

# CS 4803 / 7643: Deep Learning

Topics:

- Unsupervised Learning
- Generative Models (PixelRNNs, VAEs)

Dhruv Batra  
Georgia Tech

# Administrativa

- HW1 and HW2 solutions released
  - <https://gatech.instructure.com/courses/28059/files/>
- HW3 out
  - Due: 11/06, 11:55pm

# Overview

- Unsupervised Learning
- Generative Models
  - PixelRNN and PixelCNN
  - Variational Autoencoders (VAE)
  - Generative Adversarial Networks (GAN)

# Supervised vs Unsupervised Learning

## Supervised Learning

**Data:**  $(x, y)$

x is data, y is label

**Goal:** Learn a *function* to map  $x \rightarrow y$

**Examples:** Classification,  
regression, object detection,  
semantic segmentation, image  
captioning, etc.



# Supervised vs Unsupervised Learning

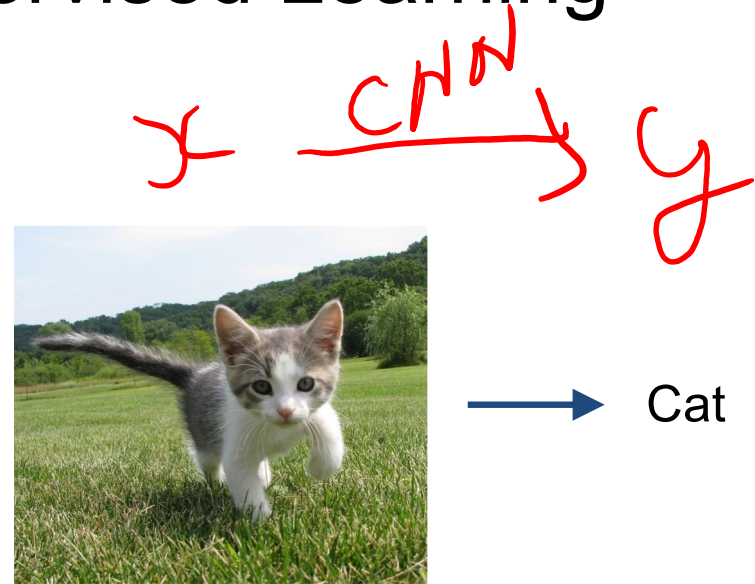
## Supervised Learning

**Data:**  $(x, y)$

$x$  is data,  $y$  is label

**Goal:** Learn a *function* to map  $x \rightarrow y$

**Examples:** Classification,  
regression, object detection,  
semantic segmentation, image  
captioning, etc.



Classification

[This image](#) is [CC0 public domain](#)

# Supervised vs Unsupervised Learning

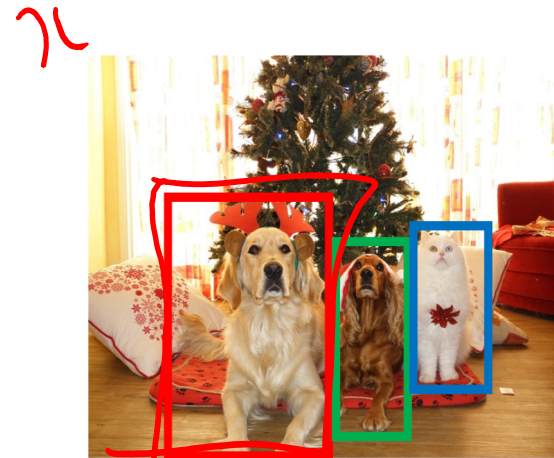
## Supervised Learning

**Data:**  $(x, y)$

$x$  is data,  $y$  is label

**Goal:** Learn a function to map  $x \rightarrow y$

**Examples:** Classification, regression, object detection, semantic segmentation, image captioning, etc.



**DOG, DOG, CAT**

Object Detection

This image is [CC0 public domain](#)

# Supervised vs Unsupervised Learning

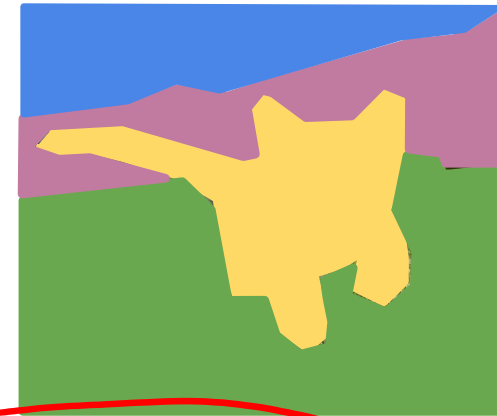
## Supervised Learning

**Data:**  $(x, y)$

$x$  is data,  $y$  is label

**Goal:** Learn a *function* to map  $x \rightarrow y$

**Examples:** Classification, regression, object detection, semantic segmentation, image captioning, etc.



GRASS, CAT,  
TREE, SKY

Semantic Segmentation

# Supervised vs Unsupervised Learning

## Supervised Learning

**Data:**  $(x, y)$

$x$  is data,  $y$  is label

**Goal:** Learn a *function* to map  $x \rightarrow y$

**Examples:** Classification, regression, object detection, semantic segmentation, image captioning, etc.



*A cat sitting on a suitcase on the floor*

Image captioning

$x$

$y$

Caption generated using [neuraltalk2](#)  
Image is [CC0 Public domain](#).

# Supervised vs Unsupervised Learning

## 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

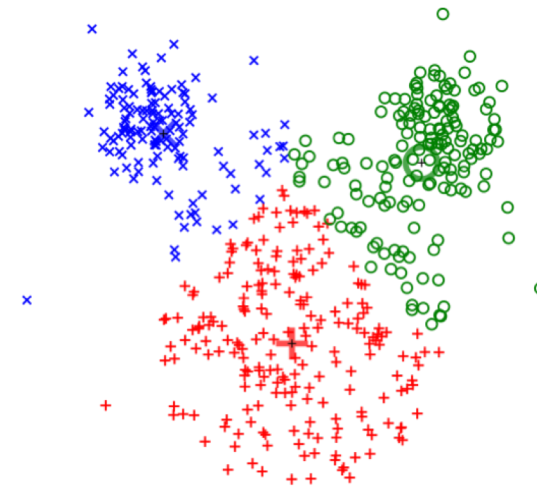
## Unsupervised Learning

**Data:**  $x$

Just data, no labels!

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

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



K-means clustering

[This image is CC0 public domain](#)

# Supervised vs Unsupervised Learning

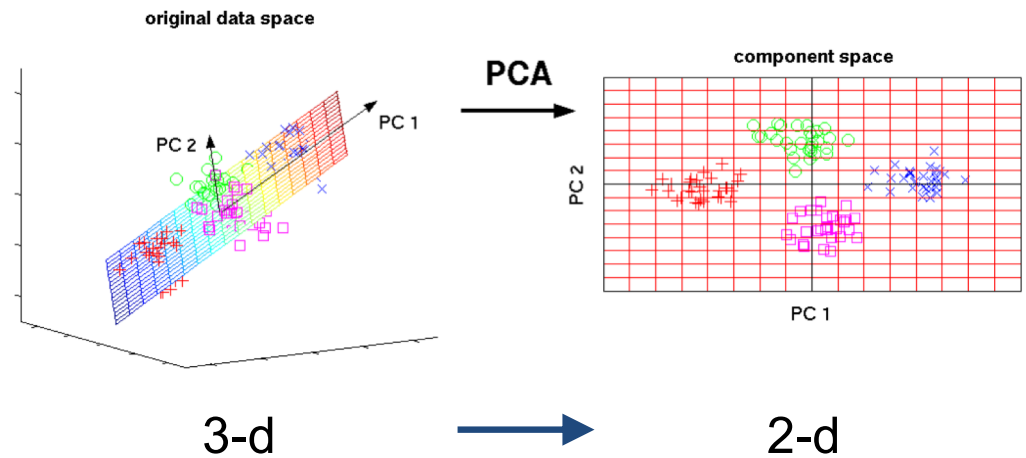
## Unsupervised Learning

**Data:**  $x$

Just data, no labels!

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

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



Principal Component Analysis  
(Dimensionality reduction)

[This image](#) from Matthias Scholz is [CC0 public domain](#)

~~$P(x)$~~

$P(x, \dots)$

# Supervised vs Unsupervised Learning

## 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.

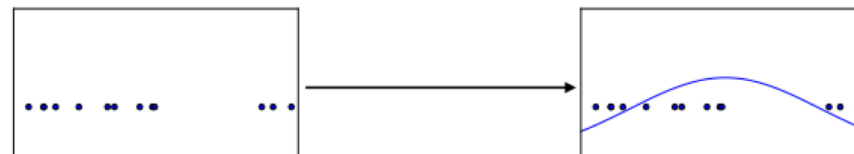
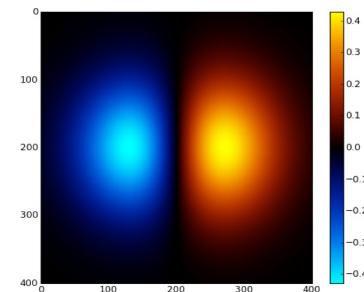
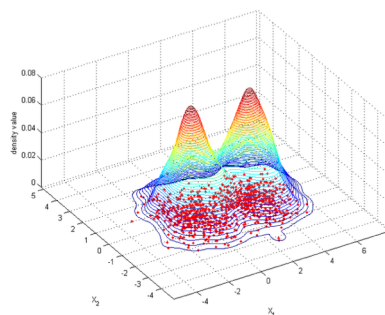


Figure copyright Ian Goodfellow, 2016. Reproduced with permission.

1-d density estimation



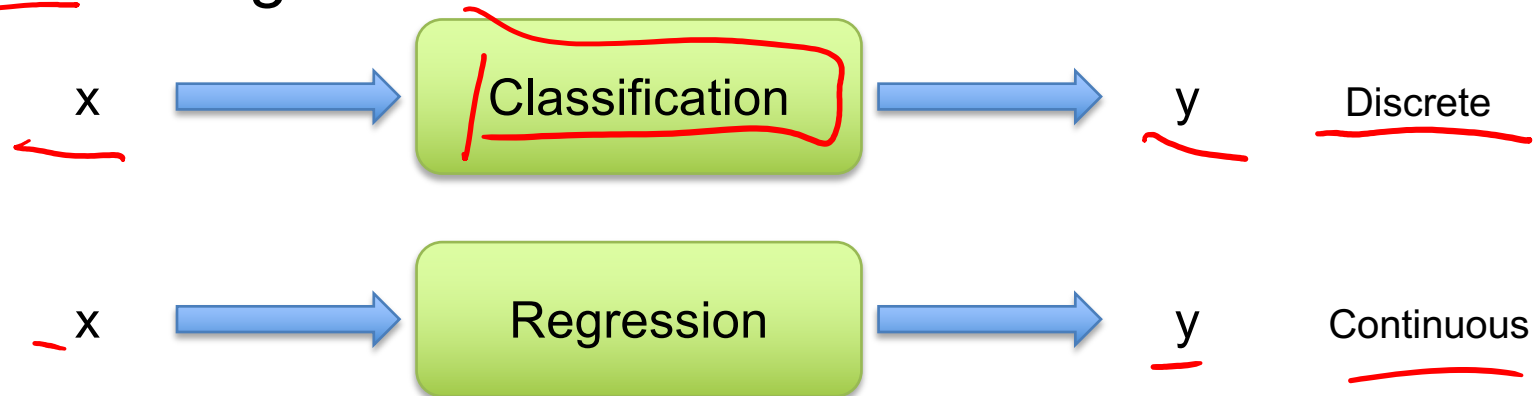
2-d density estimation

2-d density images [left](#) and [right](#) are [CC0 public domain](#).

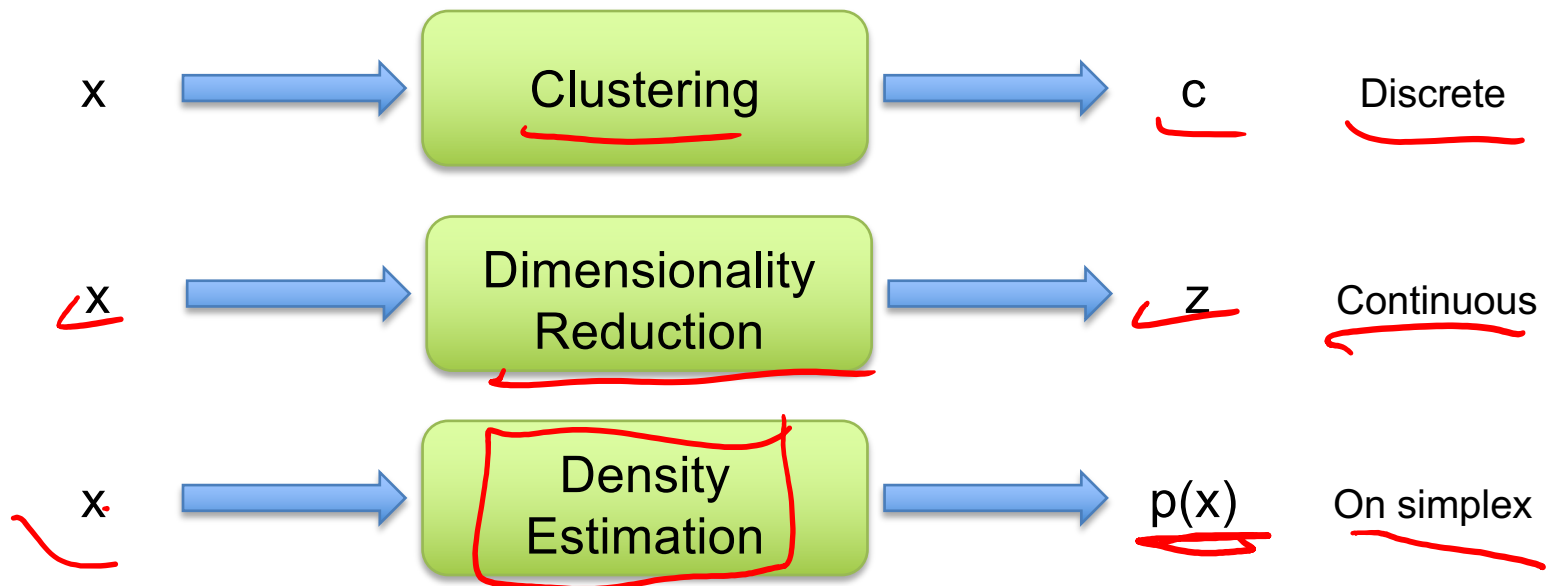


# Tasks

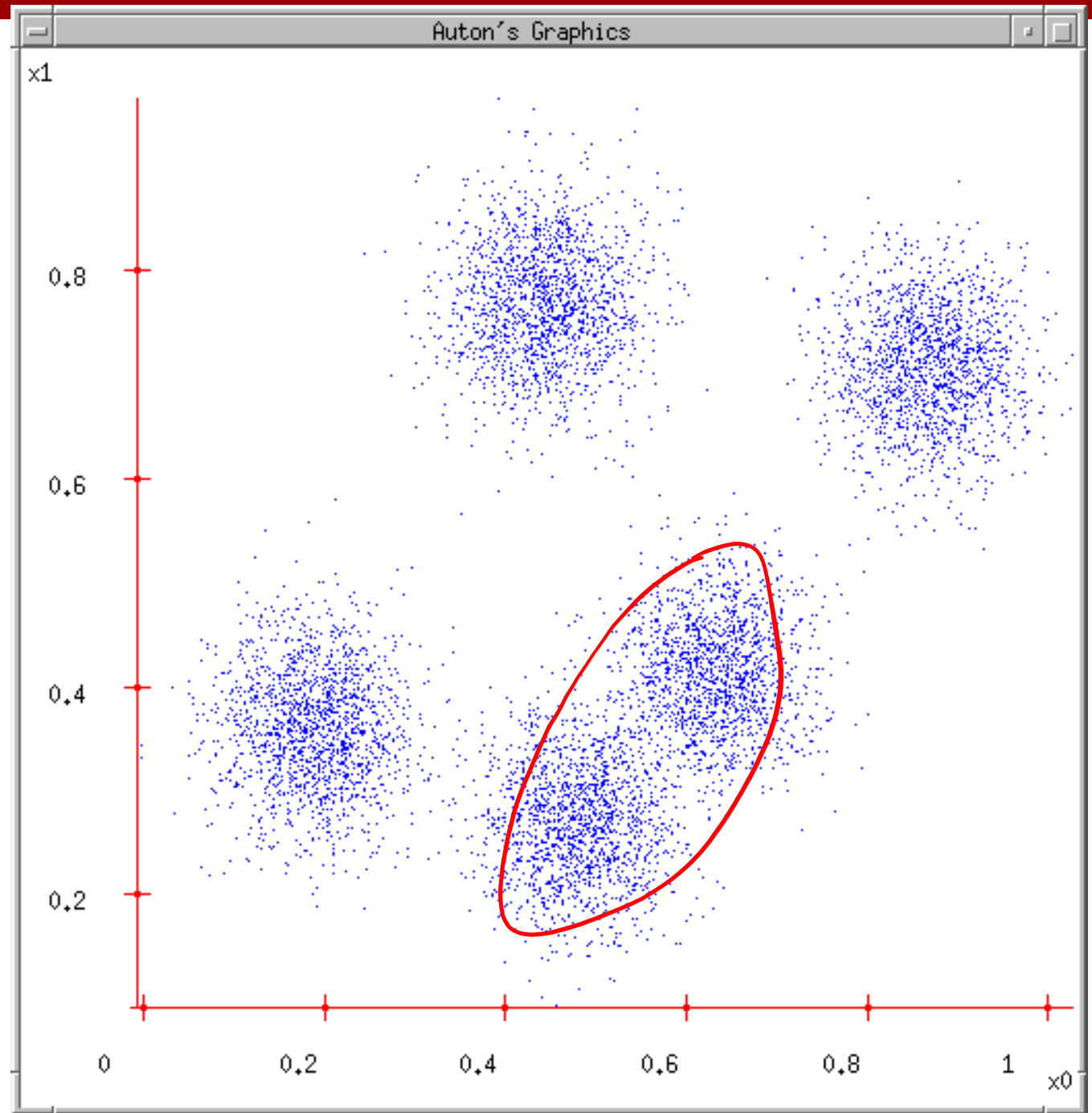
## Supervised Learning



## Unsupervised Learning

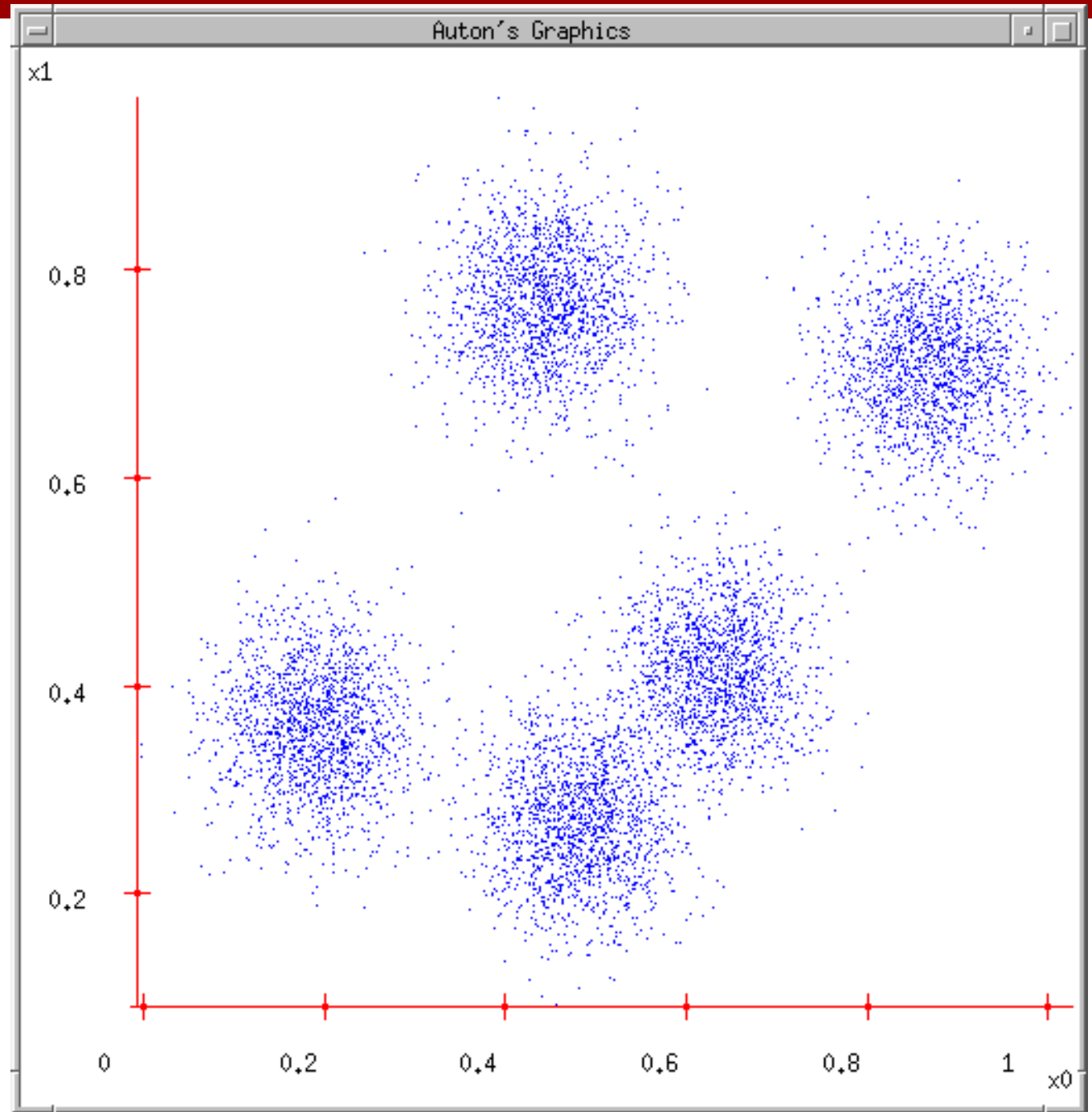


# Some Data



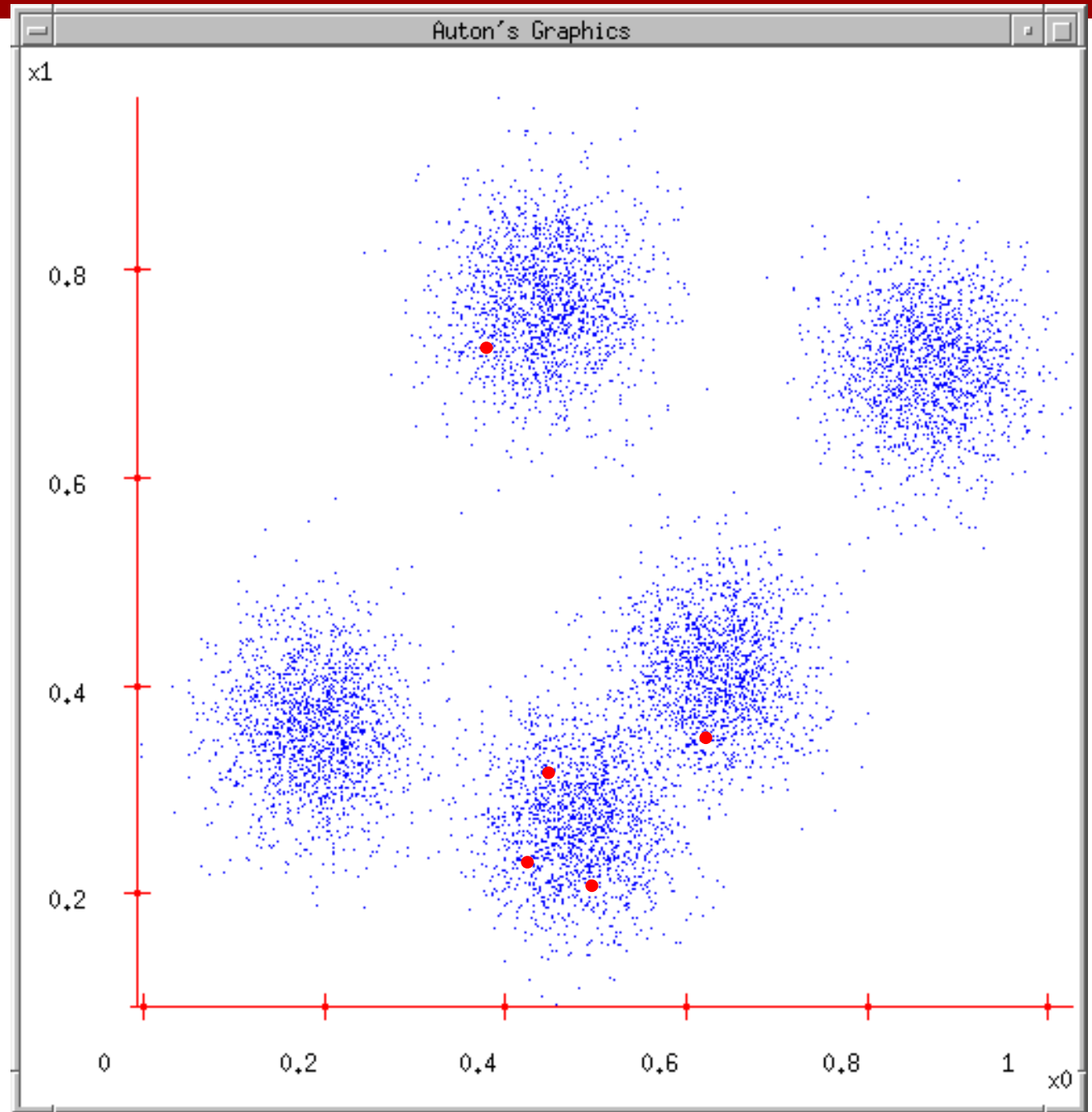
# K-means

1. Ask user how many clusters they'd like.  
*(e.g.  $k=5$ )*



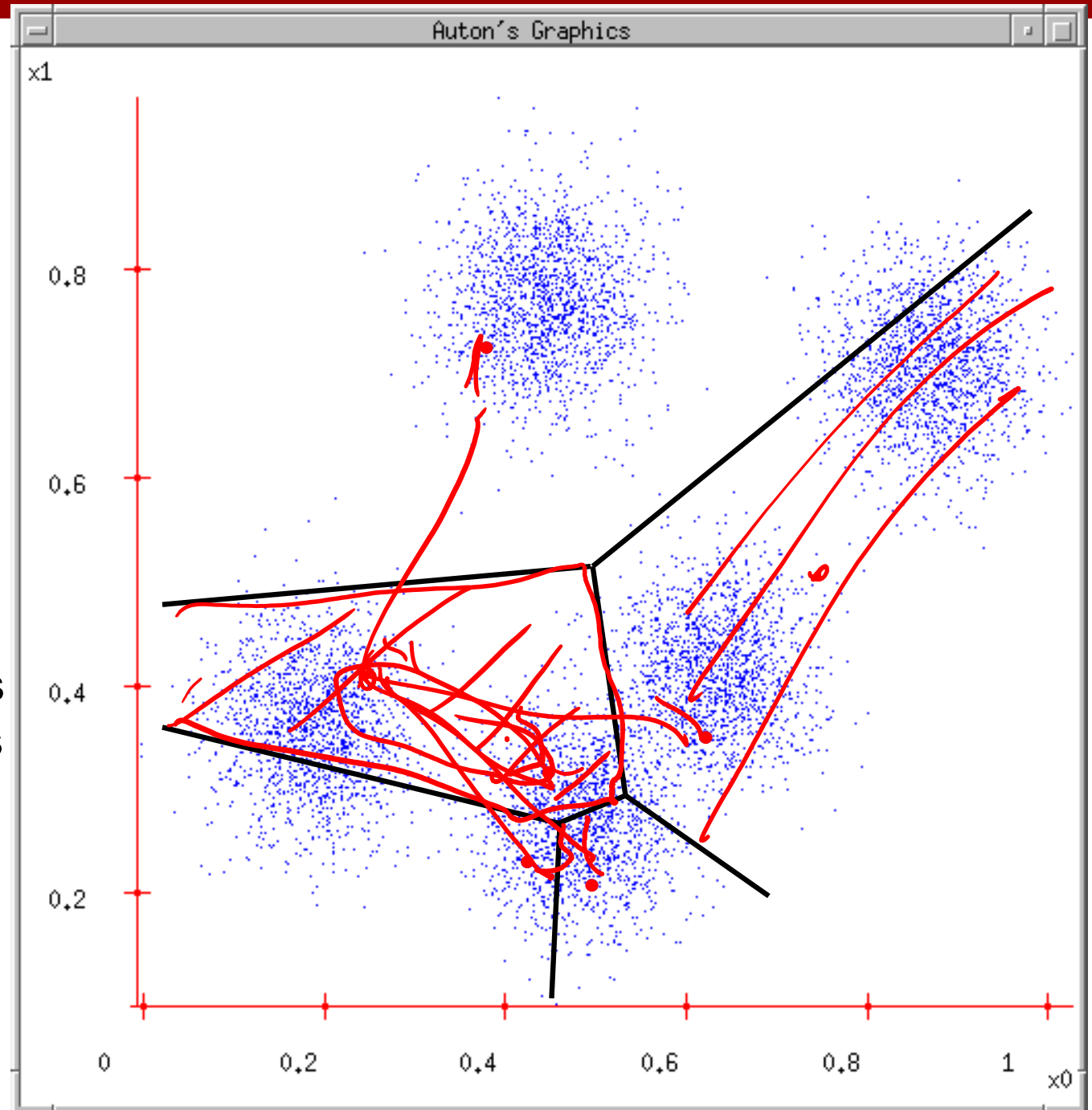
# K-means

1. Ask user how many clusters they'd like.  
(e.g.  $k=5$ )
2. Randomly guess  $k$  cluster Center locations



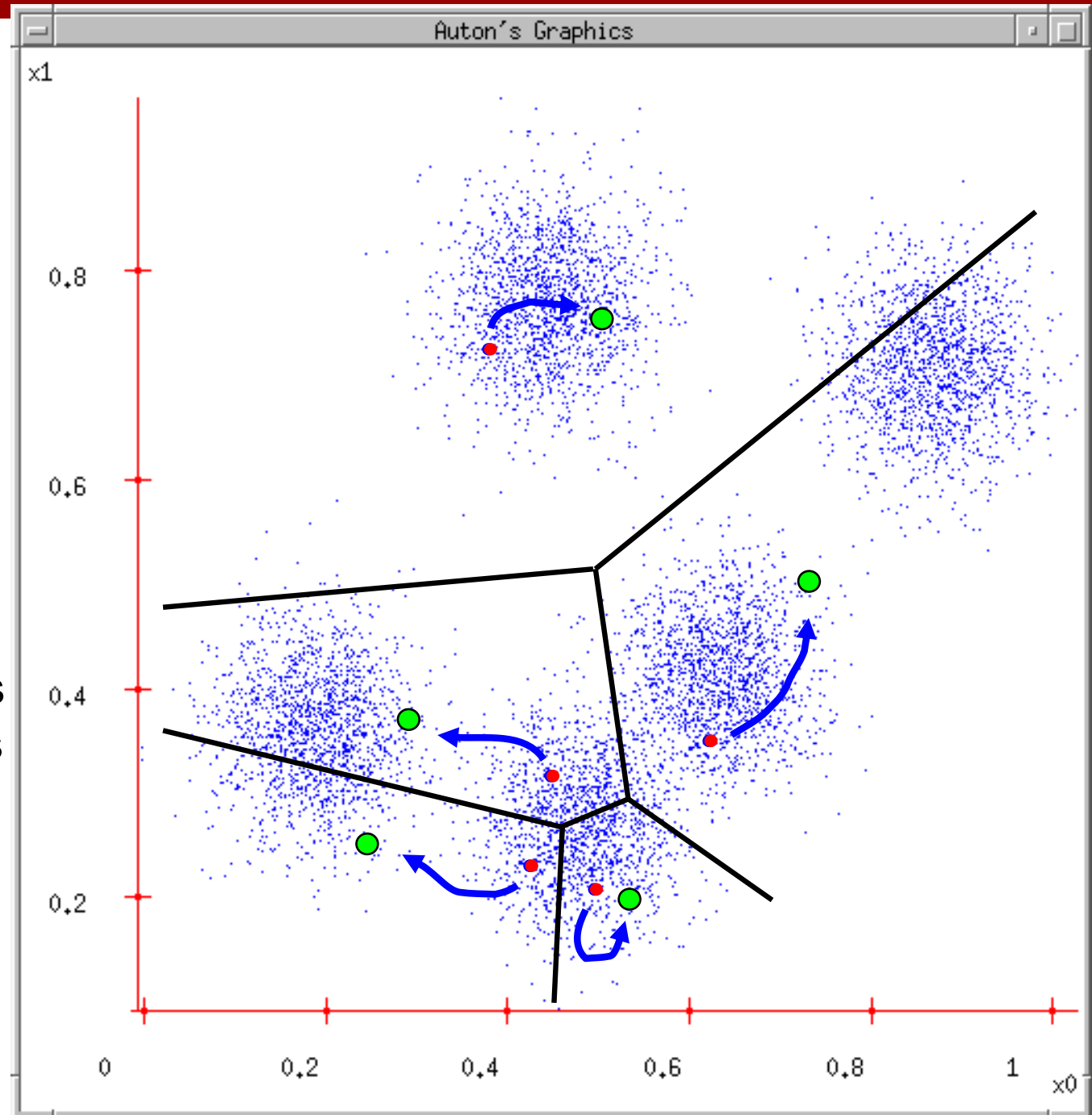
# K-means

1. Ask user how many clusters they'd like.  
(e.g.  $k=5$ )
2. Randomly guess  $k$  cluster Center locations
3. Each datapoint finds out which Center it's closest to. (Thus each Center "owns" a set of datapoints)



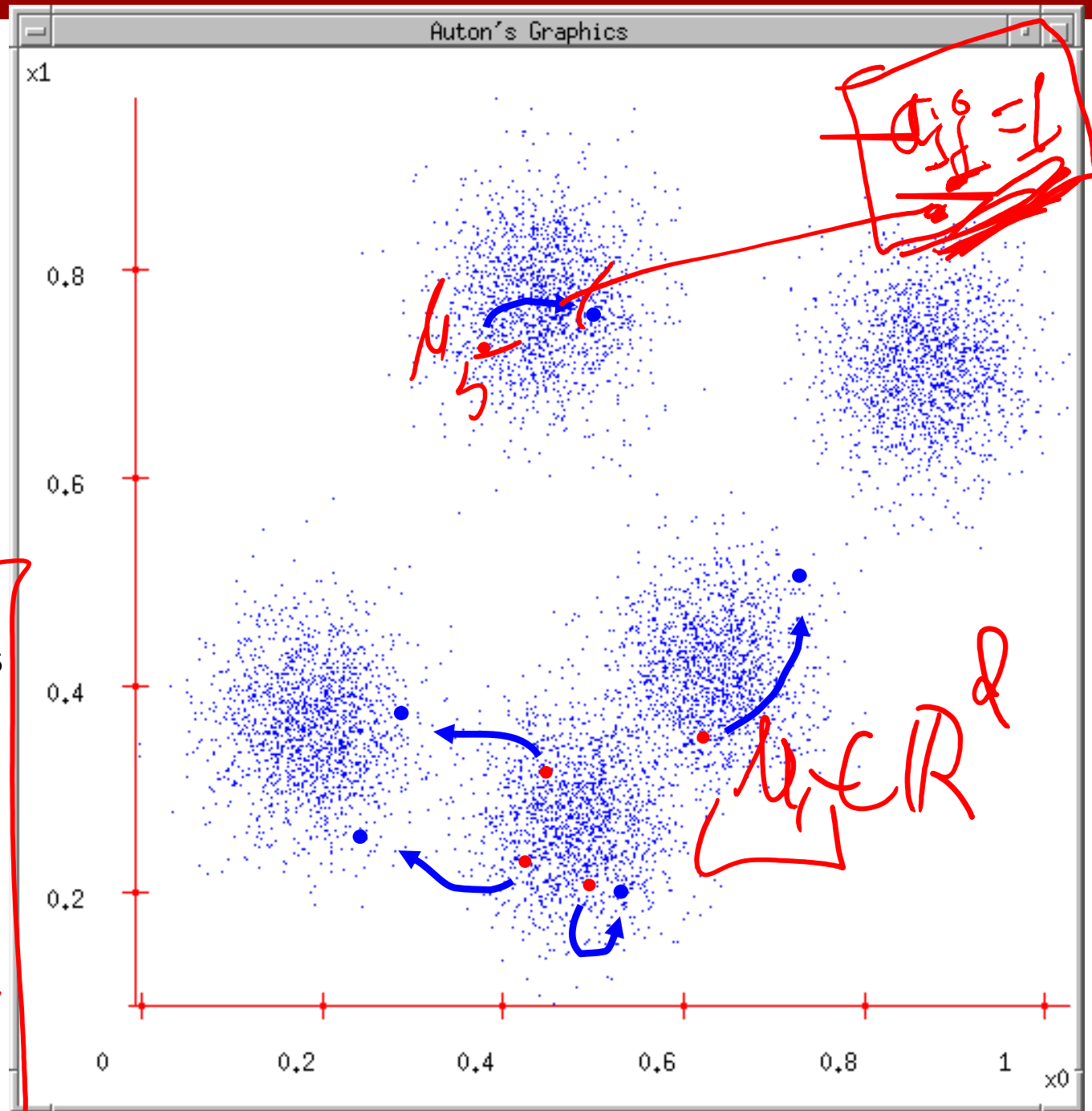
# K-means

1. Ask user how many clusters they'd like.  
(e.g.  $k=5$ )
2. Randomly guess  $k$  cluster Center locations
3. Each datapoint finds out which Center it's closest to.
4. Each Center finds the centroid of the points it owns



# K-means

1. Ask user how many clusters they'd like.  
(e.g.  $k=5$ )
2. Randomly guess  $k$  cluster Center locations
3. Each datapoint finds out which Center it's closest to.
4. Each Center finds the centroid of the points it owns...
5. ...and jumps there
6. ...Repeat until terminated!



# K-means

- Randomly initialize  $k$  centers

- $\mu^{(0)} = \mu_1^{(0)}, \dots, \mu_k^{(0)}$

- **Assign:**

- Assign each point  $i \in \{1, \dots, n\}$  to nearest center:

- $C(i) \leftarrow \underset{j}{\operatorname{argmin}} \|\mathbf{x}_i - \mu_j\|^2$

- **Recenter:**

- $\mu_j$  becomes centroid of its points



# K-means

- Demo
  - <http://stanford.edu/class/ee103/visualizations/kmeans/kmeans.html>

# Supervised vs Unsupervised Learning

## Supervised Learning

**Data:**  $(x, y)$

$x$  is data,  $y$  is label

**Goal:** Learn a *function* to map  $x \rightarrow y$

**Examples:** Classification, regression, object detection, semantic segmentation, image captioning, 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

## Supervised Learning

**Data:**  $(x, y)$

$x$  is data,  $y$  is label

**Goal:** Learn a *function* to map  $x \rightarrow y$

**Examples:** Classification, regression, object detection, semantic segmentation, image captioning, etc.

## Unsupervised Learning

Training data is cheap

**Data:**  $x$

Just data, no labels!

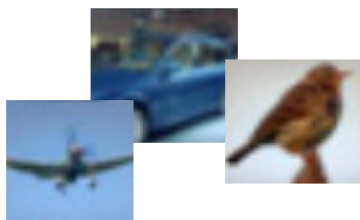
**Goal:** Learn some underlying hidden *structure* of the data

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

Holy grail: Solve unsupervised learning  
=> understand structure of visual world

# Generative Models

Given training data, generate new samples from same distribution



Training data  $\sim p_{\text{data}}(x)$



Generated samples  $\sim p_{\text{model}}(x)$

Want to learn  $p_{\text{model}}(x)$  similar to  $p_{\text{data}}(x)$

$p_{\text{model}}(x)$

# Generative Classification vs Discriminative Classification vs Density Estimation

- Generative Classification

- Model  $p(x, y)$ ; estimate  $p(x|y)$  and  $p(y)$
- Use Bayes Rule to predict  $y$
- E.g Naïve Bayes

$$\frac{p(x, y)}{p(y|x) p(x|y)}$$

- Discriminative Classification

- Estimate  $p(y|x)$  directly
- E.g. Logistic Regression

$$p(x | y = \alpha) \left| \frac{p(y|x)}{p(x)} \right| p(x)$$

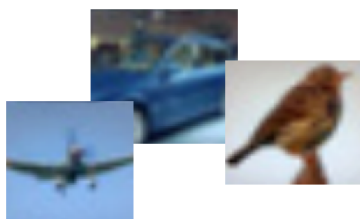
- Density Estimation

- Model  $p(x)$
- E.g. VAEs

$$p(x)$$

# Generative Models

Given training data, generate new samples from same distribution



Training data  $\sim p_{\text{data}}(x)$



Generated samples  $\sim p_{\text{model}}(x)$

Want to learn  $p_{\text{model}}(x)$  similar to  $p_{\text{data}}(x)$

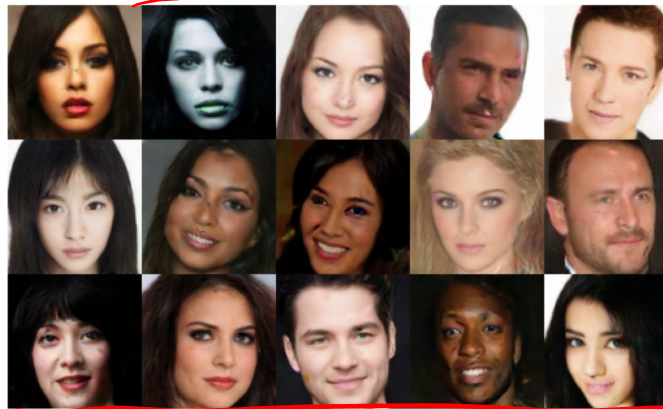
Addresses density estimation, a core problem in unsupervised learning

## Several flavors:

- Explicit density estimation: explicitly define and solve for  $p_{\text{model}}(x)$
- Implicit density estimation: learn model that can sample from  $p_{\text{model}}(x)$  w/o explicitly defining it

# Why Generative Models?

- Realistic samples for artwork, super-resolution, colorization, etc.



- Generative models of time-series data can be used for simulation and planning (reinforcement learning applications!)
- Training generative models can also enable inference of latent representations that can be useful as general features

Figures from L-R are copyright: (1) [Alec Radford et al. 2016](#); (2) [David Berthelot et al. 2017](#); [Phillip Isola et al. 2017](#). Reproduced with authors permission.

# Taxonomy of Generative Models

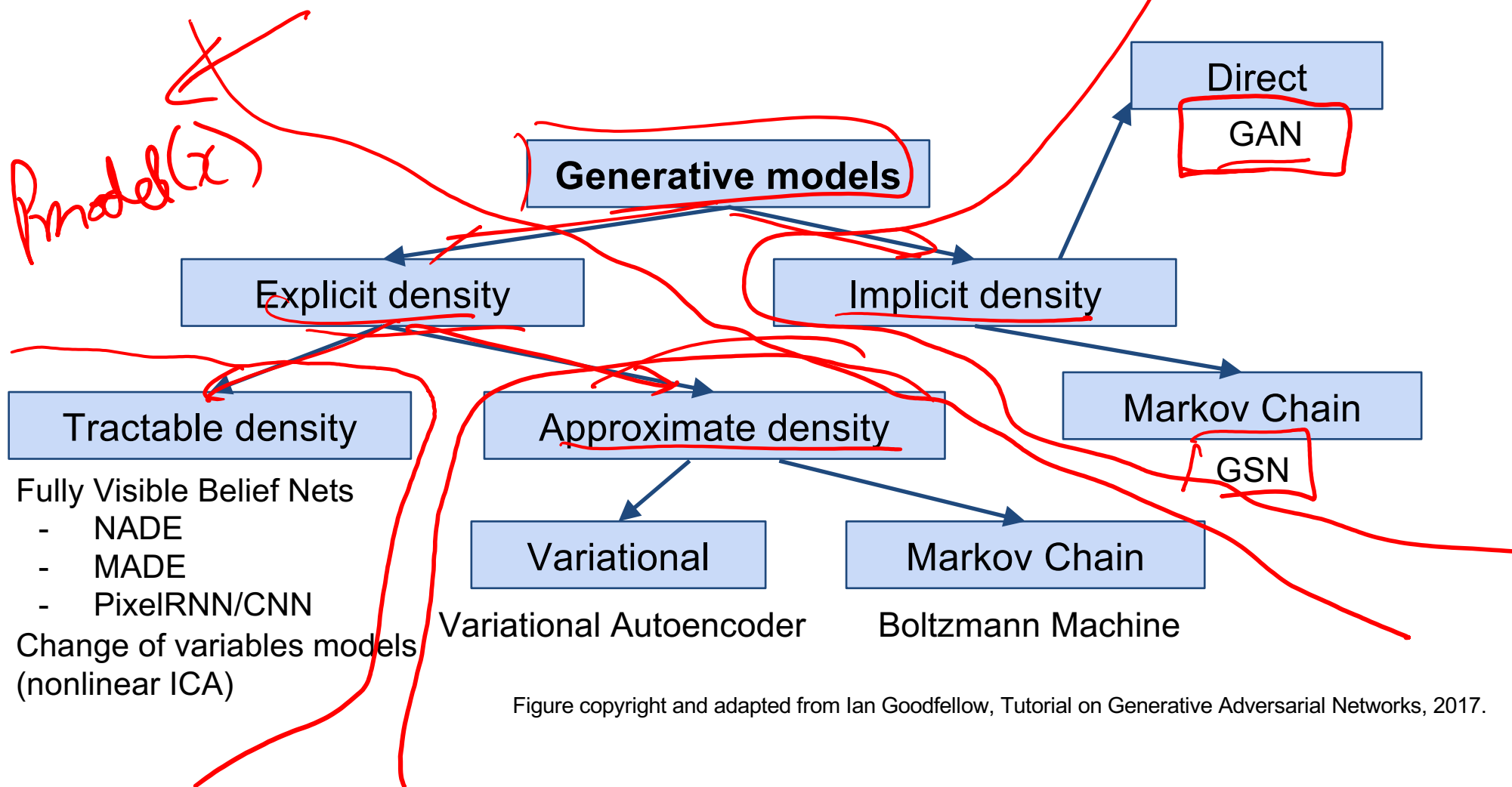
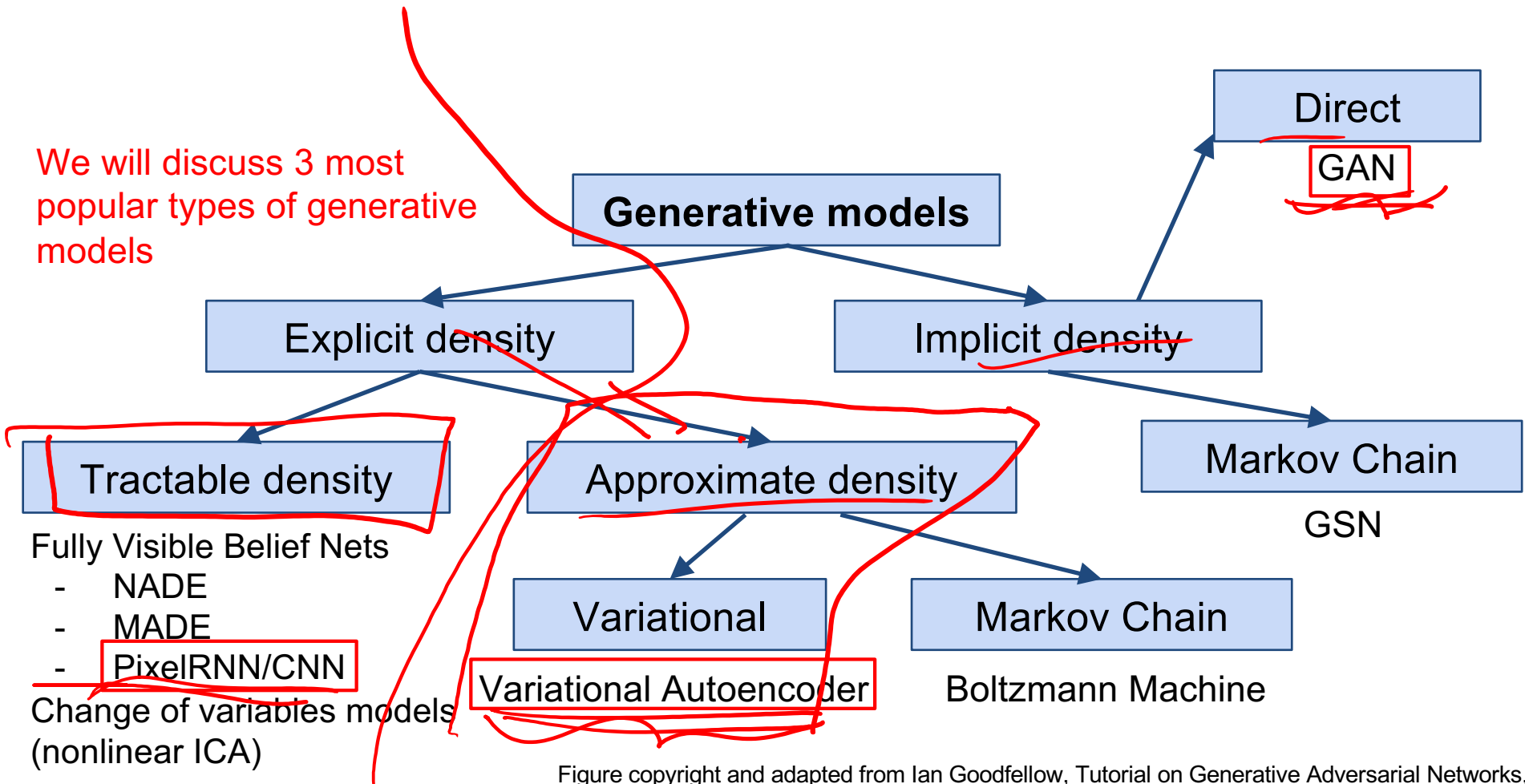


Figure copyright and adapted from Ian Goodfellow, Tutorial on Generative Adversarial Networks, 2017.



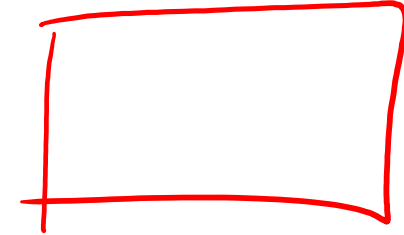
# Taxonomy of Generative Models



# PixelRNN and PixelCNN

# Fully Observable Model

Explicit density model



Use chain rule to decompose likelihood of an image  $x$  into product of 1-d distributions:

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

Likelihood of image  $x$       Probability of  $i$ 'th pixel value given all previous pixels

Then maximize likelihood of training data

$$D = \{x_1, \dots, x_n\}$$

$$x_1, \dots, x_{i-1} \rightarrow \boxed{\text{NN}}_{\theta} \rightarrow x_i \quad \max_{\theta} \log P_{\theta}(x)$$

# Fully Observable Model

## Explicit density model

Use chain rule to decompose likelihood of an image  $x$  into product of 1-d distributions:

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

↑  
Likelihood of  
image  $x$

↑  
Probability of  $i$ 'th pixel value  
given all previous pixels

Complex distribution over pixel values  
=> Express using a neural network!

Then maximize likelihood of training data

# Fully Observable Model

## Explicit density model

Use chain rule to decompose likelihood of an image  $x$  into product of 1-d distributions:

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

↑ Likelihood of image  $x$

↑ Probability of  $i$ 'th pixel value given all previous pixels

Will need to define ordering of "previous pixels"

Complex distribution over pixel values  
=> Express using a neural network!

Then maximize likelihood of training data

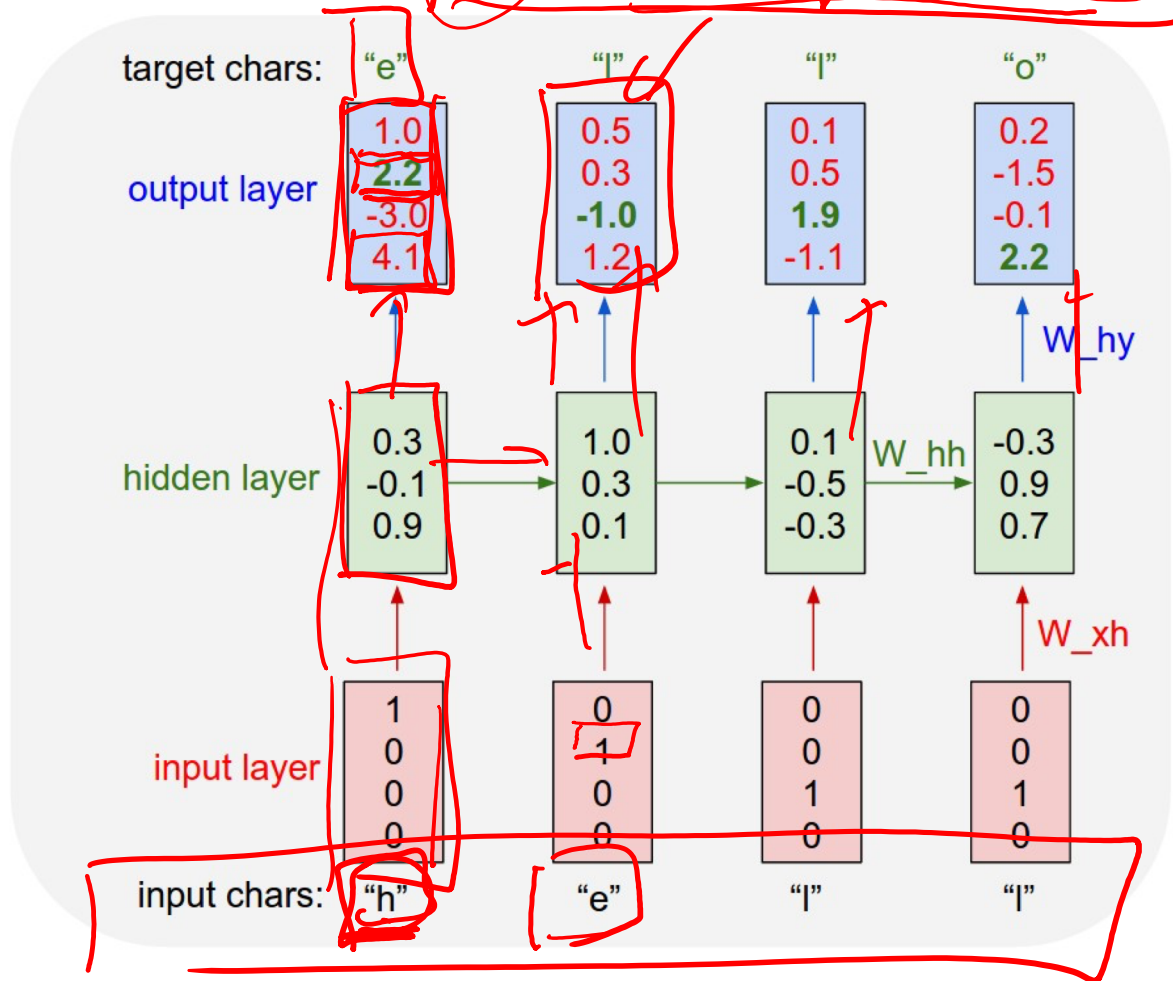
$\max_w$

$$\log P(x_1 \dots x_t | w) = \sum_t \log P(x_t | x_{1:t-1}, w)$$

# Example: Character-level Language Model

Vocabulary:  
[h,e,l,o]

Example training  
sequence:  
"hello"

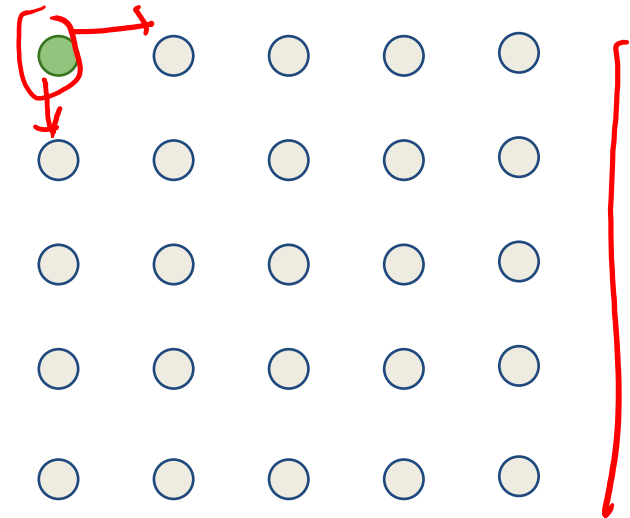


# PixelRNN [van der Oord et al. 2016]

$$P(x_1, \dots, x_{H \times W})$$

Generate image pixels starting from corner

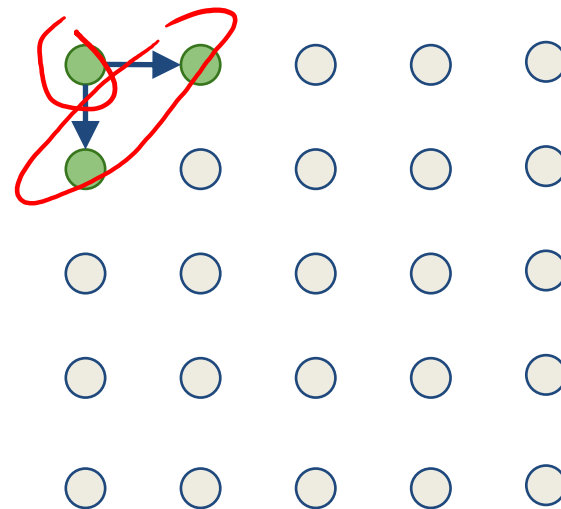
Dependency on previous pixels modeled using an RNN (LSTM)



# PixelRNN [van der Oord et al. 2016]

Generate image pixels starting from corner

Dependency on previous pixels modeled using an RNN (LSTM)

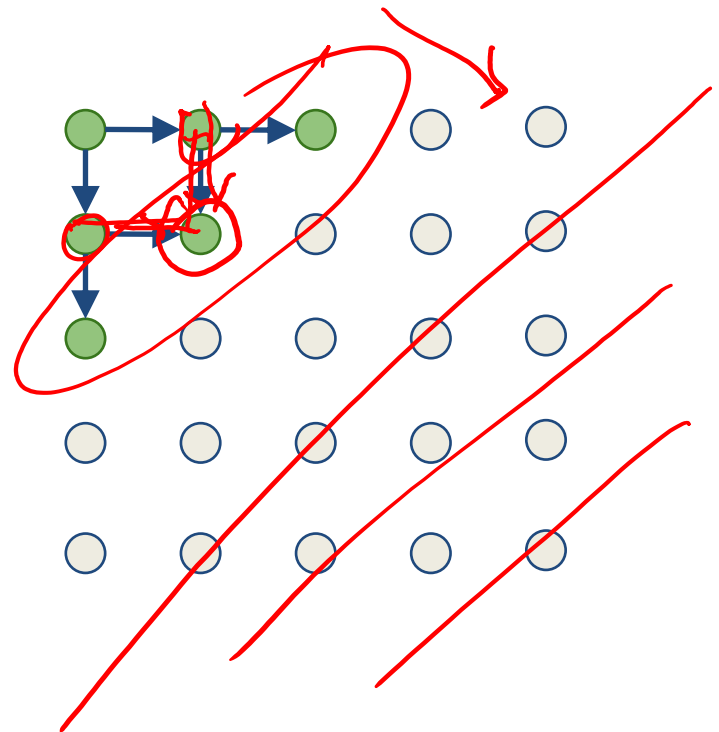




# PixelRNN [van der Oord et al. 2016]

Generate image pixels starting from corner

Dependency on previous pixels modeled using an RNN (LSTM)

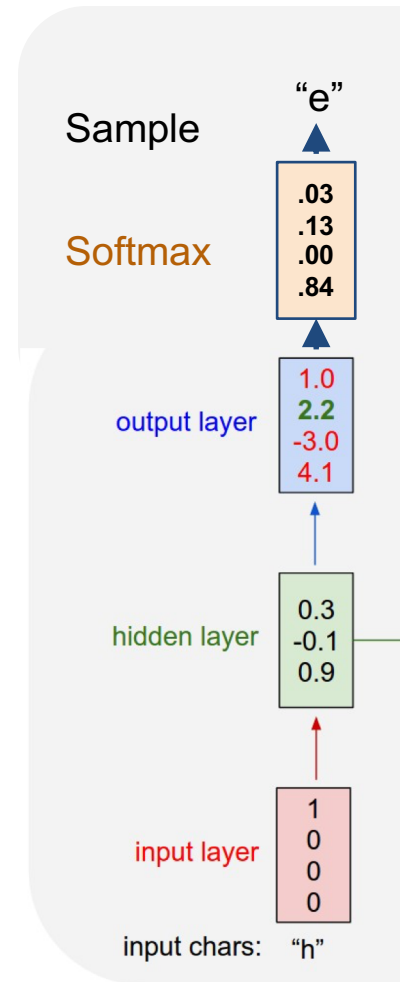


# Test Time: Sample / Argmax / Beam Search

## Example: Character-level Language Model Sampling

Vocabulary:  
[h,e,l,o]

At test-time sample  
characters one at a  
time, feed back to  
model

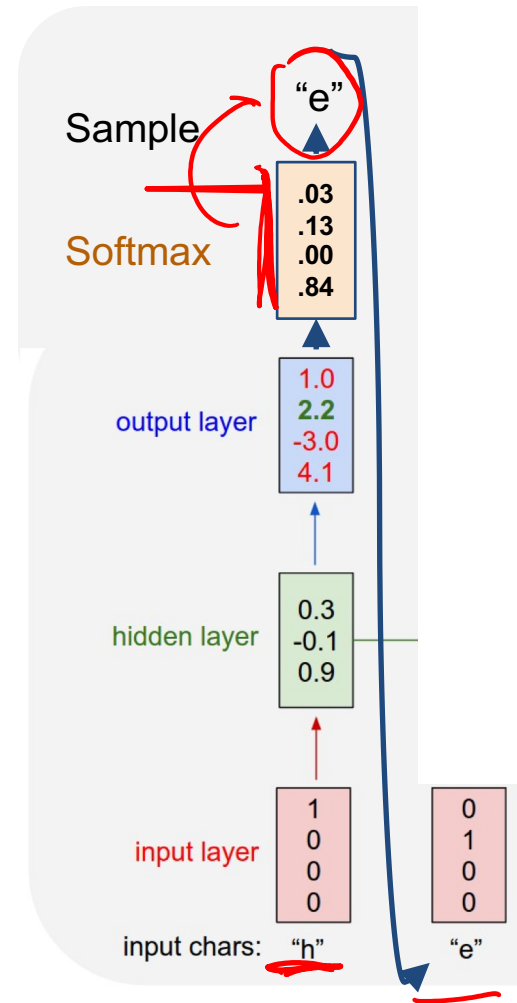


# Test Time: Sample / Argmax / Beam Search

## Example: Character-level Language Model Sampling

Vocabulary:  
[h,e,l,o]

At test-time sample  
characters one at a  
time, feed back to  
model

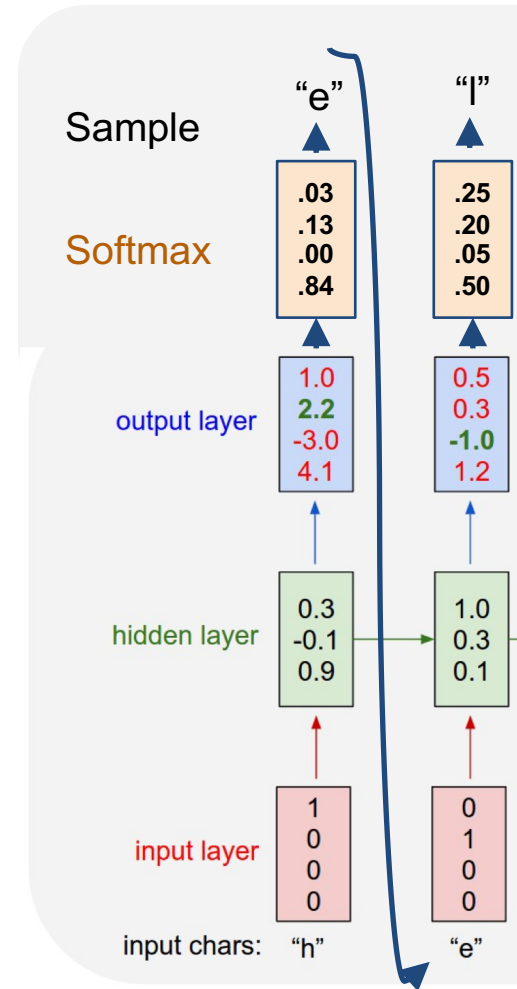


# Test Time: Sample / Argmax / Beam Search

## Example: Character-level Language Model Sampling

Vocabulary:  
[h,e,l,o]

At test-time sample  
characters one at a  
time, feed back to  
model

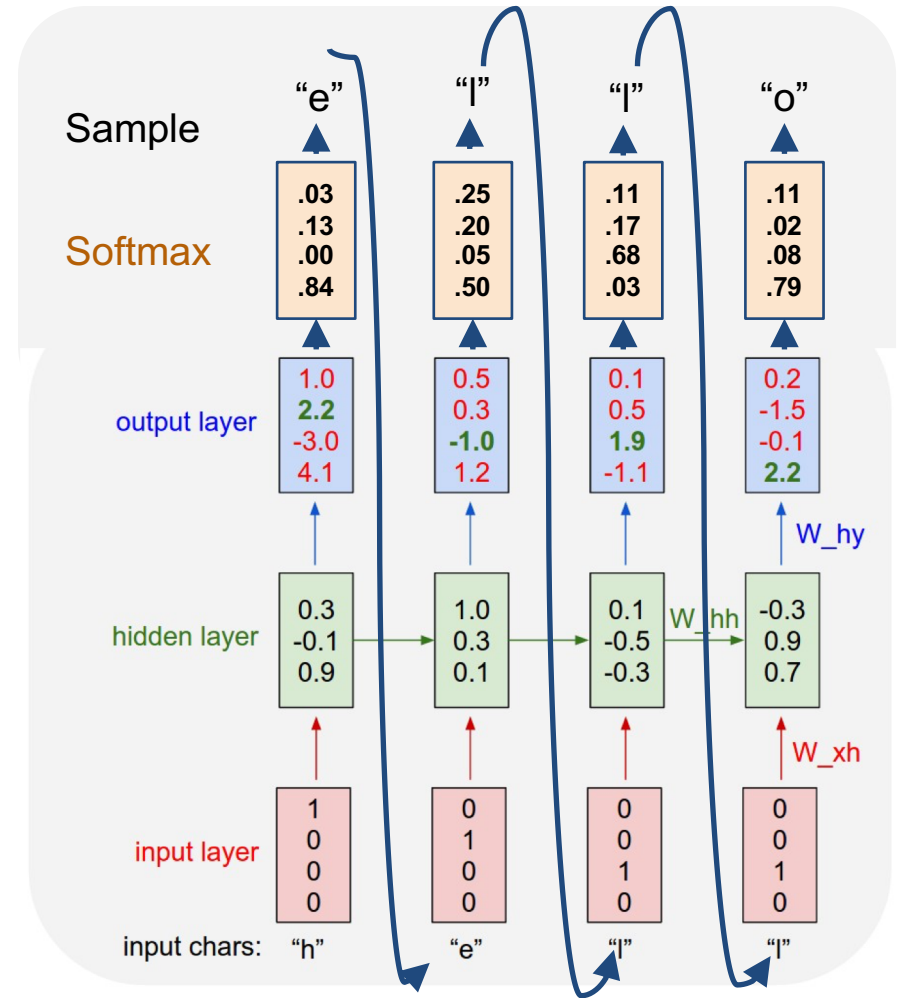


# Test Time: Sample / Argmax / Beam Search

## Example: Character-level Language Model Sampling

Vocabulary:  
[h,e,l,o]

At test-time sample  
characters one at a  
time, feed back to  
model

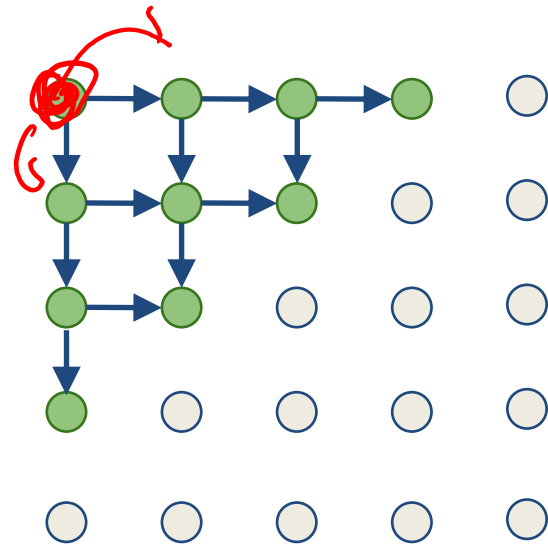


# PixelRNN [van der Oord et al. 2016]

Generate image pixels starting from corner

Dependency on previous pixels modeled using an RNN (LSTM)

Drawback: sequential generation is slow!



# PixelCNN *[van der Oord et al. 2016]*

Still generate image pixels starting from corner

Dependency on previous pixels now modeled using a CNN over context region

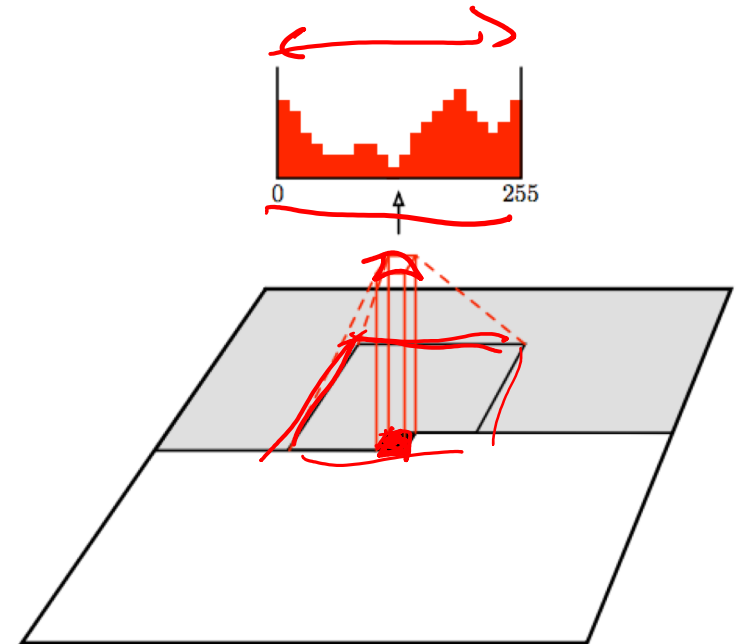
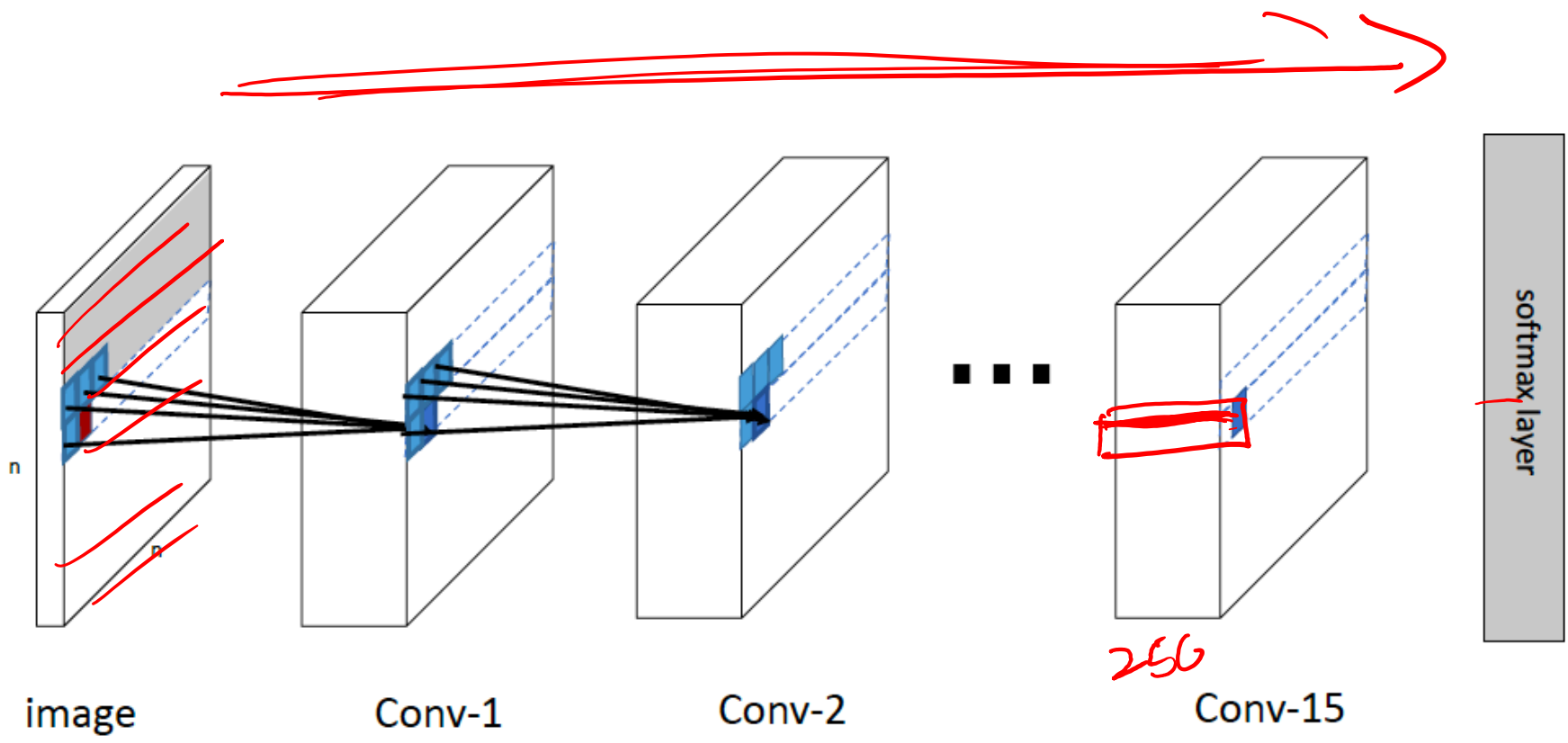


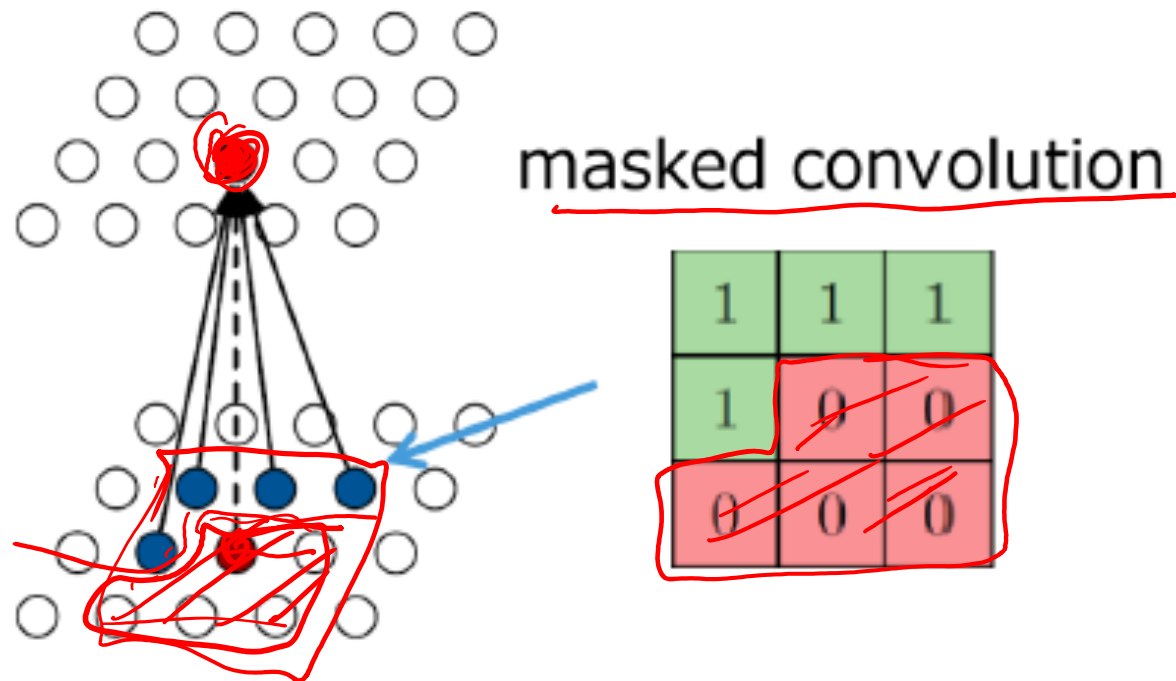
Figure copyright van der Oord et al., 2016. Reproduced with permission.





# Masked Convolutions

- Apply masks so that a pixel does not see “future” pixels



# PixelCNN [van der Oord et al. 2016]

Still generate image pixels starting from corner

Dependency on previous pixels now modeled using a CNN over context region

Training: maximize likelihood of training images

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

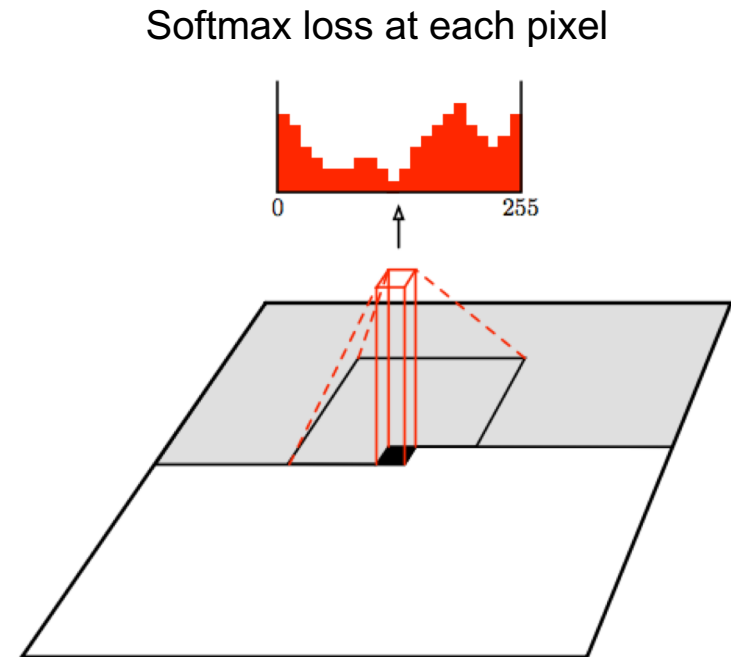


Figure copyright van der Oord et al., 2016. Reproduced with permission.

# PixelCNN *[van der Oord et al. 2016]*

Still generate image pixels starting from corner

Dependency on previous pixels now modeled using a CNN over context region

Training is faster than PixelRNN  
(can parallelize convolutions since context region values known from training images)

Generation must still proceed sequentially  
=> still slow

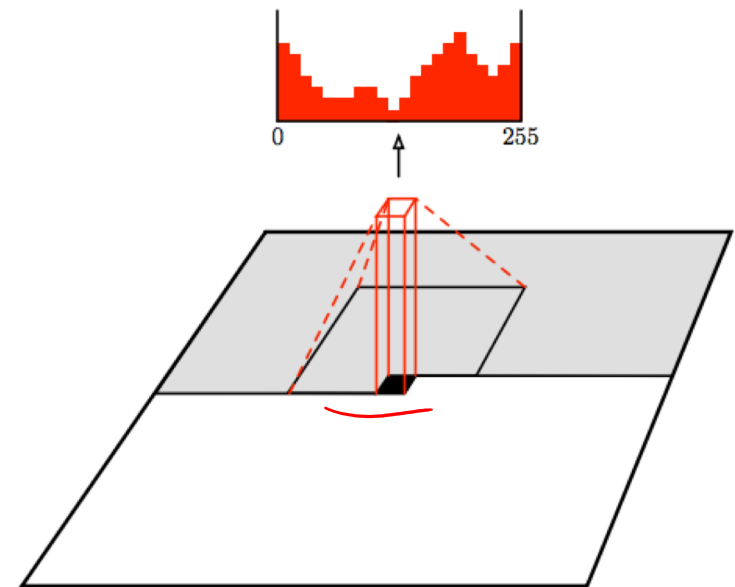
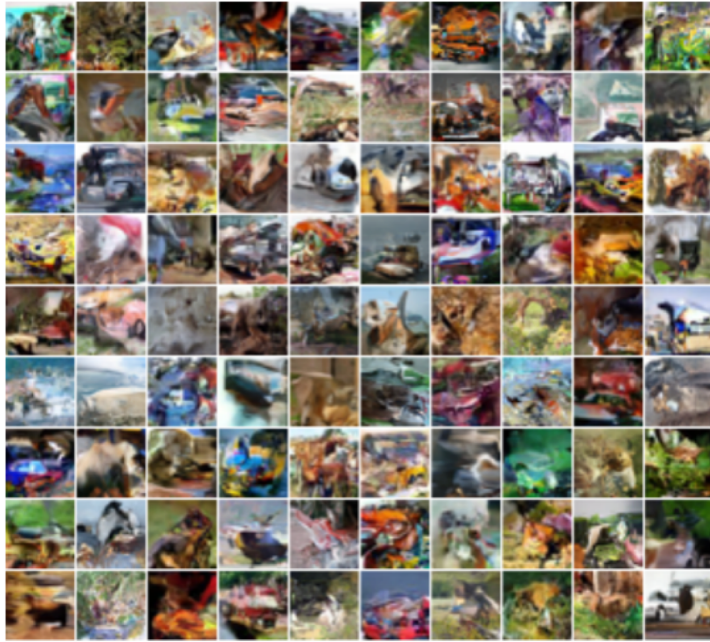
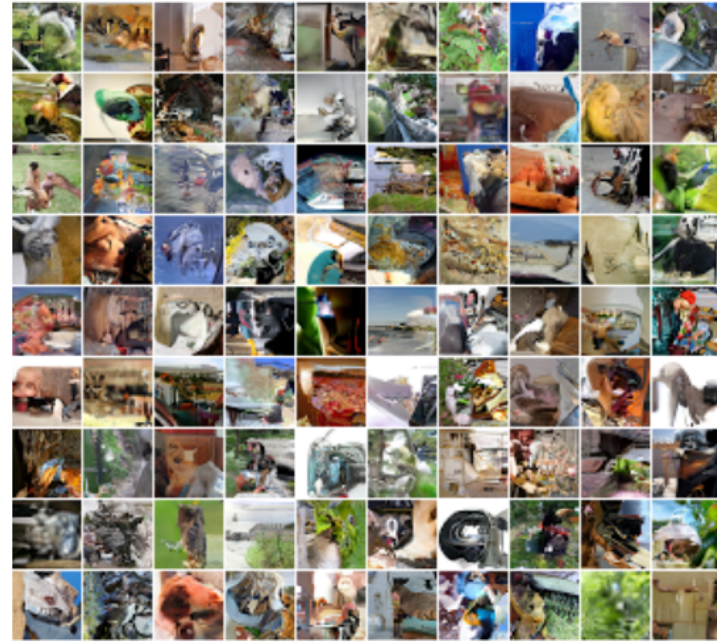


Figure copyright van der Oord et al., 2016. Reproduced with permission.

# Generation Samples



32x32 CIFAR-10



32x32 ImageNet

Figures copyright Aaron van der Oord et al., 2016. Reproduced with permission.

# Image Completion

$p(\tilde{x})$

occluded

completions

original

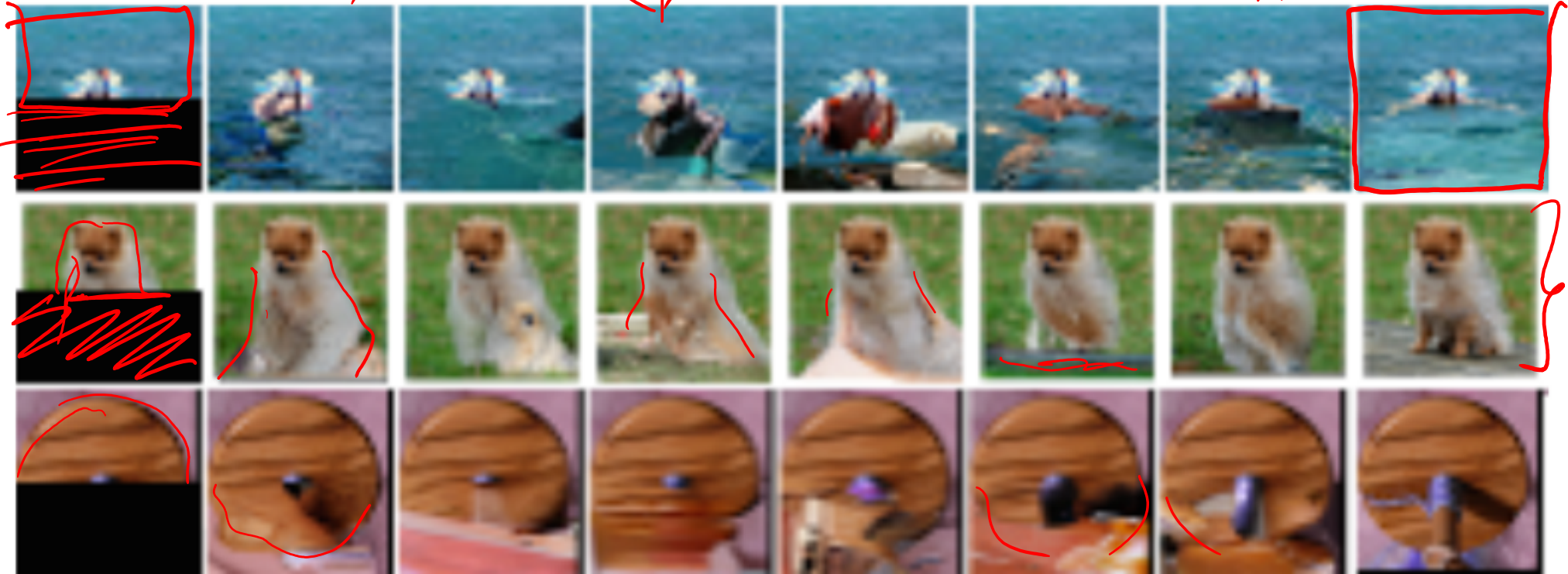


Figure 1. Image completions sampled from a PixelRNN.

# Results from generating sounds

- <https://deepmind.com/blog/wavenet-generative-model-raw-audio/>

# PixelRNN and PixelCNN

## Pros:

- Can explicitly compute likelihood  $p(x)$
- Explicit likelihood of training data gives good evaluation metric
- Good samples

## Con:

- Sequential generation  
=> slow

## Improving PixelCNN performance

- Gated convolutional layers
- Short-cut connections
- Discretized logistic loss
- Multi-scale
- Training tricks
- Etc...

## See

- Van der Oord et al. NIPS 2016
- Salimans et al. 2017  
(PixelCNN++)