

Deep Generative Models

Part1: Introduction



KAUST

King Abdullah University of
Science and Technology

Mohamed Elhoseiny

KAUST

mohamed.elhoseiny@kaust.edu.sa

Yen
na Yeung





What to Expect ?

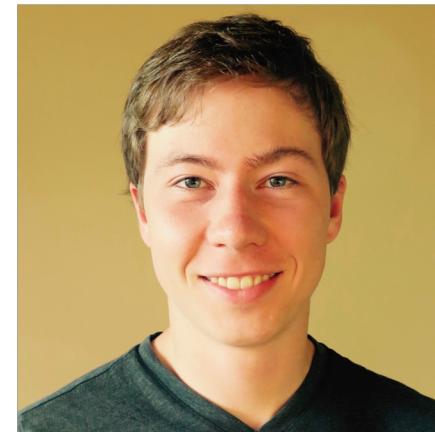
- **Part 1:** Learn what are generative models and about key modeling approaches including Generative Adversarial Networks and Variational Auto Encoders.
- **Part 2:** Learn more about Deep Generative Models capabilities to solve problems like recognizing unseen objects and creating unseen likable images.



Credit for Part 1



Fei-Fei Li



Justin Johnson



Serena Yeung

Part1: Introduction to Deep Generative Models



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.



→ Cat

Classification

[This image is CC0 public domain](#)



Supervised vs Unsupervised Learning

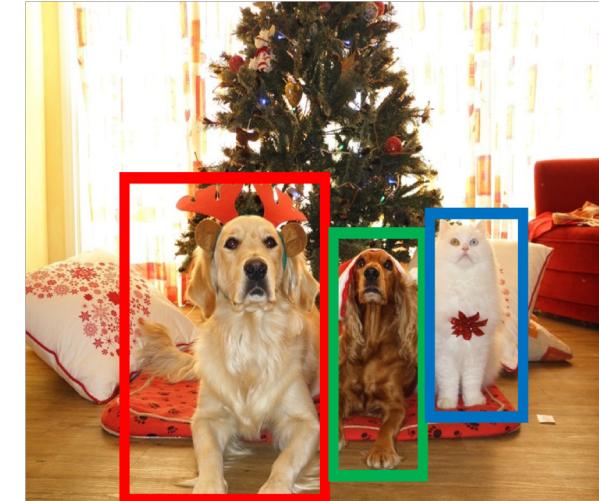
Supervised Learning

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

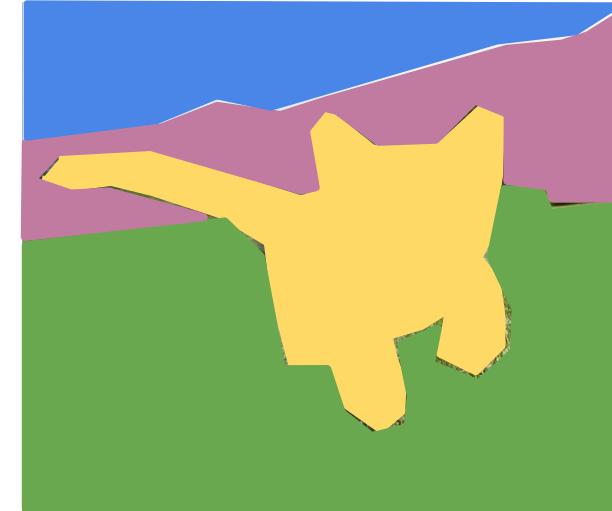
Supervised Learning

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

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

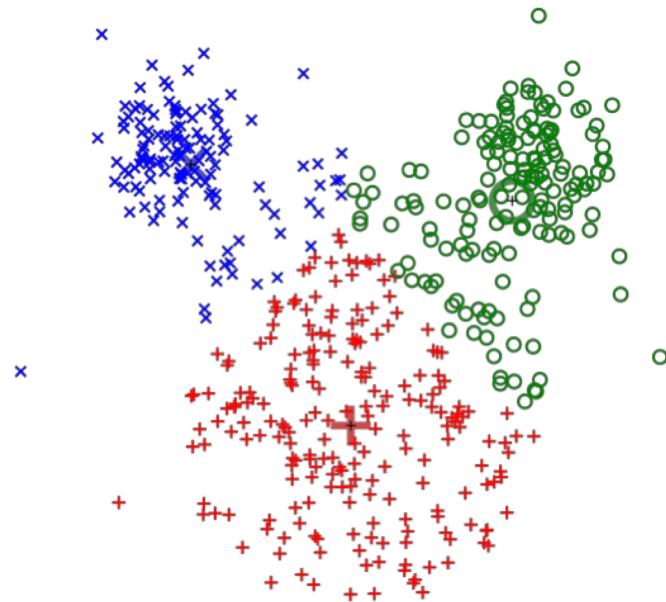
Unsupervised Learning

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

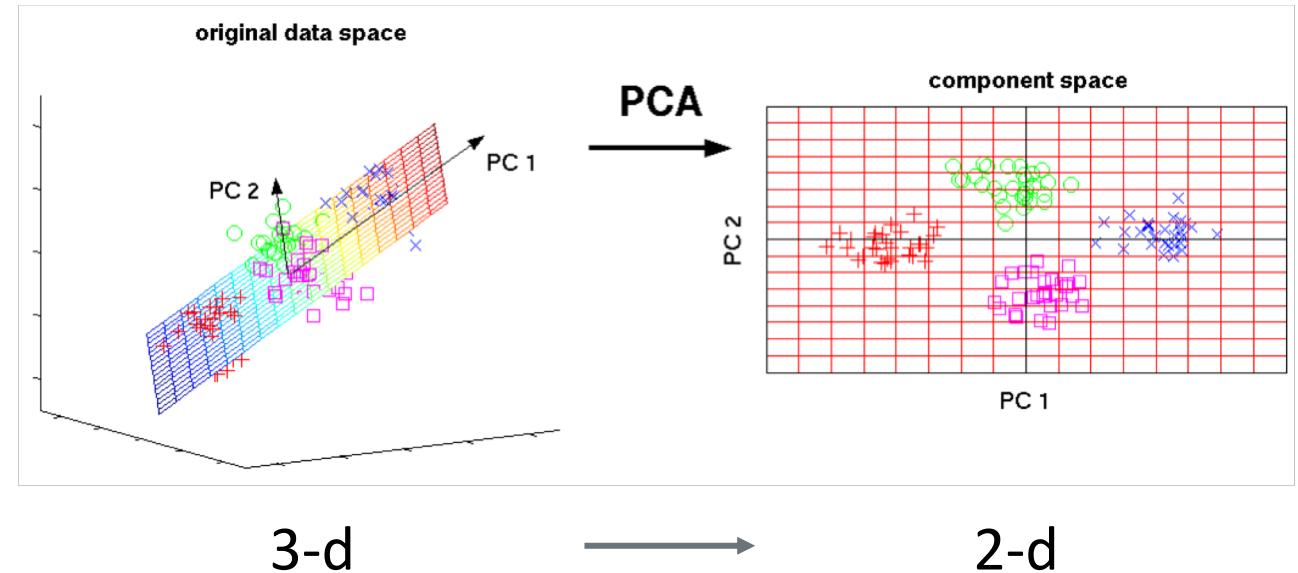
Unsupervised Learning

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



Supervised vs Unsupervised Learning

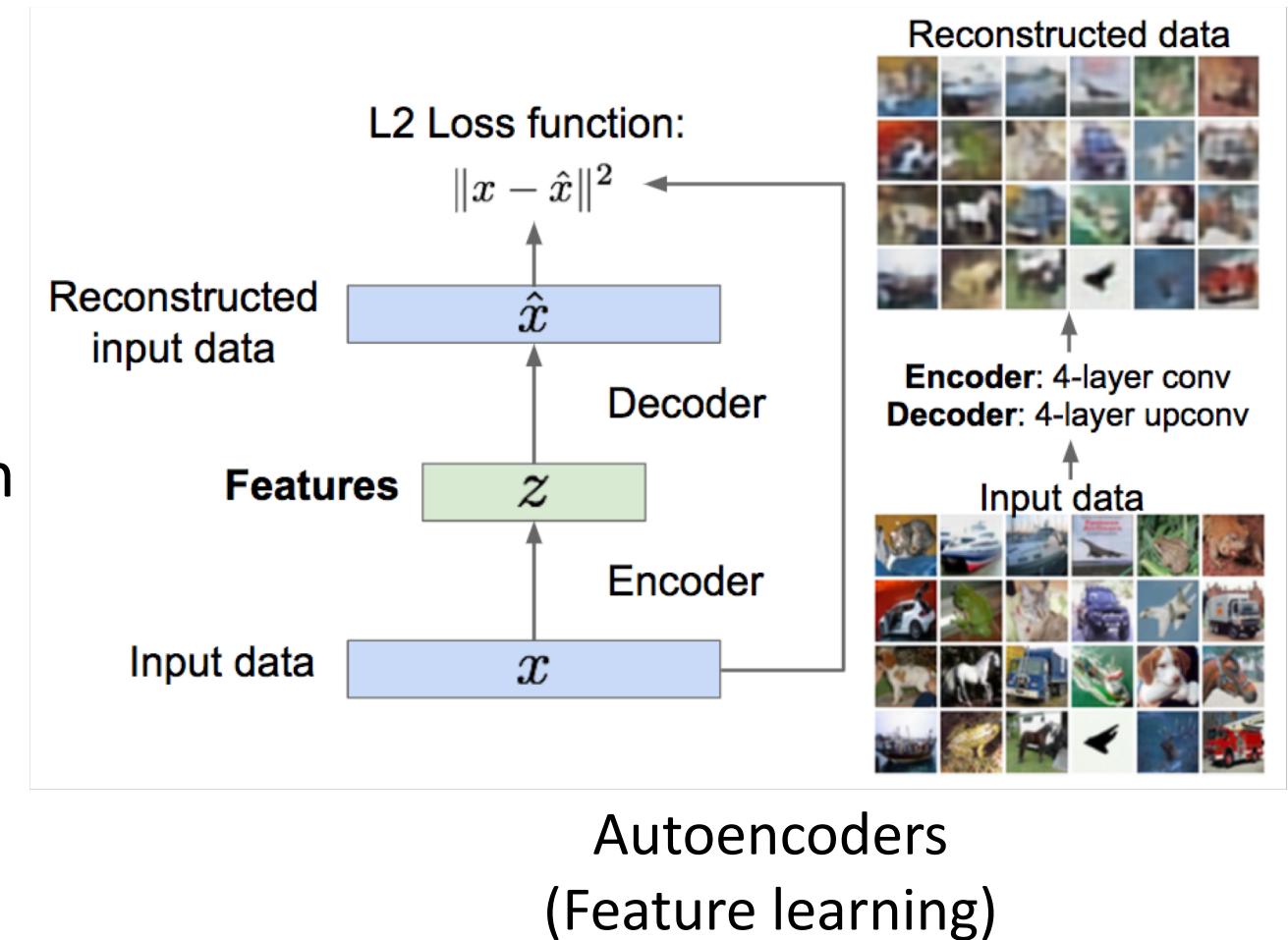
Unsupervised Learning

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

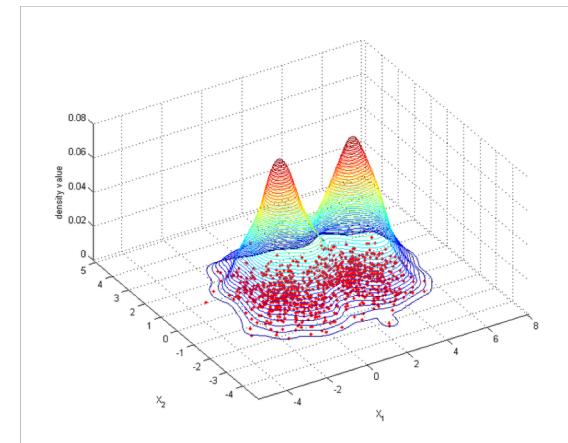
Just data, no labels!

Goal: Learn some underlying hidden
structure of the data

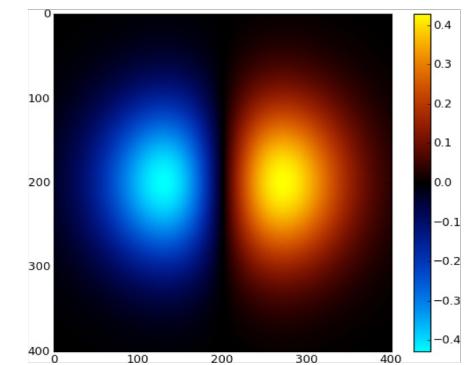
Examples: Clustering,
dimensionality reduction, feature
learning, density estimation, etc.



1-d density estimation



2-d density estimation



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

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

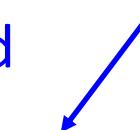
Data: x

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



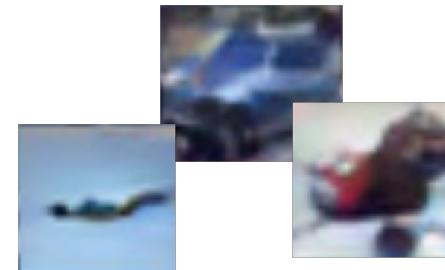


Generative Models

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



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

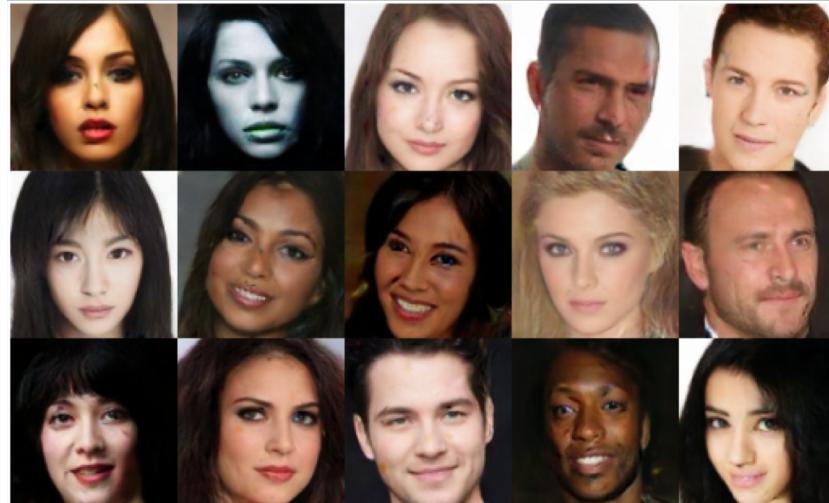
Several flavors:

- 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) [Phillip Isola et al. 2017](#). Reproduced with authors permission (3) [BAIR Blog](#).



Taxonomy of Generative Models

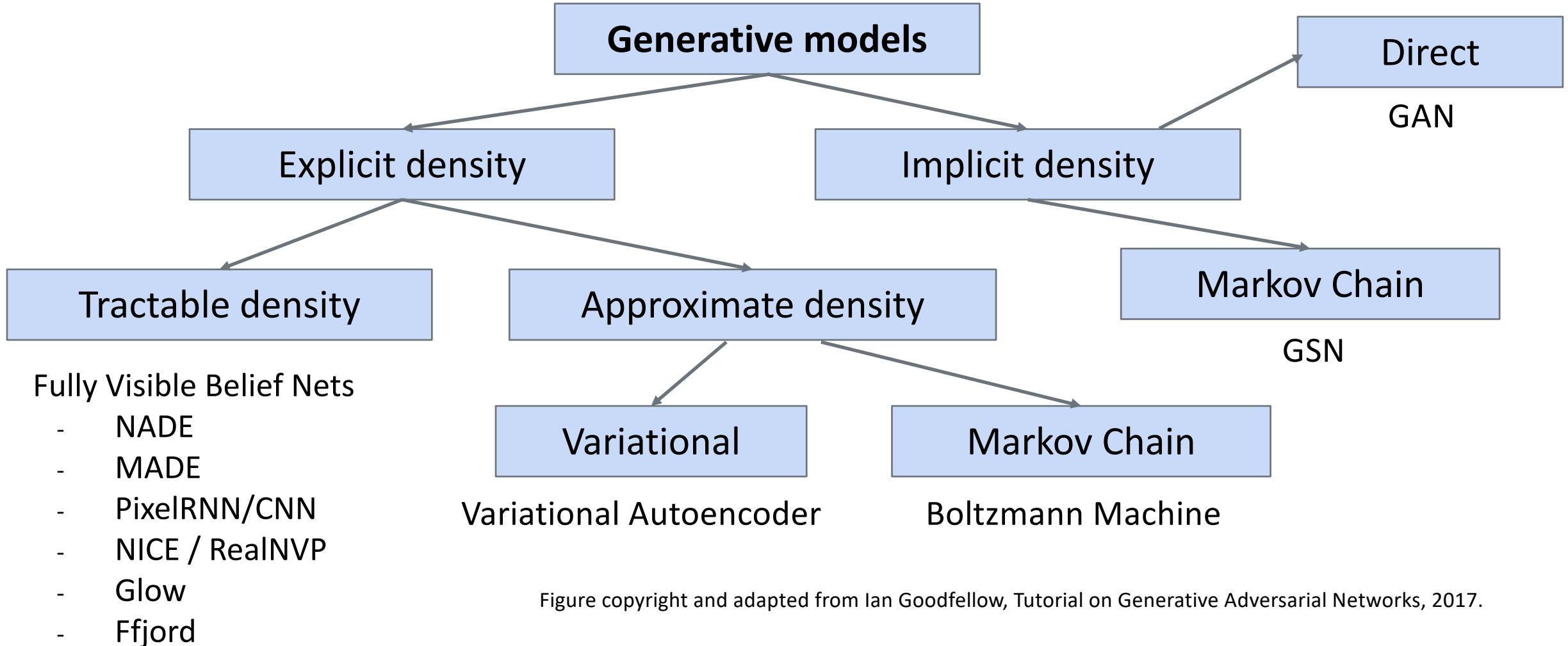


Figure copyright and adapted from Ian Goodfellow, Tutorial on Generative Adversarial Networks, 2017.



Taxonomy of Generative Models

Today: discuss 3 most popular types of generative models today

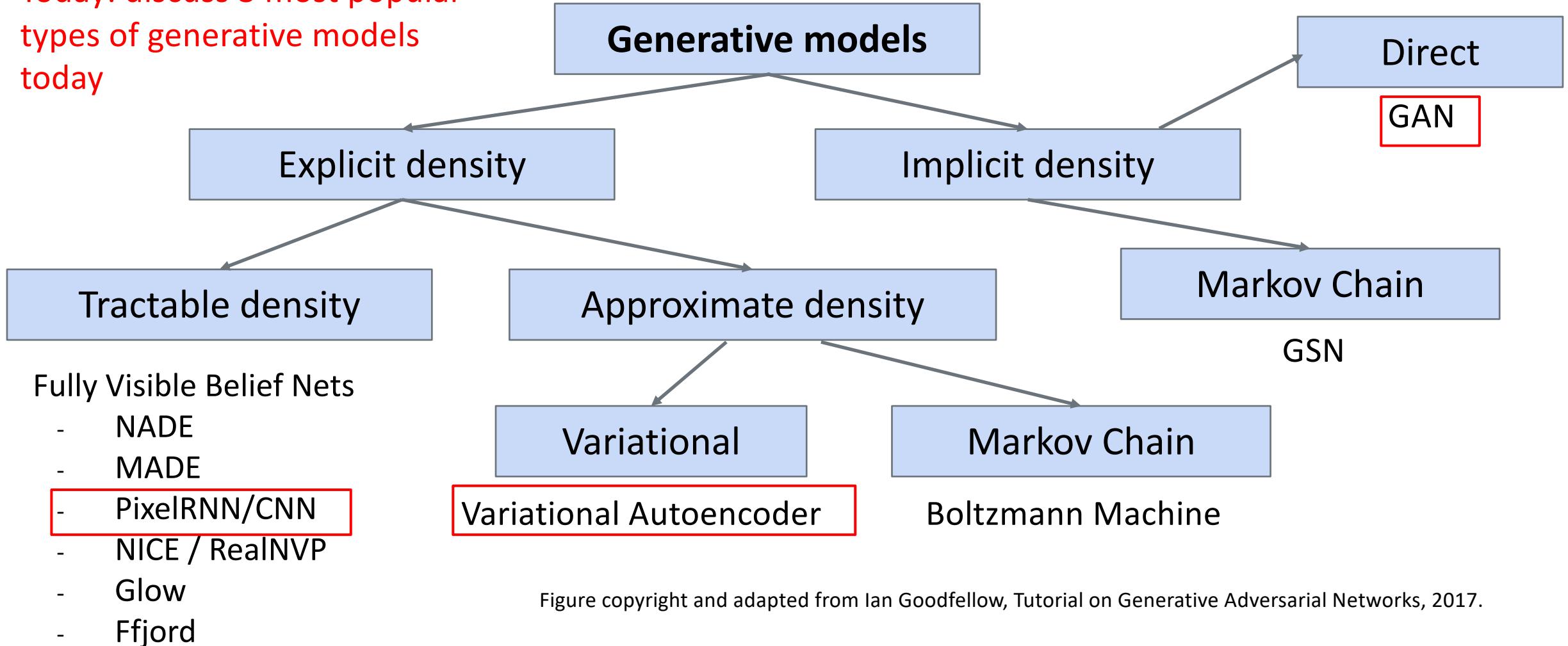


Figure copyright and adapted from Ian Goodfellow, Tutorial on Generative Adversarial Networks, 2017.



PixelRNN and PixelCNN



Fully visible belief network

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



Fully visible belief network

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 visible belief network

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!

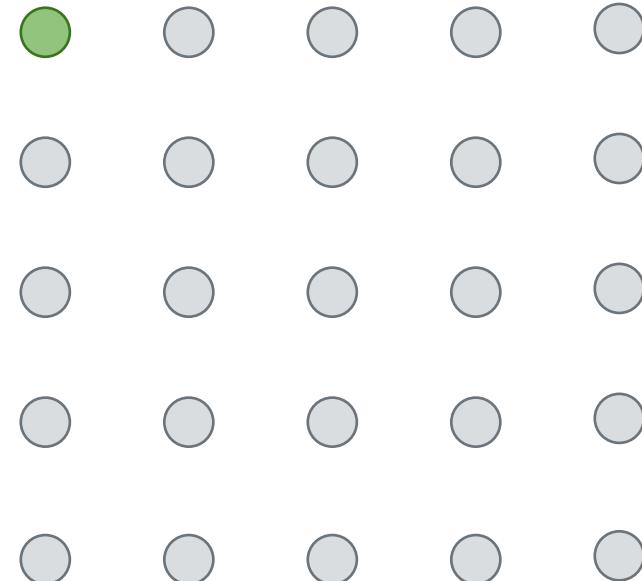
Will need to define ordering of “previous pixels”



PixelRNN

[van der Oord et al. 2016]

Generate image pixels starting from corner



Dependency on previous pixels modeled using
an RNN (LSTM)

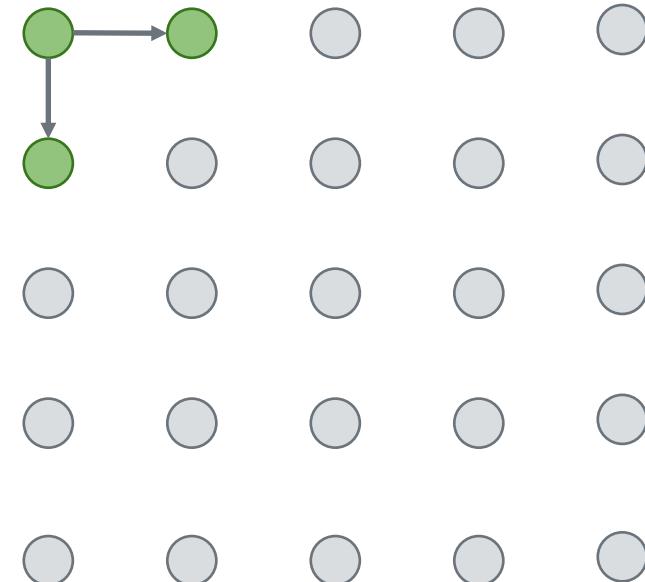


PixelRNN

[van der Oord et al. 2016]

Generate image pixels starting from corner

Dependency on previous pixels modeled using
an RNN (LSTM)



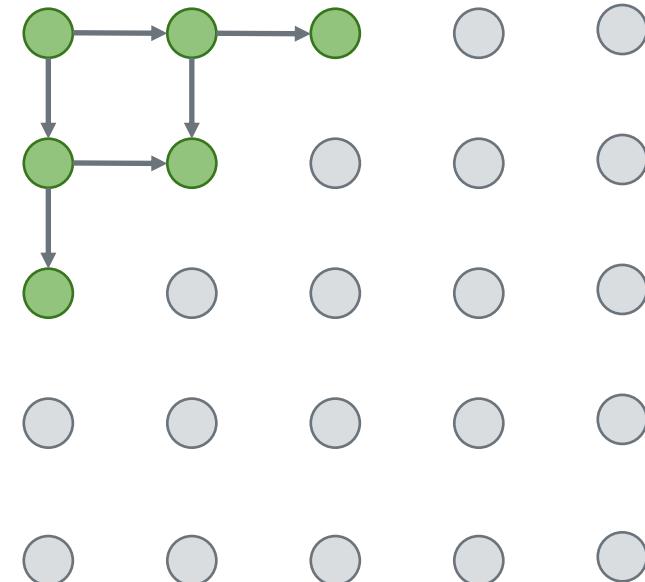


PixelRNN

[van der Oord et al. 2016]

Generate image pixels starting from corner

Dependency on previous pixels modeled using
an RNN (LSTM)





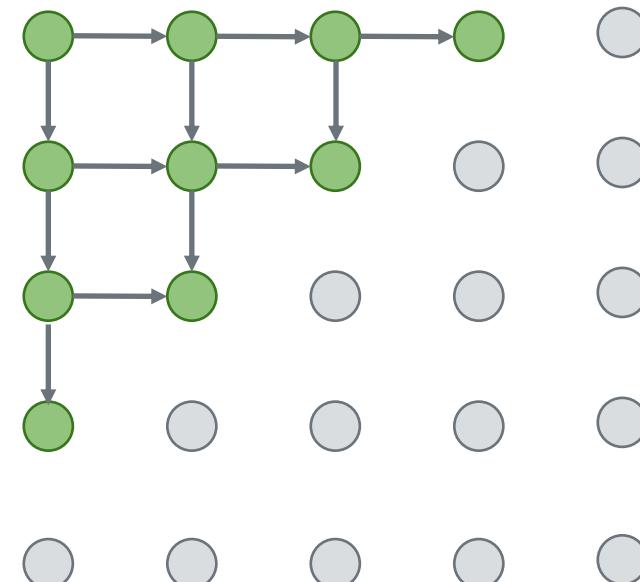
PixelRNN

[van der Oord et al. 2016]

Generate image pixels starting from corner

Dependency on previous pixels modeled using
an RNN (LSTM)

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

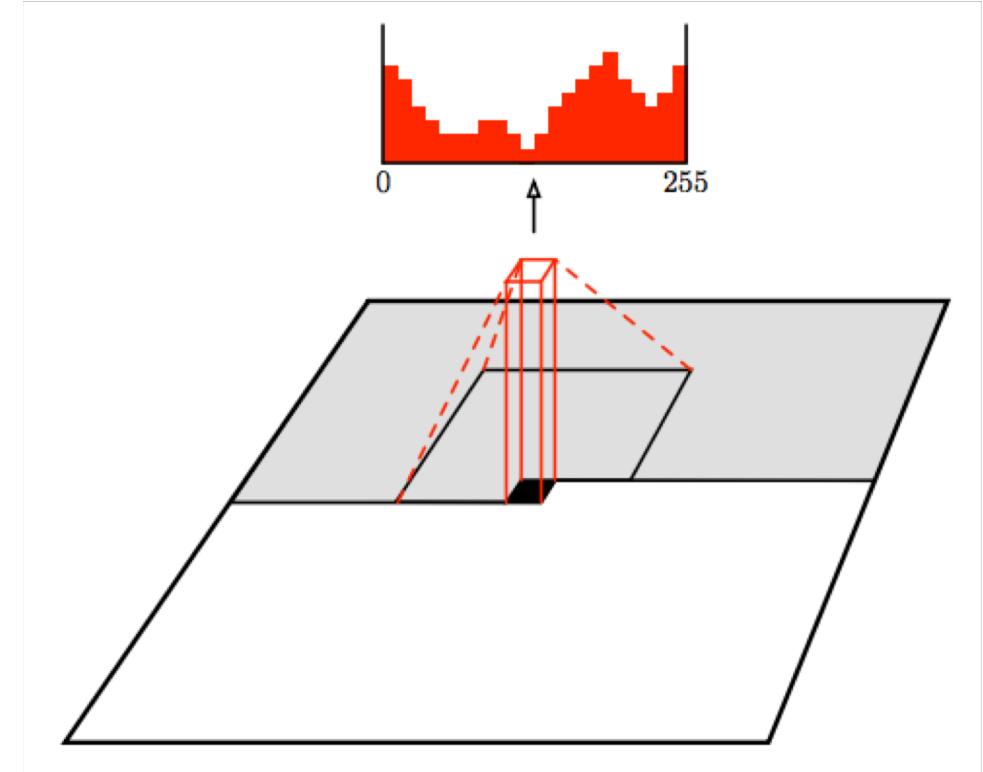


Figure copyright van der Oord et al., 2016. Reproduced with permission.



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

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

Softmax loss at each pixel

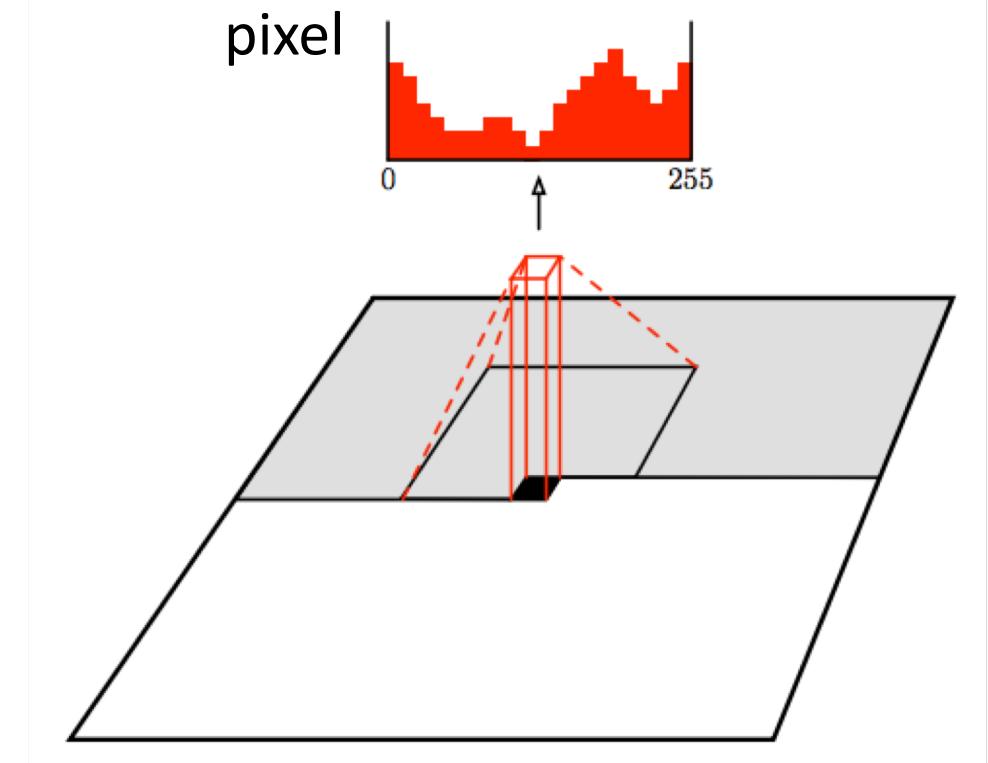


Figure copyright van der Oord et al., 2016. Reproduced with permission.



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

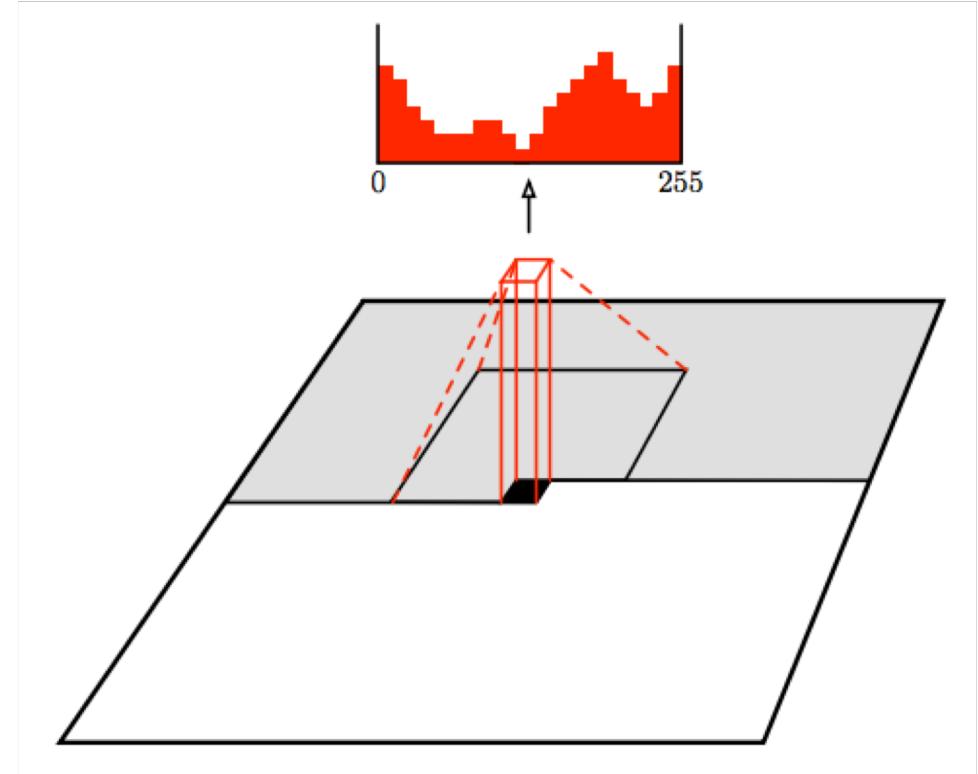
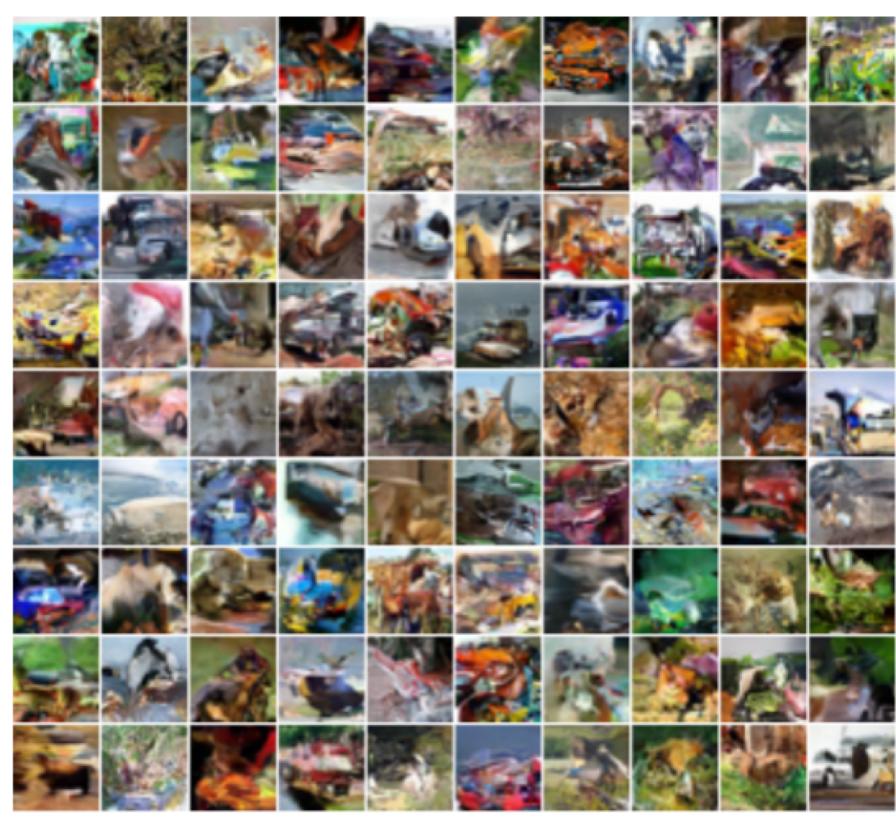


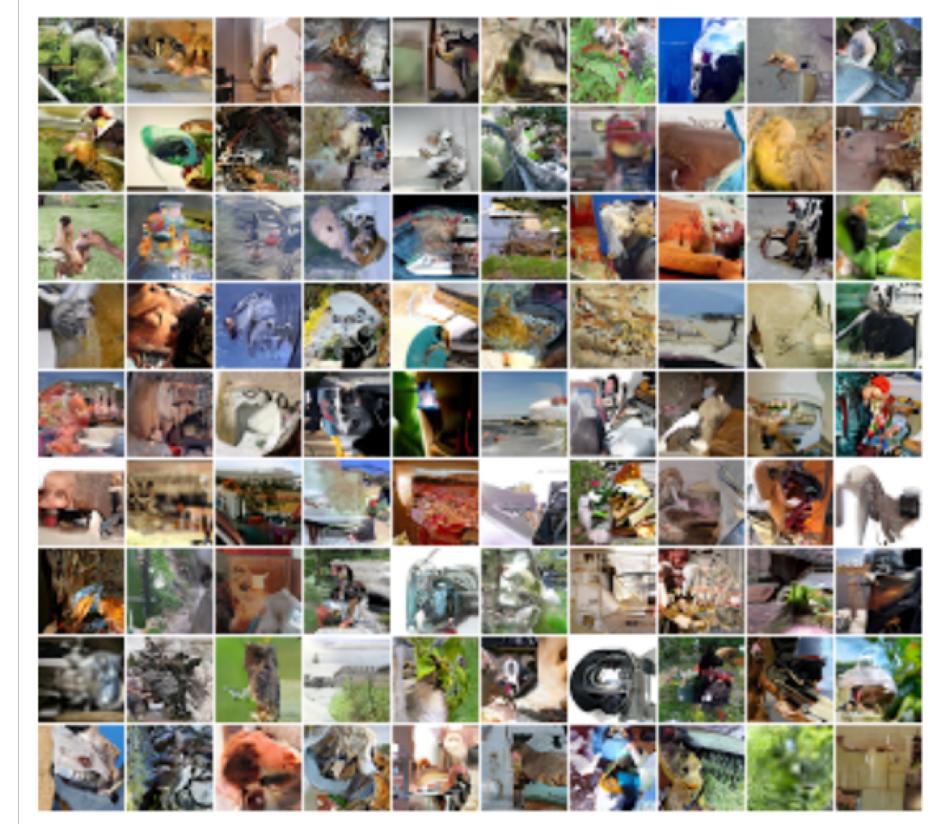
Figure copyright van der Oord et al., 2016. Reproduced with permission.



Generation Samples



32x32 CIFAR-10



32x32
ImageNet

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



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++)



Variational Autoencoders (VAE)



Variational Autoencoders

So far...

PixelCNNs define tractable density function, optimize likelihood of training data:

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



Variational Autoencoders

So far...

PixelCNNs define tractable density function, optimize likelihood of training data:

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

VAEs define intractable density function with latent \mathbf{z} :

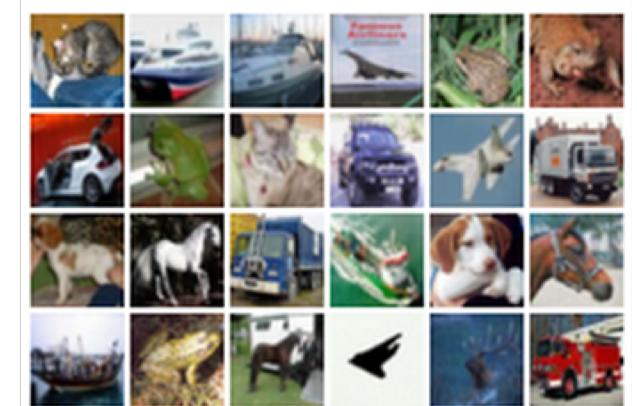
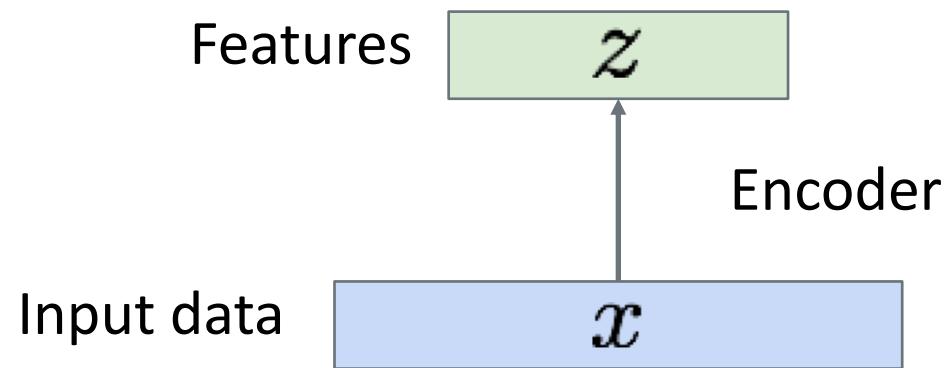
$$p_{\theta}(x) = \int p_{\theta}(z)p_{\theta}(x|z)dz$$

Cannot optimize directly, derive and optimize lower bound on likelihood instead



Some background first: Autoencoders

Unsupervised approach for learning a lower-dimensional feature representation from unlabeled training data

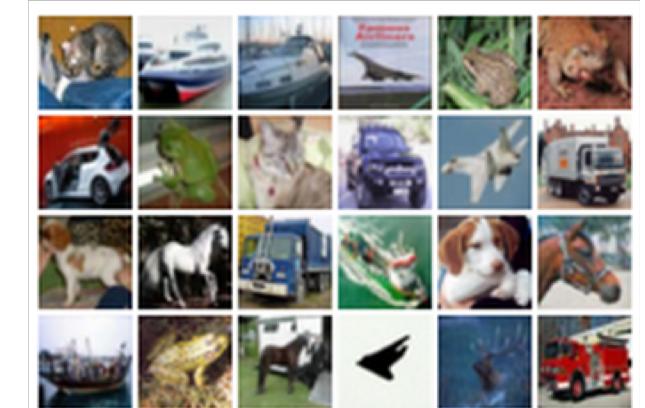
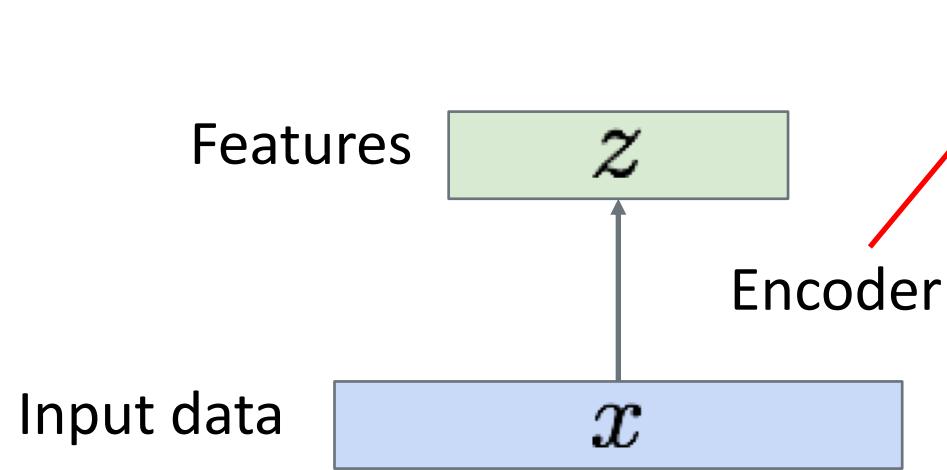




Some background first: Autoencoders

Unsupervised approach for learning a lower-dimensional feature representation from unlabeled training data

Originally: Linear + nonlinearity
(sigmoid)
Later: Deep, fully-connected
Later: ReLU CNN



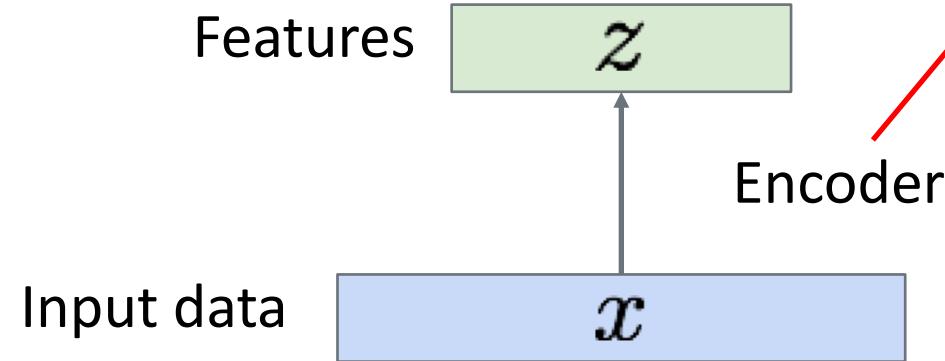


Some background first: Autoencoders

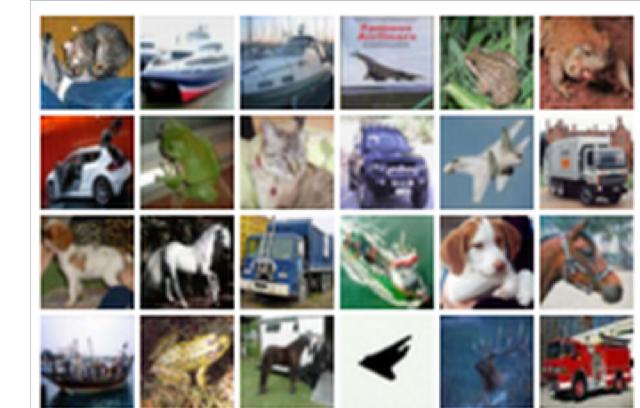
Unsupervised approach for learning a lower-dimensional feature representation from unlabeled training data

z usually smaller than x
(dimensionality reduction)

Q: Why dimensionality reduction?



Originally: Linear + nonlinearity (sigmoid)
Later: Deep, fully-connected
Later: ReLU CNN





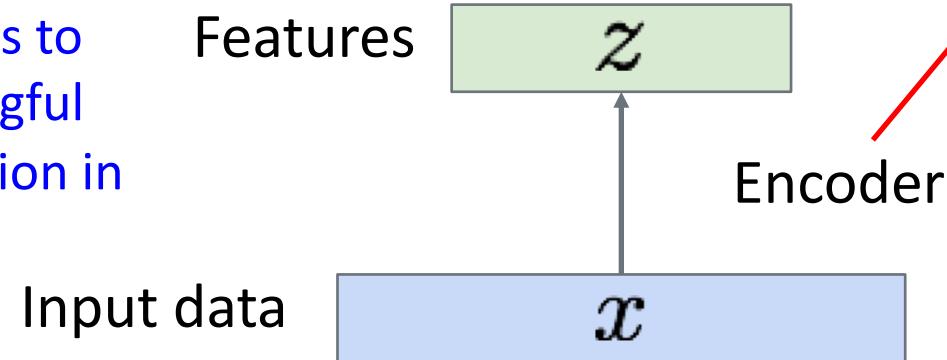
Some background first: Autoencoders

Unsupervised approach for learning a lower-dimensional feature representation from unlabeled training data

z usually smaller than x
(dimensionality reduction)

Q: Why dimensionality reduction?

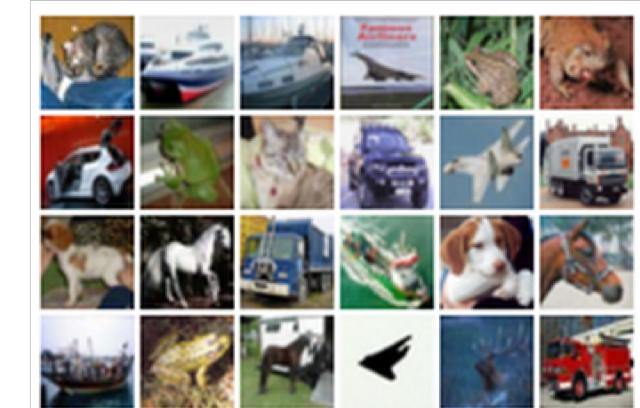
A: Want features to capture meaningful factors of variation in data



Originally: Linear + nonlinearity
(sigmoid)

Later: Deep, fully-connected

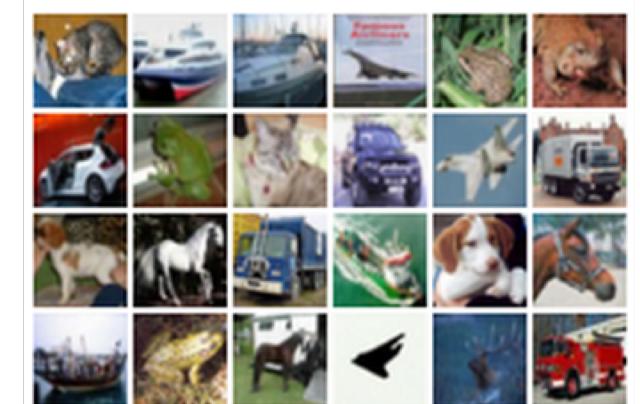
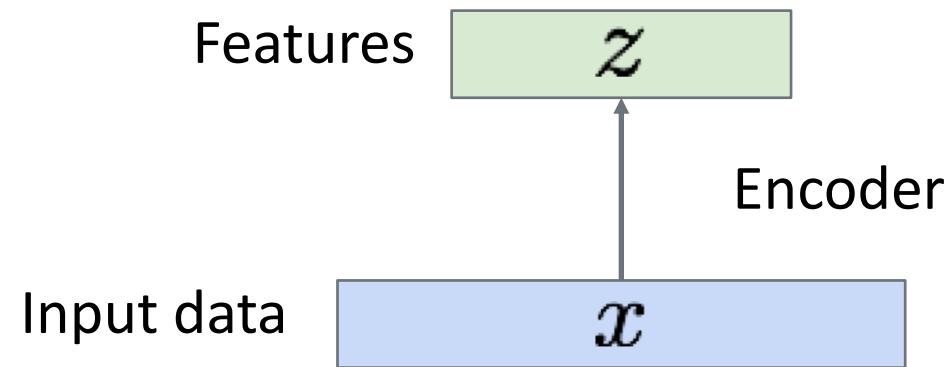
Later: ReLU CNN





Some background first: Autoencoders

How to learn this feature representation?



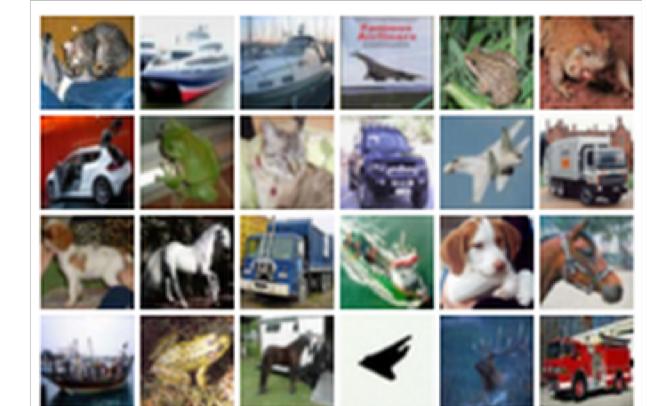
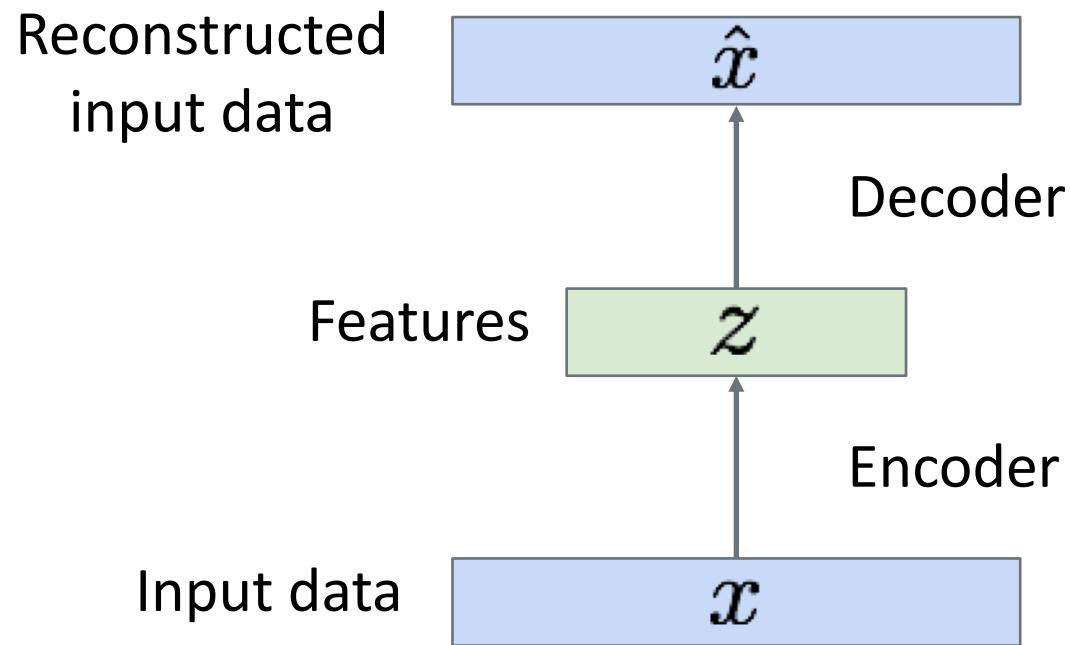


Some background first: Autoencoders

How to learn this feature representation?

Train such that features can be used to reconstruct original data

“Autoencoding” - encoding itself



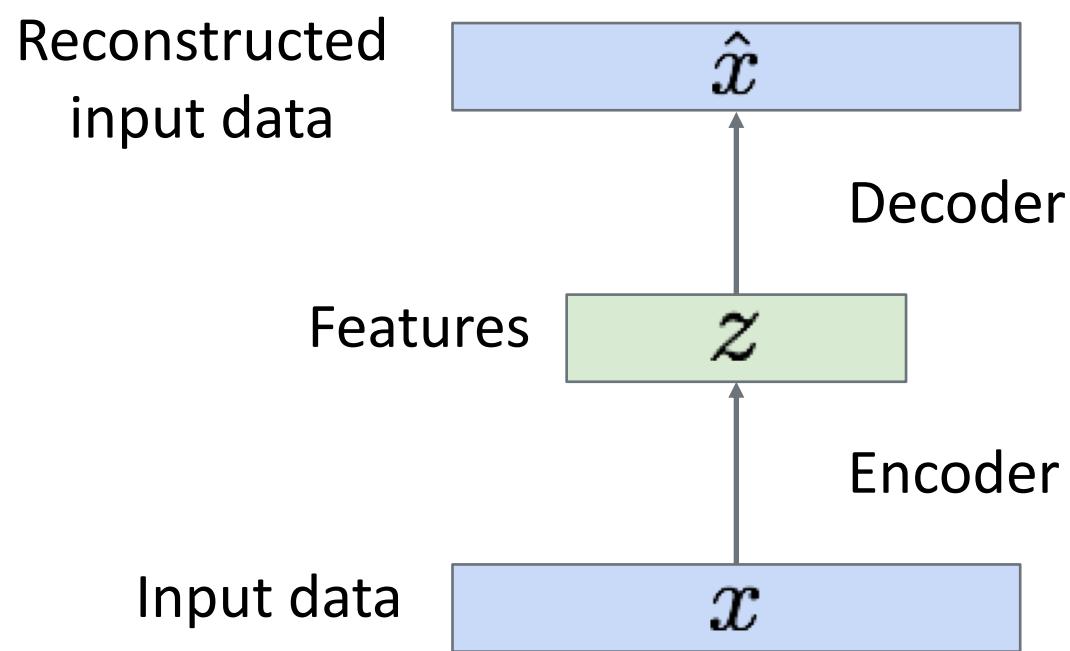


Some background first: Autoencoders

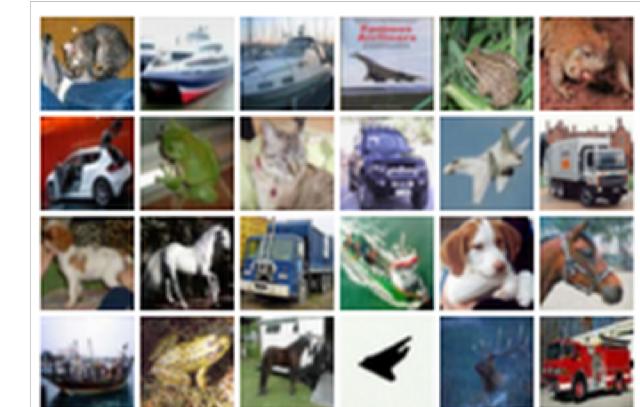
How to learn this feature representation?

Train such that features can be used to reconstruct original data

“Autoencoding” - encoding itself



Originally: Linear +
nonlinearity (sigmoid)
Later: Deep, fully-connected
Later: ReLU CNN (upconv)



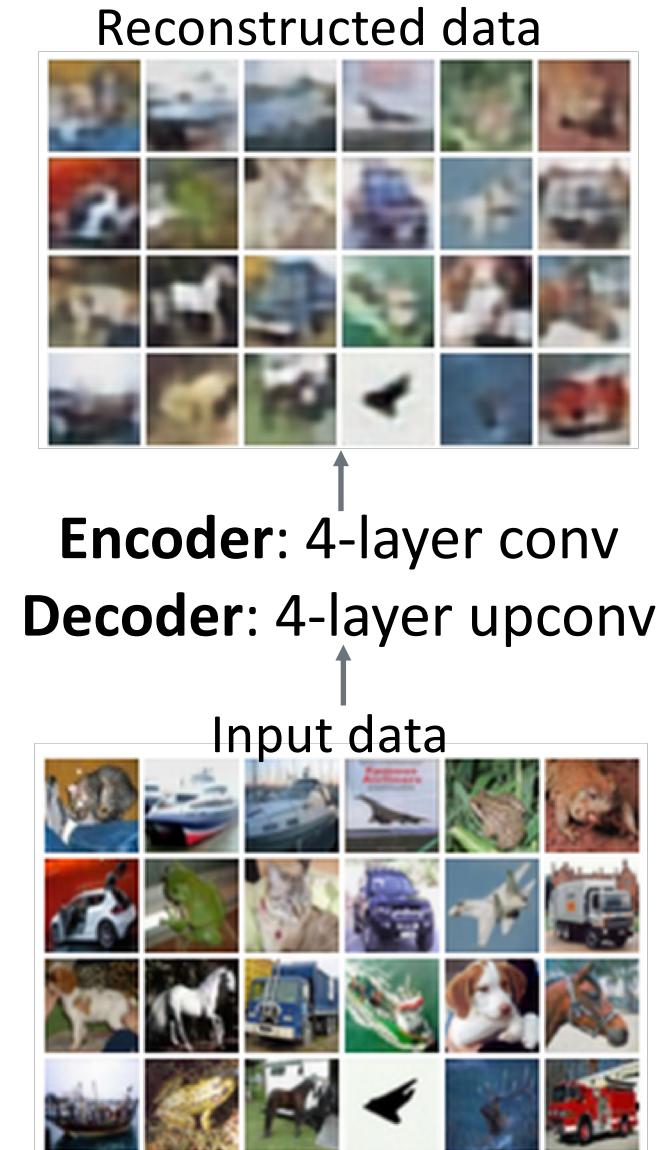
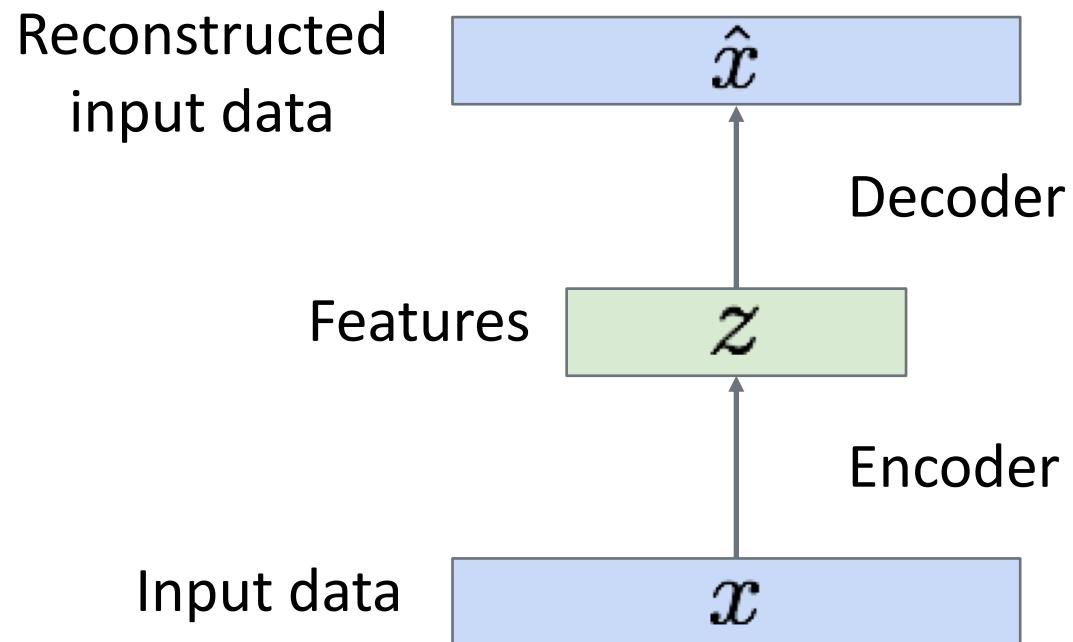


Some background first: Autoencoders

How to learn this feature representation?

Train such that features can be used to reconstruct original data

“Autoencoding” - encoding itself





Some background first: Autoencoders

Train such that features can
be used to reconstruct
original data

Reconstructed
input data

Features

Input data

L2 Loss function:

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

\hat{x}

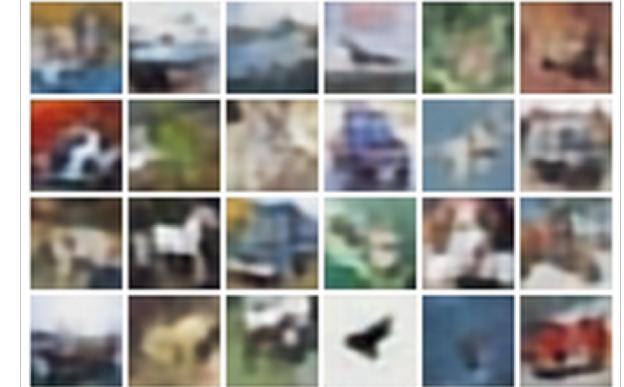
Decoder

z

Encoder

x

Reconstructed data



Encoder: 4-layer conv

Decoder: 4-layer upconv

Input data





Some background first: Autoencoders

Train such that features can
be used to reconstruct
original data

Reconstructed
input data

Features

Input data

L2 Loss function:

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

\hat{x}

Decoder

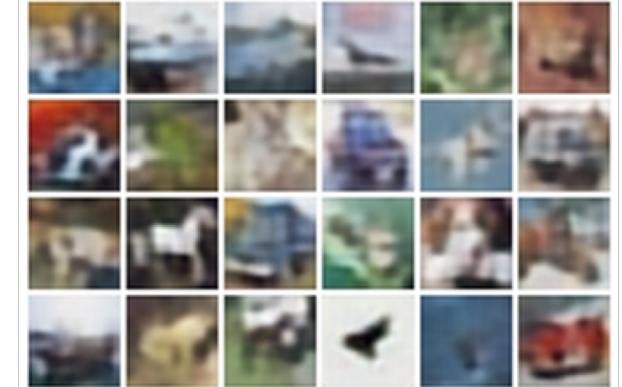
z

Encoder

x

Doesn't use labels!

Reconstructed data



Encoder: 4-layer conv

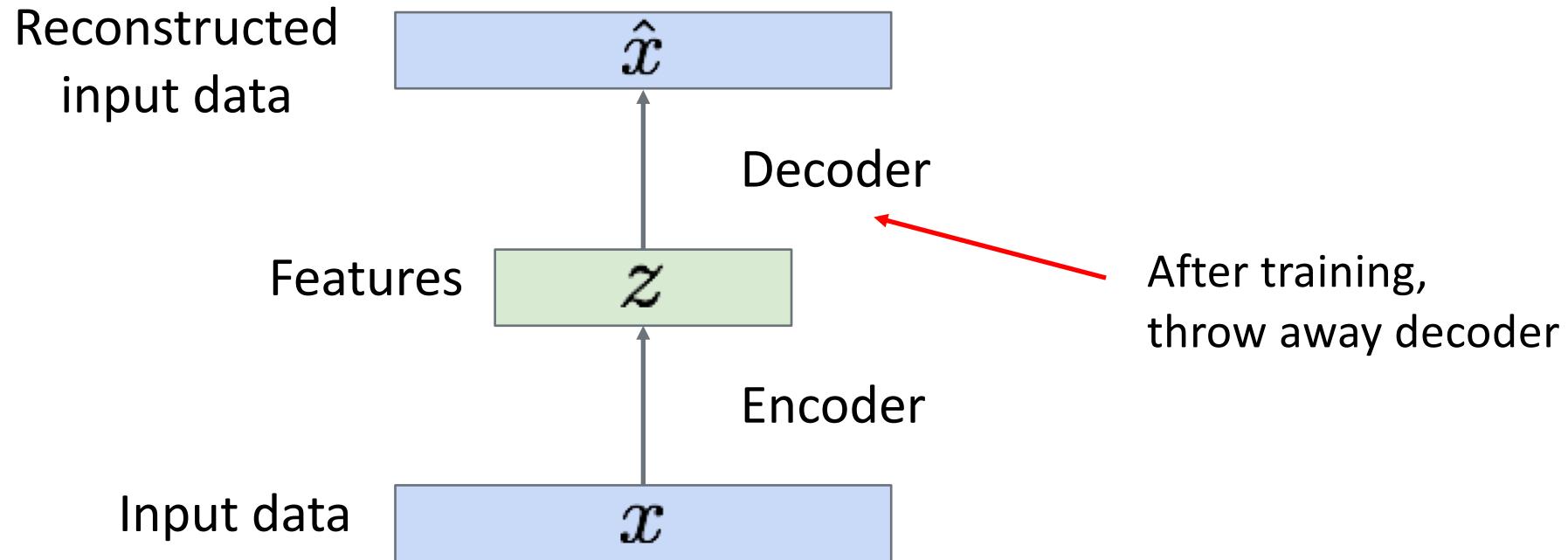
Decoder: 4-layer upconv

Input data





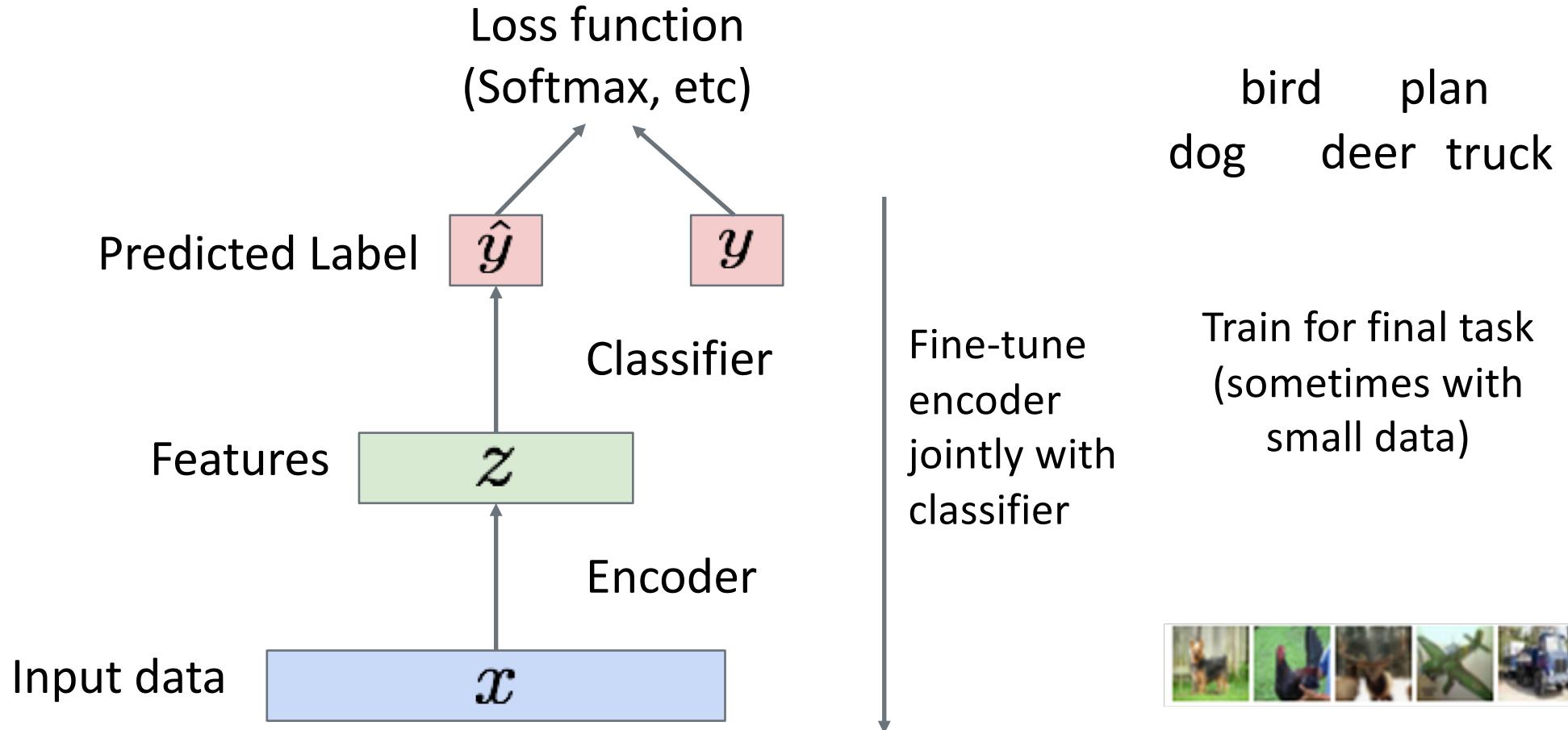
Some background first: Autoencoders





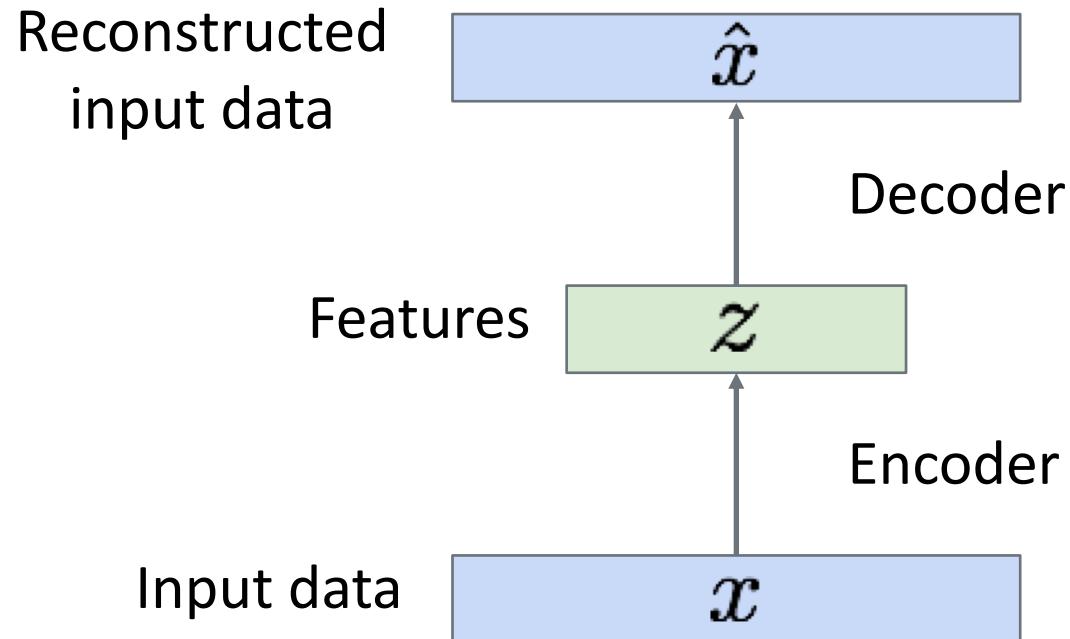
Some background first: Autoencoders

Encoder can be used
to initialize a
supervised model





Some background first: Autoencoders



Autoencoders can reconstruct data, and can learn features to initialize a supervised model

Features capture factors of variation in training data. Can we generate new images from an autoencoder?