might be learning a sparse code of visual stimulus

  - Since then, many other models have been shown to learn similar features

  - they usually all incorporate a notion of sparsity



Emergence of simple-cell receptive field properties by learning a sparse code of natural images. Olshausen and Field, 1996

# Interpreting Sparse Coding

$$\min_{\mathbf{a},\phi} \sum_{n=1}^{N} \left\| \mathbf{x}_n - \sum_{k=1}^{K} a_{nk}\phi_k \right\|_2^2 + \lambda \sum_{n=1}^{N} \sum_{k=1}^{K} |a_{nk}|$$

a    Sparse features

g(a)

Explicit
Linear
Decoding

x'

# Interpreting Sparse Coding

$$\min_{\mathbf{a},\phi} \sum_{n=1}^{N} \left\| \mathbf{x}_n - \sum_{k=1}^{K} a_{nk} \phi_k \right\|_2^2 + \lambda \sum_{n=1}^{N} \sum_{k=1}^{K} |a_{nk}|$$

a          Sparse features

a

g(a)

Explicit Linear Decoding

f(x)

Implicit Nonlinear Encoding

x'

x

- Sparse, over-complete representation **a.**
- **Encoding** a = f(**x**) is implicit and nonlinear function of **x**.
- **Reconstruction** (or decoding) **x'** = g(**a**) is linear and explicit.

# Autoencoder

# Autoencoder

| Feature Representation |
|:---:|

Feed-back,
generative,
top-down

| Decoder |
|:---:|

| Encoder |
|:---:|

Feed-forward,
bottom-up

| Input Image |
|:---:|

- Details of what goes insider the encoder and decoder matter!
- Need constraints to avoid learning an identity.

# Autoencoder



Binary Features z

Decoder Filters D

Linear function

$$Dz$$

$$z = \sigma(Wx)$$

Encoder filters W

Sigmoid function

$$\sigma(x) = \frac{1}{1 + \exp(-x)}$$

Input Image x

# Autoencoder

**Binary Features z**

$$\sigma(W^Tz)$$  Decoder Filters D

$$z=\sigma(Wx)$$

Encoder filters W

Sigmoid function

$$\sigma(x) = \frac{1}{1 + \exp(-x)}$$

**Input Image x**

- Need additional constraints to avoid learning an identity.
- Relates to Restricted Boltzmann Machines (later).

# Autoencoder

- Feed-forward neural network trained to reproduce its input at the output layer



**Decoder**

$$\widehat{\mathbf{x}} = o(\widehat{\mathbf{a}}(\mathbf{x}))$$

$$= \mathrm{sigm}(\mathbf{c} + \mathbf{W}^*\mathbf{h}(\mathbf{x}))$$

for binary units

**Encoder**

$$\mathbf{h}(\mathbf{x}) = g(\mathbf{a}(\mathbf{x}))$$

$$= \mathrm{sigm}(\mathbf{b} + \mathbf{W}\mathbf{x})$$

# Loss Function

$\widehat{\mathbf{x}} = \sigma(\mathbf{c} + \mathbf{W}^*\mathbf{h}(\mathbf{x}))$

$\mathbf{h}(\mathbf{x}) = g(\mathbf{a}(\mathbf{x}))$

- **Loss function for binary inputs**

$= \text{sigm}(\mathbf{c} + \mathbf{W}^*\mathbf{h}(\mathbf{x}))$

$= \text{sigm}(\mathbf{b} + \mathbf{W}\mathbf{x})$

$l(f(\mathbf{x})) = -\sum_k \left( x_k \log(\widehat{x}_k) + (1 - x_k)\log(1 - \widehat{x}_k) \right)$

 – Cross-entropy error function (reconstruction loss)

$l(f(\mathbf{x})) = \sum_k (\widehat{x}_k - x_k)^2$

$\widehat{\mathbf{x}} = \sigma(\mathbf{c} + \mathbf{W}^*\mathbf{h}(\mathbf{x}))$

$= \text{sigm}(\mathbf{c} + \mathbf{W}^*\mathbf{h}(\mathbf{x}))$

- **Loss function for real-valued inputs**

$l(f(\mathbf{x})) = \frac{1}{2}\sum_k (\widehat{x}_k - x_k)^2$

$l(f(\mathbf{x})) = -\sum_k \left( x_k \log(\widehat{x}_k) + (1 - x_k)\log(1 - \widehat{x}_k) \right)$

 – onstruction loss)

 – linear activation function at the output

$\nabla_{\mathbf{a}(\mathbf{x}^{(t)})} l(f(\mathbf{x}^{(t)})) = \widehat{\mathbf{x}}^{(t)} - \mathbf{x}^{(t)}$

$\mathbf{a}(\mathbf{x}^{(t)}) \Leftarrow \mathbf{b} + \mathbf{W}\mathbf{x}^{(t)}$

$\mathbf{h}(\mathbf{x}^{(t)}) \Leftarrow \text{sigm}(\mathbf{a}(\mathbf{x}^{(t)}))$

# Autoencoder

**Linear Features z**

↓

**Wz**

↓

**Input Image x**

↑

**z=Wx**

↑

- If the **hidden and output layers are linear**, it will learn hidden units that are a linear function of the data and minimize the squared error.

- The K hidden units will span the same space as the first k principal components. The weight vectors may not be orthogonal.

- With nonlinear hidden units, we have a nonlinear generalization of PCA.

# Denoising Autoencoder

- **Idea:** Representation should be robust to introduction of noise:
  - random assignment of subset of inputs to 0, with probability $\nu$
  - Gaussian additive noise

- $\widehat{\mathbf{x}}$ Similar to dropouts on the input layer

- Reconstruction $\widehat{\mathbf{x}}$ computed from the corrupted input $\widetilde{\mathbf{x}}$

- **Loss function** compares $\widehat{\mathbf{x}}$ reconstruction with the noiseless input $\mathbf{x}$

# Denoising Autoencoder

$$\widehat{\mathbf{x}} = \mathrm{sigm}(\mathbf{c} + \mathbf{W}^* \mathbf{h}(\widetilde{\mathbf{x}}))$$

$$\widetilde{\mathbf{x}}$$

$$\mathbf{x}$$

# Learned Filters

Non-corrupted                                          25% corrupted input

# Learned Filters

Non-corrupted

50% corrupted input



(a) No destroyed inputs

(b) 25% destruction

# Predictive Sparse Decomposition



Binary Features z

L$_1$ Sparsity

Decoder filters D

Dz

z=σ(Wx)

Real-valued Input x

Encoder filters W

Sigmoid function

$$\min_{D,W,\mathbf{z}} ||D\mathbf{z} - \mathbf{x}||_2^2 + \lambda|\mathbf{z}|_1 + ||\sigma(W\mathbf{x}) - \mathbf{z}||_2^2$$

(Kavukcuoglu, Ranzato, Fergus, LeCun, 2009)

# Predictive Sparse Decomposition

Binary Features z

L₁ Sparsity

Dz

z=σ(Wx)

Encoder filters W

Decoder filters D

Sigmoid function

Real-valued Input x

At training time

$$\min_{D,W,\mathbf{z}} ||D\mathbf{z} - \mathbf{x}||_2^2 + \lambda|\mathbf{z}|_1 + ||\sigma(W\mathbf{x}) - \mathbf{z}||_2^2$$

Decoder

Encoder

(Kavukcuoglu, Ranzato, Fergus, LeCun, 2009)

# Stacked Autoencoders



Features

Sparsity

Decoder

Encoder

Input x

# Stacked Autoencoders

# Stacked Autoencoders

| | | |
|---|---|---|
| | **Class Labels** | |
| **Decoder** | | **Encoder** |
| | **Features** | |
| **Sparsity** | **Decoder** | **Encoder** |
| | **Features** | |
| **Sparsity** | **Decoder** | **Encoder** |
| | **Input x** | |

# Stacked Autoencoders

| Class Labels |
| :---: |

Decoder     Encoder

| Features |
| :---: |

Sparsity     Decoder     Encoder

Features

## Greedy Layer-wise Learning

| Input x |
| :---: |

# Stacked Autoencoders

- Remove decoders and use feed-forward part.

Class Labels

Encoder

Features

Encoder

Features

Encoder

Input x

# Stacked Autoencoders

- Remove decoders and use feed-forward part.

- Standard, or convolutional neural network architecture.

- Parameters can be fine-tuned using backpropagation.

Class Labels

Encoder

Features

Encoder

Features

Encoder

Input x

# Deep Autoencoder

**30**

$W_4$

**500**

Top RBM

**500**

$W_3$

**1000**

RBM

**1000**

$W_2$

**2000**

RBM

**2000**

$W_1$

RBM

**Pretraining**

Decoder

$W_1^T$

**2000**

$W_2^T$

**1000**

$W_3^T$

**500**

$W_4^T$

**30** Code layer

$W_4$

**500**

$W_3$

**1000**

$W_2$

**2000**

$W_1$

Encoder

**Unrolling**

$W_1^T + \varepsilon_8$

**2000**

$W_2^T + \varepsilon_7$

**1000**

$W_3^T + \varepsilon_6$

**500**

$W_4^T + \varepsilon_5$

**30**

$W_4 + \varepsilon_4$

**500**

$W_3 + \varepsilon_3$

**1000**

$W_2 + \varepsilon_2$

**2000**

$W_1 + \varepsilon_1$

**Fine−tuning**

# Deep Autoencoders

- 25x25 – 2000 – 1000 – 500 – 30 autoencoder to extract 30-D real-valued codes for Oliver face patches.



- **Top:** Random samples from the test dataset.
- **Middle:** Reconstructions by the 30-dimensional deep autoencoder.
- **Bottom:** Reconstructions by the 30-dimensional PCA.

(Hinton and Salakhutdinov, Science 2006)

# Information Retrieval



Interbank Markets · European Community Monetary/Economic · Energy Markets · Disasters and Accidents · Leading Economic Indicators · Legal/Judicial · Accounts/Earnings · Government Borrowings

2-D LSA space

- The Reuters Corpus Volume II contains 804,414 newswire stories (randomly split into **402,207 training** and **402,207 test**).

- "Bag-of-words" representation: each article is represented as a vector containing the counts of the most frequently used 2000 words in the training set.

# Semantic Hashing



- Learn to map documents into semantic 20-D binary codes.
- Retrieve similar documents stored at the nearby addresses with no search at all.

(Hinton and Salakhutdinov, Science 2006)

# Searching Large Image Database using Binary Codes

- Map images into binary codes for fast retrieval.



| Input image | 30–RBM | 64–RBM | 128–RBM | 256–RBM |

- Small Codes, Torralba, Fergus
- Spectral Hashing, Y. Weiss, A. Torralba, R. Fergus, NIPS 2008
- Kulis and Darrell, NIPS 2009, Gong and Lazebnik, CVPR 2011
- Norouzi and Fleet, ICML 2011

Unsupervised Learning

Non-probabilistic Models
- Sparse Coding
- Autoencoders
- Others (e.g. k-means)

Probabilistic (Generative) Models

Tractable Models
- Fully observed Belief Nets
- NADE
- PixelRNN

Non-Tractable Models
- BoltzmannMachines
- Variational Autoencoders
- Helmholtz Machines
- Many others...

- Generative Adversarial Networks
- Moment Matching Networks

Explicit Density p(x)

Implicit Density

# Generative Adversarial Networks

# Genetive Adversarial Networks (GANs)

(Goodfellow et al., 2014)

Noise
(random input)

Generative
Model

$z \sim \text{Uniform}_{100}$

think of this as
a transformation

- A game-theoretic likelihood free model

## Advantages:

- Uses a latent code

- No Markov chains needed

- Produces the best looking samples

# Genetive Adversarial Networks (GANs)



(Goodfellow et al., 2014)

$$\{x_1, \dots, x_n\} \sim p_{\text{data}}$$

- A game between a generator $G_\theta(z)$ and a discriminator $D_\omega(x)$
  - Generator tries to fool discriminator (i.e. generate realistic samples)
  - Discriminator tries to distinguish fake from real samples

# Intuition behind GANs



$D_{\omega}$ : Discriminator (Art Critic)

$\boldsymbol{x}_{real}$

$\boldsymbol{x}_{fake}$

$G_{\theta}$ : Generator (Forger)

# Training Procedure

- Use SGD on two minibatches simultaneously:
  - A minibatch of training examples
  - A minibatch of generated samples

# GAN Training: Minimax Game (Goodfellow et al., 2014)

$$\min_\theta \max_\omega \mathbb{E}_{\boldsymbol{x} \sim p_{\text{data}}} \left[\log D_\omega(\boldsymbol{x})\right] + \mathbb{E}_{\boldsymbol{z} \sim p_{\boldsymbol{z}}} \left[\log \left(1 - D_\omega(G_\theta(\boldsymbol{z}))\right)\right]$$

Real data

Noise vector used
to generate data

$$J^{(D)} = -\frac{1}{2}\mathbb{E}_{\boldsymbol{x} \sim p_{\text{data}}} \log D(\boldsymbol{x}) - \frac{1}{2}\mathbb{E}_{\boldsymbol{z}} \log \left(1 - D\left(G(\boldsymbol{z})\right)\right)$$

Cross-entropy loss for binary classification

$$J^{(G)} = -\frac{1}{2}\mathbb{E}_{\boldsymbol{z}} \log D\left(G(\boldsymbol{z})\right)$$

Generator maximizes the log-probability of the discriminator being mistaken

- Equilibrium of the game
- **Minimizes the Jensen-Shannon divergence between $p_{\text{data}}$ and $p_x$**

# GAN Training: Minimax Game (Goodfellow et al., 2014)

$$\min_{\theta} \max_{\omega} \mathbb{E}_{\boldsymbol{x} \sim p_{\text{data}}} \left[ \log D_{\omega}(\boldsymbol{x}) \right] + \mathbb{E}_{\boldsymbol{z} \sim p_{\boldsymbol{z}}} \left[ \log \left( 1 - D_{\omega}(G_{\theta}(\boldsymbol{z})) \right) \right]$$

Real data

Noise vector used

$$J^{(D)} = -\frac{1}{2}\mathbb{E}_{\boldsymbol{x} \sim}$$

$$J^{(G)} = -\frac{1}{2}\mathbb{E}_{\boldsymbol{z}} \log$$

Important question is
"Does this converge??"

ross-entropy
ss for binary
assification

g-probability
of the discriminator being mistaken

- Equilibrium of the game
- Minimizes the Jensen-Shannon divergence

# Training Procedure

(Goodfellow et al., 2014)



Source: Alec Radford



Source: OpenAI blog

Generating 1D points

Generating images

# Training Procedure

(Goodfellow et al., 2014)

- Use SGD on two minibatches simultaneously:
  - A minibatch of training examples
  - A minibatch of generated samples

# Training Procedure

- Updating the discriminator:



$$\{\boldsymbol{x}_1, \ldots, \boldsymbol{x}_n\} \sim p_{\text{data}}$$

update the discriminator weights using backprop on the classification objective

# Training Procedure

- Updating the generator:



update the generator weights using backprop

flip the sign of the derivatives

backprop the derivatives, but don't modify the discriminator weights

# Results

(Goodfellow et al., 2014)

- The generator uses a mixture of rectifier linear activations and/or sigmoid activations

- The discriminator net used maxout activations.



MNIST samples



TFD samples



CIFAR10 samples
(fully-connected model)



CIFAR10 samples
(convolutional discriminator, deconvolutional generator)

# Deep Convolutional GANs (DCGAN)

- **Idea:** Tricks to make GAN training more stable



- No fully connected layers

- Batch Normalization
  (Ioffe and Szegedy, 2015)

- Leaky Rectifier in D

- Use Adam (Kingma and Ba, 2015)

- Tweak Adam hyperparameters a bit
  (lr=0.0002, b1=0.5)

# DCGAN for LSUN Bedrooms

64×64 pixels
~3M images

(Radford et al., 2015)

# Walking over the latent space

(Radford et al., 2015)

- Interpolation suggests non-overfitting behavior

# Walking over the latent space    (Radford et al., 2015)

# Vector Space Arithmetic

(Radford et al., 2015)



man
with glasses — man
without glasses + woman
without glasses = woman with glasses

# Vector Space Arithmetic

smiling woman $-$ neutral woman $+$ neutral man $=$ smiling man

# Cartoon of the Image manifold

# What makes GANs special?



more traditional max-likelihood approach

GAN

# GAN Failures: Mode Collapse

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

- $D$ in inner loop: convergence to correct distribution
- $G$ in inner loop: place all mass on most likely point



Target

Step 0    Step 5k    Step 10k    Step 15k    Step 20k    Step 25k

(Metz et al., 2016)

92

# Mode Collapse: Solutions

- **Unrolled GANs** (Metz et al 2016): Prevents mode collapse by backproping through a set of (k) updates of the discriminator to update generator parameters



Step 0    Step 5k    Step 10k    Step 15k    Step 20k    Step 25k

- **VEEGAN** (Srivastava et al 2017): Introduce a reconstructor network which is learned both to map the true data distribution $p(x)$ to a Gaussian and to approximately invert the generator network.

# Mode Collapse: Solutions

- **Minibatch Discrimination** (Salimans et al 2016): Add minibatch features that classify each example by comparing it to other members of the minibatch (Salimans et al 2016)

- **PacGAN:** The power of two samples in generative adversarial networks (Lin et al 2017): Also uses multisample discrimination.



**GAN Discriminator**

Input Layer

**PacGAN2 Discriminator**

Input Layer

# Mode Collapse: Solutions

- **PacGAN:** The power of two samples in generative adversarial networks (Lin et al 2017): Also uses multisample discrimination.



Figure 2: Scatter plot of the 2D samples from the true distribution (left) of 2D-grid and the learned generators using GAN (middle) and PacGAN2 (right). PacGAN2 captures all of the 25 modes.

95

# GAN Evaluation

- Quantitatively evaluating GANs is not straightforward:
  - Max Likelihood is a poor indication of sample quality

- Some evaluation metrics

  - **Inception Score (IS)**:
    $y$ = labels given gen. image. p(y|x) is from classifier - InceptionNet

    $$\text{IS}(\mathbb{P}_g) = e^{\mathbb{E}_{\mathbf{x} \sim \mathbb{P}_g}[KL(p_{\mathcal{M}}(y|\mathbf{x})||p_{\mathcal{M}}(y))]}$$

  - **Fréchet inception distance (FID):** (Currently most popular)
    Estimate mean $m$ and covariance $C$ from classifier output - InceptionNet

    $$d^2((\boldsymbol{m}, \boldsymbol{C}), (\boldsymbol{m}_w, \boldsymbol{C}_w)) = \|\boldsymbol{m} - \boldsymbol{m}_w\|_2^2 + \text{Tr}\big(\boldsymbol{C} + \boldsymbol{C}_w - 2\big(\boldsymbol{C}\boldsymbol{C}_w\big)^{1/2}\big)$$

  - **Kernel MMD** (Maximum Mean Discrepancy):

    $$\text{MMD}(\mathbb{P}_r, \mathbb{P}_g) = \left( \mathbb{E}_{\substack{\mathbf{x}_r, \mathbf{x}_r' \sim \mathbb{P}_r, \\ \mathbf{x}_g, \mathbf{x}_g' \sim \mathbb{P}_g}} \Big[ k(\mathbf{x}_r, \mathbf{x}_r') - 2k(\mathbf{x}_r, \mathbf{x}_g) + k(\mathbf{x}_g, \mathbf{x}_g') \Big] \right)^{\frac{1}{2}}$$

# Subclasses of GANs



**Vanilla GAN**

**Vanilla GAN**
(Goodfellow, et al., 2014)

**Discriminator Looks at Latent Variables**

**Conditional GAN**
(Mirza & Osindero, 2014)

**Bidirectional GAN**
(Donahue, et al., 2016;  Dumoulin, et al., 2016)

**Discriminator Predicts Latent Variables**

**Semi-Supervised GAN**
(Odena, 2016;  Salimans, et al., 2016)

**InfoGAN**
(Chen, et al., 2016)

**Auxiliary Classifier GAN**
(Odena, et al., 2016)

97

# Vanilla GAN (Goodfellow et al., 2014)



DCGAN (Radford et al., 2015)

# Conditional GAN (Mirza and Osindero, 2014)



- Add conditional variables $\boldsymbol{y}$ into $G$ and $D$

$$\min_G \max_D V(D, G) = \mathbb{E}_{\boldsymbol{x} \sim p_{\text{data}}(\boldsymbol{x})}[\log D(\boldsymbol{x}|\boldsymbol{y})] + \mathbb{E}_{\boldsymbol{z} \sim p_z(\boldsymbol{z})}[\log(1 - D(G(\boldsymbol{z}|\boldsymbol{y})))]$$

# Auxiliary Classifier GAN (Odena et al., 2016)



- Every generated sample has a corresponding class label

$$L_S = E[\log P(S = real \mid X_{real})] + E[\log P(S = fake \mid X_{fake})]$$
$$L_C = E[\log P(C = c \mid X_{real})] + E[\log P(C = c \mid X_{fake})]$$

- $D$ is trained to maximize $L_S + L_C$
- $G$ is trained to maximize $L_C - L_S$

- Learns a representation for $z$ that is independent of class label

# Auxiliary Classifier GAN (Odena et al., 2016)

128×128 resolution samples from 5 classes taken from an AC-GAN trained on the ImageNet



monarch butterfly       goldfinch       daisy       redshank       grey whale

# Bidirectional GAN (Donahue et al., 2016; Dumoulin et al., 2016)



- Jointly learns a generator network and an inference network using an adversarial process.

$$\min_{G} \max_{D} V(D, G) = \mathbb{E}_{q(\boldsymbol{x})}[\log(D(\boldsymbol{x}, G_z(\boldsymbol{x})))] + \mathbb{E}_{p(\boldsymbol{z})}[\log(1 - D(G_x(\boldsymbol{z}), \boldsymbol{z}))]$$

$$= \iint q(\boldsymbol{x})q(\boldsymbol{z} \mid \boldsymbol{x})\log(D(\boldsymbol{x}, \boldsymbol{z}))d\boldsymbol{x}d\boldsymbol{z}$$

$$+ \iint p(\boldsymbol{z})p(\boldsymbol{x} \mid \boldsymbol{z})\log(1 - D(\boldsymbol{x}, \boldsymbol{z}))d\boldsymbol{x}d\boldsymbol{z}.$$



CelebA reconstructions



SVNH reconstructions

# Bidirectional GAN (Donahue et al., 2016; Dumoulin et al., 2016)

LSUN bedrooms

Tiny ImageNet

# Wasserstein GAN (Arjovsky et al., 2016)

- Objective based on Earth-Mover or Wassertein distance:

$$\min_{\theta} \max_{\omega} \mathbb{E}_{\boldsymbol{x} \sim p_{\mathrm{data}}} \left[ D_\omega(\boldsymbol{x}) \right] - \mathbb{E}_{\boldsymbol{z} \sim p_{\boldsymbol{z}}} \left[ D_\omega(G_\theta(\boldsymbol{z})) \right]$$

- Provides nice gradients over real and fake samples

# Wasserstein GAN (Arjovsky et al., 2016)

- Wasserstein loss seems to correlate well with image quality.

# WGAN with gradient penalty (Gulraani et al., 2017)

$$L = \underbrace{\mathop{\mathbb{E}}_{\tilde{\boldsymbol{x}} \sim \mathbb{P}_g} [D(\tilde{\boldsymbol{x}})] - \mathop{\mathbb{E}}_{\boldsymbol{x} \sim \mathbb{P}_r} [D(\boldsymbol{x})]}_{\text{Original critic loss}} + \underbrace{\lambda \mathop{\mathbb{E}}_{\hat{\boldsymbol{x}} \sim \mathbb{P}_{\hat{\boldsymbol{x}}}} \left[ (\|\nabla_{\hat{\boldsymbol{x}}} D(\hat{\boldsymbol{x}})\|_2 - 1)^2 \right]}_{\text{Our gradient penalty}}$$

- Faster convergence and higher-quality samples than WGAN with weight clipping

- Train a wide variety of GAN architectures with almost no hyperparameter tuning, including discrete models

Samples from a character-level GAN language model on Google Billion Word

**WGAN with gradient penalty**

| |
|---|
| Busino game camperate spent odea |
| In the bankaway of smarling the |
| SingersMay , who kill that imvic |
| Keray Pents of the same Reagun D |
| Manging include a tudancs shat " |
| His Zuith Dudget , the Denmbern |
| In during the Uitational questio |
| Divos from The ' noth ronkies of |
| She like Monday , of macunsuer S |
| The investor used ty the present |
| A papees are cointry congress oo |
| A few year inom the group that s |
| He said this syenn said they wan |
| As a world 1 88 ,for Autouries |
| Foand , th Word people car , Il |
| High of the upseader homing pull |
| The guipe is worly move dogsfor |
| The 1874 incidested he could be |
| The allo tooks to security and c |

| |
|---|
| Solice Norkedin pring in since |
| ThiS record ( 31. ) UBS ) and Ch |
| It was not the annuas were plogr |
| This will be us , the ect of DAN |
| These leaded as most-worsd p2 a0 |
| The time I paidOa South Cubry i |
| Dour Fraps higs it was these del |
| This year out howneed allowed lo |
| Kaulna Seto consficutes to repor |
| A can teal , he was schoon news |
| In th 200. Pesish picriers rega |
| Konney Panice rimimber the teami |
| The new centuct cut Denester of |
| The near , had been one injostie |
| The incestion to week to shorted |
| The company the high product of |
| 20 - The time of accomplete , wh |
| John WVuderenson seqiivic spends |
| A ceetens in indestredly the Wat |

**Standard GAN objective**

ddddddddddddddddddddddddddddddddd
ddddddddddddddddddddddddddddddddd

ddddddddddddddddddddddddddddddddd
ddddddddddddddddddddddddddddddddd

# Least Squares GAN (LSGAN) (Mao et al., 2017)

- Use a loss function that provides smooth and non-saturating gradient in discriminator D

$$\min_D V_{\text{LSGAN}}(D) = \frac{1}{2}\mathbb{E}_{\boldsymbol{x}\sim p_{\text{data}}(\boldsymbol{x})}\big[(D(\boldsymbol{x}) - b)^2\big] + \frac{1}{2}\mathbb{E}_{\boldsymbol{z}\sim p_{\boldsymbol{z}}(\boldsymbol{z})}\big[(D(G(\boldsymbol{z})) - a)^2\big]$$

$$\min_G V_{\text{LSGAN}}(G) = \frac{1}{2}\mathbb{E}_{\boldsymbol{z}\sim p_{\boldsymbol{z}}(\boldsymbol{z})}\big[(D(G(\boldsymbol{z})) - c)^2\big],$$



Decision boundaries of Sigmoid & Least Squares loss functions

Sigmoid decision boundary

Least Squares decision boundary

# Least Squares GAN (LSGAN) (Mao et al., 2017)



Church

Kitchen

# Boundary Equilibrium GAN (BEGAN)

(Berthelot et al., 2017)

- A loss derived from the Wasserstein distance for training auto-encoder based GANs

$$\mathcal{L}(v) = |v - D(v)|^\eta \text{ where } \begin{cases} D : \mathbb{R}^{N_x} \mapsto \mathbb{R}^{N_x} & \text{is the autoencoder function.} \\ \eta \in \{1, 2\} & \text{is the target norm.} \\ v \in \mathbb{R}^{N_x} & \text{is a sample of dimension } N_x. \end{cases}$$

  - Wasserstein distance btw. the reconstruction losses of real and generated data

- Convergence measure:

$$\mathcal{M}_{global} = \mathcal{L}(x) + |\gamma \mathcal{L}(x) - \mathcal{L}(G(z_G))|$$

- Objective:

$$\begin{cases} \mathcal{L}_D = \mathcal{L}(x) - k_t.\mathcal{L}(G(z_D)) & \text{for } \theta_D \\ \mathcal{L}_G = \mathcal{L}(G(z_G)) & \text{for } \theta_G \\ k_{t+1} = k_t + \lambda_k(\gamma\mathcal{L}(x) - \mathcal{L}(G(z_G))) & \text{for each training} \\ & \text{step } t \end{cases}$$



(a) Generator/Decoder        (b) Encoder



109

# BEGANs for CelebA

360K celebrity face images
128x128 with 128 filters

(Berthelot et al., 2017)



Interpolations in the latent space



Mirror interpolation example

# Progressive GANs (Karras et al., 2018)

- Progressively generate high-res images

- Multi-step training from low to high resolutions

# Progressive GANs (Karras et al., 2018)



- Training process

4x4

Progressive GANs (Karras et al., 2018)

CelebA-HQ
random interpolations

## High resolution, class-conditional samples generated by the model



- BigGANs trained with 2-4x as many parameters and 8x the batch size compared to prior art.
- Uses Gaussian truncation to sample z (avoid sampling from the tail of the Gaussian distribution)
- Uses multiple other tricks including multiple regularizations including a Gradient penalty regularization and an Orthogonal Regularization.

# BigGANs (Brock et al., 2019)



Easy classes

Hard classes

Resolution: 512x512

115

# StyleGAN (Karras et al., 2019)



Latent $\mathbf{z} \in \mathcal{Z}$

Normalize

Fully-connected

PixelNorm

Conv 3×3

PixelNorm

4×4

Upsample

Conv 3×3

PixelNorm

Conv 3×3

PixelNorm

8×8

...

(a) Traditional

Latent $\mathbf{z} \in \mathcal{Z}$

Normalize

Mapping network $f$

FC

FC

FC

FC

FC

FC

FC

FC

$\mathbf{w} \in \mathcal{W}$

Synthesis network $g$

Noise

Const 4×4×512

B

A — style → AdaIN

Conv 3×3

B

A — style → AdaIN

4×4

Upsample

Conv 3×3

B

A — style → AdaIN

Conv 3×3

B

A — style → AdaIN

8×8

...

(b) Style-based generator

Feature map affine transformation:

$$\mathrm{AdaIN}(\mathbf{x}_i, \mathbf{y}) = \mathbf{y}_{s,i} \frac{\mathbf{x}_i - \mu(\mathbf{x}_i)}{\sigma(\mathbf{x}_i)} + \mathbf{y}_{b,i},$$

Samples (trained on the FFHQ dataset)

116

# StyleGAN (Karras et al., 2019)

- Swapping out the **destination** style for the **source** style



source

destination

Coarse styles copied

Latent $\mathbf{z} \in \mathcal{Z}$

Synthesis network $g$

Noise

Normalize

Mapping network $f$

Const 4×4×512

FC

A → style → AdaIN

B

FC

Conv 3×3

FC

B

FC

A → style → AdaIN

4×4

FC

FC

Upsample

FC

Conv 3×3

FC

B

FC

A → style → AdaIN

$\mathbf{w} \in \mathcal{W}$

Conv 3×3

B

A → style → AdaIN

8×8

...

# Some Applications of GANs

# Semi-supervised Classification

## SVNH

| Model | Misclassification rate |
|---|---|
| VAE (M1 + M2) (Kingma et al., 2014) | 36.02 |
| SWWAE with dropout (Zhao et al., 2015) | 23.56 |
| DCGAN + L2-SVM (Radford et al., 2015) | 22.18 |
| SDGM (Maaløe et al., 2016) | 16.61 |
| **GAN (feature matching) (Salimans et al., 2016)** | $\mathbf{8.11 \pm 1.3}$ |
| ALI (ours, L2-SVM) | $19.14 \pm 0.50$ |
| **ALI (ours, no feature matching)** | $\mathbf{7.42 \pm 0.65}$ |

# Class-specific Image Generation (Nguyen et al., 2016)

- Generates 227x227 realistic images from all ImageNet classes

- Combines adversarial training, moment matching, denoising autoencoders, and Langevin sampling



Noiseless joint PPGN-$h$    Image classifier

classes



redshank          ant          monastery          volcano

# Video Generation (Vondrick et al., 2016)



Beach          Golf          Train Station

# Generative Shape Modeling (Wu et al., 2016)



512×4×4×4

256×8×8×8

128×16×16×16

64×32×32×32

z

G(z) in 3D Voxel Space
64×64×64

Chairs

Sofas

The small bird has a red head with feathers that fade from red to gray from head to tail

This bird is black with green and has a very short beak

# Single Image Super-Resolution (Ledig et al., 2016)

- Combine content loss with adversarial loss



bicubic        SRResNet        SRGAN        original

4× upscaling

# Image Inpainting (Pathak et al., 2016)

# Unsupervised Domain Adaptation (Bousmalis et al., 2016)



Image examples from the Linemod dataset

RGDB image samples
(conditioned on a synthetic image)

# Image to Image Translation (Pix2Pix)



(Isola et al. 2016)

$$\arg\min_{G}\max_{D} \; \mathbb{E}_{\mathbf{x},\mathbf{y}}[\; \log D(G(\mathbf{x})) \;\; + \;\; \log(1 - D(\mathbf{y}))\;]$$

$$\arg\min_{G}\max_{D}\ \mathbb{E}_{\mathbf{x},\mathbf{y}}\big[\ \log D(\mathbf{x}, G(\mathbf{x})) + \log(1 - D(\mathbf{x}, \mathbf{y}))\ \big]$$

$$\arg \min_G \max_D \ \mathbb{E}_{\mathbf{x},\mathbf{y}}\big[\ \log D(\mathbf{x}, G(\mathbf{x})) + \log(1 - D(\mathbf{x}, \mathbf{y}))\ \big]$$
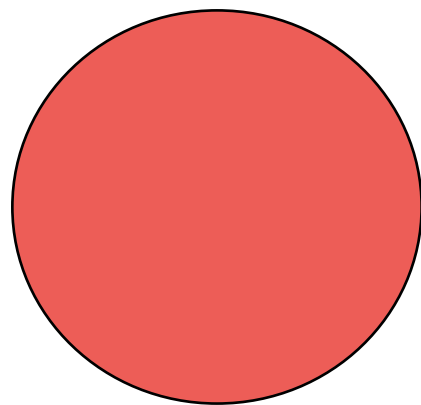
$$\arg \min_G \max_D \; \mathbb{E}_{\mathbf{x},\mathbf{y}}\big[\; \log D(\mathbf{x}, G(\mathbf{x})) + \log(1 - D(\mathbf{x}, \mathbf{y})) \;\big]$$

$$\arg\min_{G}\max_{D} \; \mathbb{E}_{\mathbf{x},\mathbf{y}}\big[ \; \log D(\mathbf{x}, G(\mathbf{x})) + \log(1 - D(\mathbf{x}, \mathbf{y})) \; \big]$$

# BW → Color

Input  Output     Input  Output     Input  Output



Data from [Russakovsky et al. 2015]

# #edges2cats [Chris Hesse]

INPUT

OUTPUT

pix2pix

process

Ivy Tasi @ivymyt

Vitaly Vidmirov @vvid

# Shrinking the capacity: Patch Discriminator



Rather than penalizing if output image looks fake, penalize if each overlapping patch in output looks fake

- Faster, fewer parameters
- More supervised observations
- Applies to arbitrarily large images

[Li & Wand 2016]
[Shrivastava et al. 2017]
[Isola et al. 2017]

# Labels → Facades

Input

1x1 Discriminator



Data from [Tylecek, 2013]

# Labels → Facades

Input

Data from [Tylecek, 2013]

# Labels → Facades

Input

Data from [Tylecek, 2013]

# Labels → Facades

Input

Full image Discriminator



Data from [Tylecek, 2013]

# Pix2Pix w/o input-output pairs



(Zhu et al. 2017)

# Paired data

$$x_i \qquad y_i$$

# Paired data

$x_i$     $y_i$



# Unpaired data

$X$     $Y$

$$\arg \min_{G} \max_{D} \; \mathbb{E}_{\mathbf{x},\mathbf{y}}\big[\; \log D(\mathbf{x}, G(\mathbf{x})) + \log(1 - D(\mathbf{x}, \mathbf{y}))\; \big]$$

$$\arg \min_{G} \max_{D} \; \mathbb{E}_{\mathbf{x},\mathbf{y}} \big[ \; \log D(\mathbf{x}, G(\mathbf{x})) + \log(1 - D(\mathbf{x}, \mathbf{y})) \; \big]$$

No input-output pairs!

$$\arg \min_{G} \max_{D} \; \mathbb{E}_{\mathbf{x},\mathbf{y}}[ \; \log D(G(\mathbf{x})) \; + \; \log(1 - D(\mathbf{y})) \; ]$$

Usually loss functions check if output matches a target instance

GAN loss checks if output is part of an admissible set

Gaussian

Target distribution



z

Y

Horses → Zebras

X

Y

$\mathbf{x}$     $G$     $G(\mathbf{x})$     $D$     Real!

$\mathbf{x}$  $G$  $G(\mathbf{x})$  $D$  Real too!

Nothing to force output to correspond to input

# Cycle-Consistent Adversarial Networks



[Zhu et al. 2017], [Yi et al. 2017], [Kim et al. 2017]

# Cycle-Consistent Adversarial Networks

# Cycle Consistency Loss

# Cycle Consistency Loss

# Collection Style Transfer
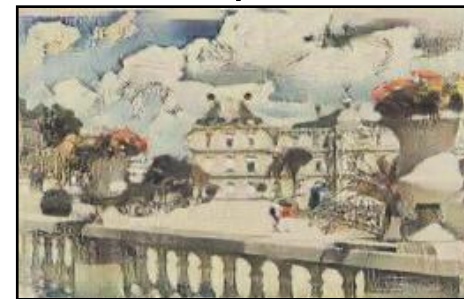


Photograph
@ Alexei Efros

Monet

Van Gogh

Cezanne

Ukiyo-e

| Input | Monet | Van Gogh | Cezanne | Ukiyo-e |

# Monet's paintings → photos

# Monet's paintings → photos

# Failure case

# Failure case

# Semantic Image Synthesis (SPADE) (Park et al., 2019)

- Image generation conditioned on semantic layouts



Semantic Manipulation Using Segmentation Map

Style Manipulation using Style Images

**Manipulating Attributes of Natural Scenes via Hallucination.**
Levent Karacan, Zeynep Akata, Aykut Erdem & Erkut Erdem.
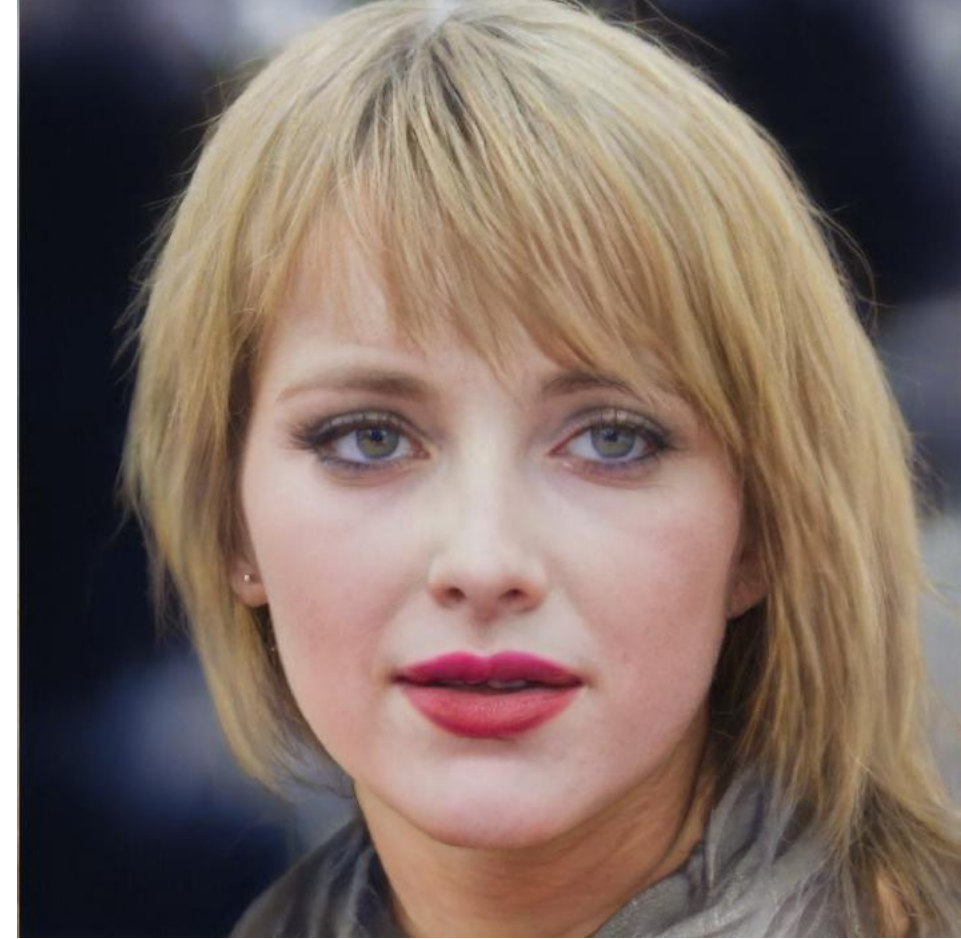ACM Trans. on Graphics, Vol. 39, Issue 1, Article 7, February 2020.

snow

night

prediction

166

Spring
+
Clouds

prediction

A young woman with bangs wearing lipstick

168

An old and grumpy British shorthair

**CLIP-Guided StyleGAN Inversion for Text-Driven Real Image Editing.**
Canberk Baykal, Abdul Basit Anees, Duygu Ceylan, Aykut Erdem, Erkut Erdem, Deniz Yuret
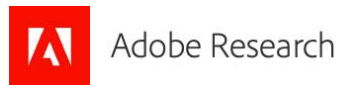ACM Transactions on Graphics, 2023

169

green jacket

Sleeveless blue blouse

black short

170

**Audio-based Image Editing and Generation using Latent Diffusion Models**
Burak Can Biner, Farrin Marouf Sofian,
Umur Berkay Karakaş, Duygu Ceylan, Erkut Erdem,
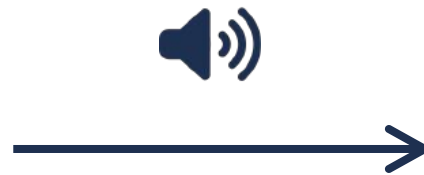Aykut Erdem. In progress

171

**Audio-based Image Editing and Generation using Latent Diffusion Models**
Burak Can Biner, Farrin Marouf Sofian,
Umur Berkay Karakaş, Duygu Ceylan, Erkut Erdem,
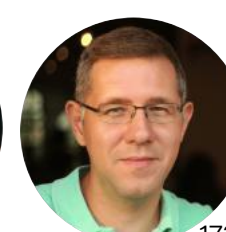Aykut Erdem. In progress

**Audio-based Image Editing and Generation using Latent Diffusion Models**
Burak Can Biner, Farrin Marouf Sofian,
Umur Berkay Karakaş, Duygu Ceylan, Erkut Erdem,
Aykut Erdem. In progress

# Next lecture:
# Autoregressive and Flow Models