

Generative Models

Generative Model

What you will learn in this course

Low-Level Vision

Image Processing:

Photometric Image formation
Point Operations
Histogram Equalization
Linear & Non-Linear filters
Convolution & Cross correlation

Feature Detection:

Edge detection
Corner detection
Line detection

Feature Description:

SIFT
Feature Matching
Image Stitching
Panorama



Geometric Vision

Transformation and Camera Mode:

Geometric Image Formation
Pinhole Camera Model

Camera Calibration:

Estimating intrinsics
Estimating extrinsics



Visual Understanding

Object Detection & Tracking:

Single-stage detectors
Two-stage detectors
YOLO
Object tracking

Image Segmentation:

Traditional Image Segmentation
Learning-based Image Segmentation



Machine Learning for CV

Machine Learning for Classification:

k-NN
Linear Classifier

Deep Learning:

Optimization
Neural Network
CNN
CNN architectures:



Generative & Geometric Learning

Generative Models:

AE, VAE, GAN



3D Vision and 3D Deep Learning:

3D shape representations
Geometric deep learning overview
3D Deep Learning models : pointnet

What we will learn today

- ❑ Autoencoder
- ❑ Variational Autoencoder
- ❑ Generative Adversarial Network

Supervised vs. Unsupervised Learning

Supervised Learning

Data: (x, y)

x is data, y is label

Data Labeling is needed

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

$$y = f(x)$$

Classification: Input x



$$y = 8$$

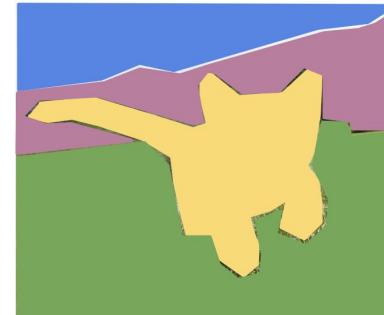
Object Detection



DOG, DOG, CAT

Examples: classification, regression, object detection, semantic segmentation, image segmentation 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 $x \rightarrow y$

$$y = f(x)$$

Examples: classification, regression, object detection, semantic segmentation, image segmentation etc.

Unsupervised Learning

Data: x

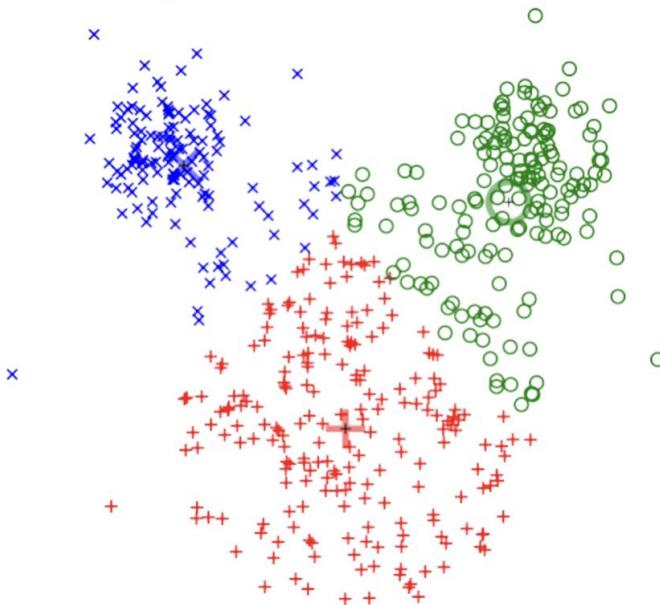
Just data, No labels!

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

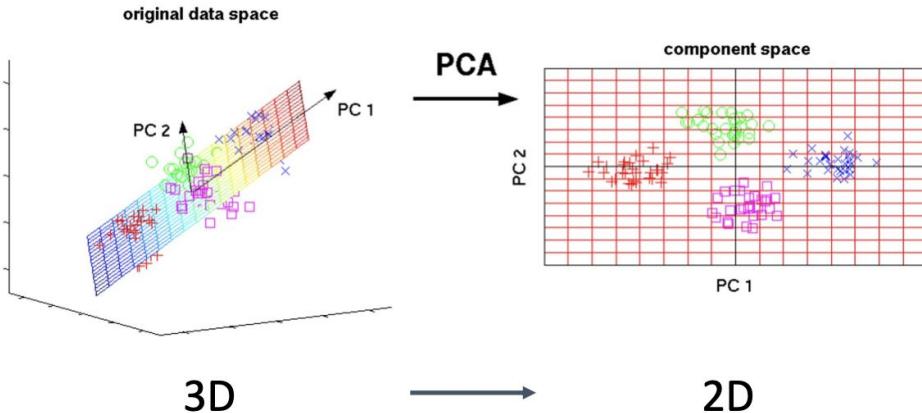
Just data, No labels!

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 Component Analysis- PCA)



Unsupervised Learning

Data: x

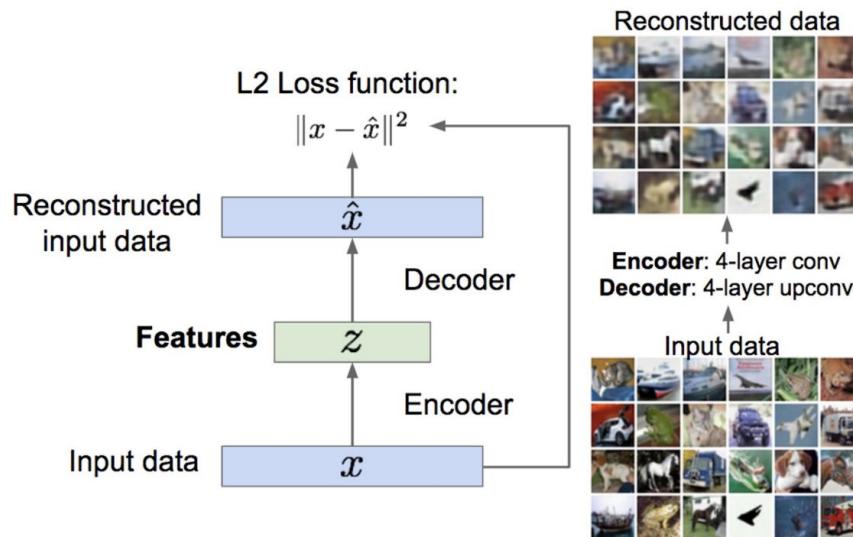
Just data, No labels!

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)



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 a probability distribution $P(x|y)$

Data: x



Label: $y = \text{Cat}$

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 **complete** for density

Discriminative vs. Generative Models

Discriminative Model:

Learn a probability distribution $P(y|x)$



$P(\text{cat} | \text{kitten})$

$P(\text{dog} | \text{kitten})$

Generative Model:

Learn a probability distribution $P(x)$

Conditional Generative

Model: Learn a probability distribution $P(x|y)$

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 **complete** for density

Discriminative vs. Generative Models

Discriminative Model:

Learn a probability distribution $P(y|x)$



$$P(\text{cat} | \text{kitten})$$

$$P(\text{dog} | \text{kitten})$$

Generative Model:

Learn a probability distribution $P(x)$



$$P(\text{dog} | \text{dog})$$

$$P(\text{cat} | \text{dog})$$

Conditional Generative

Model: Learn a probability distribution $P(x|y)$

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 $P(y|x)$



$P(\text{dog} | \text{dog})$

$P(\text{cat} | \text{dog})$



Generative Model:

Learn a probability distribution $P(x)$



$P(\text{cat} | \text{monkey})$

$P(\text{dog} | \text{monkey})$



Conditional Generative

Model: Learn a probability distribution $P(x|y)$

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

Discriminative vs. Generative Models

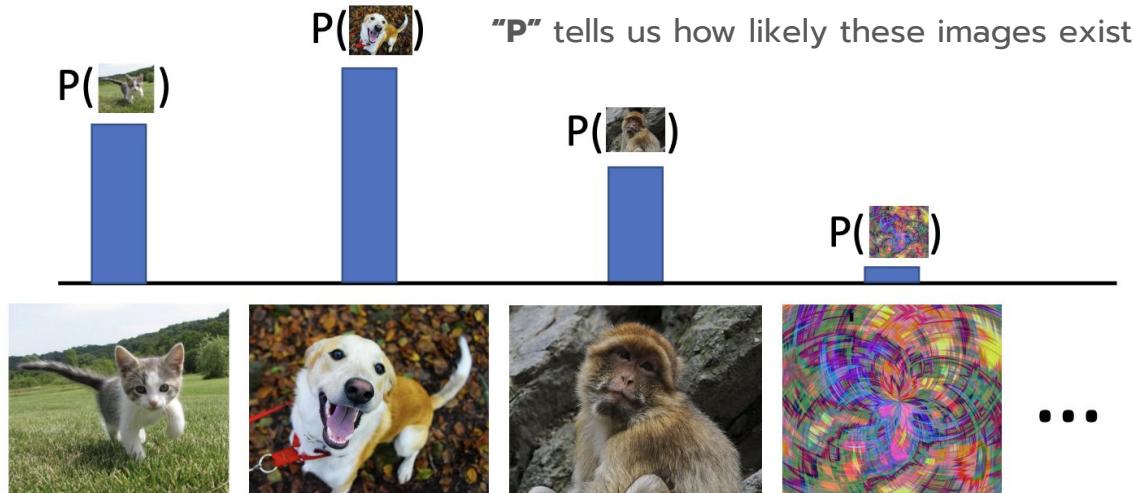
Discriminative Model:

Learn a probability distribution $P(y|x)$

Generative Model:

Learn a probability distribution $P(x)$

Conditional Generative Model: Learn a probability distribution $P(x|y)$



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

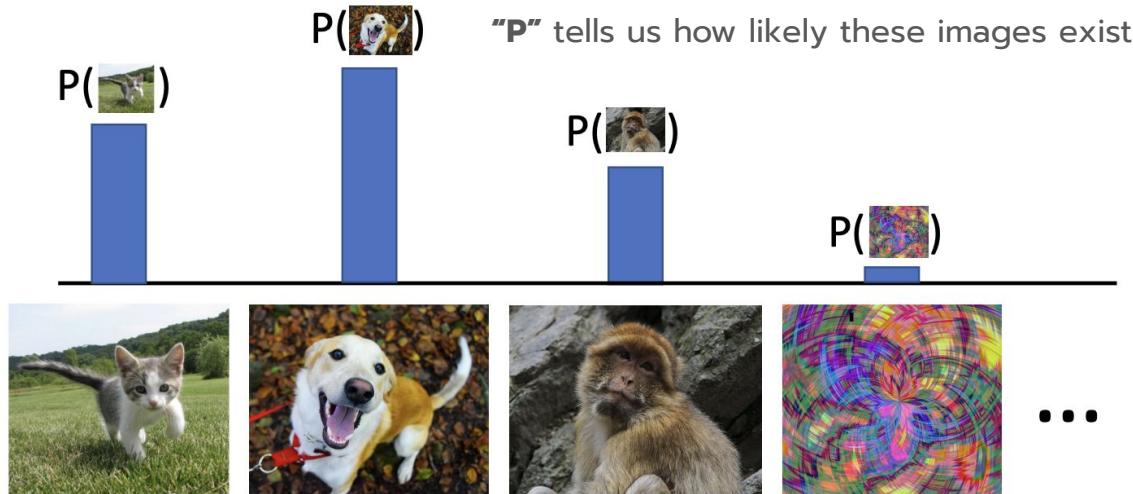
Discriminative Model:

Learn a probability distribution $P(y|x)$

Generative Model:

Learn a probability distribution $P(x)$

Conditional Generative Model: Learn a probability distribution $P(x|y)$



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

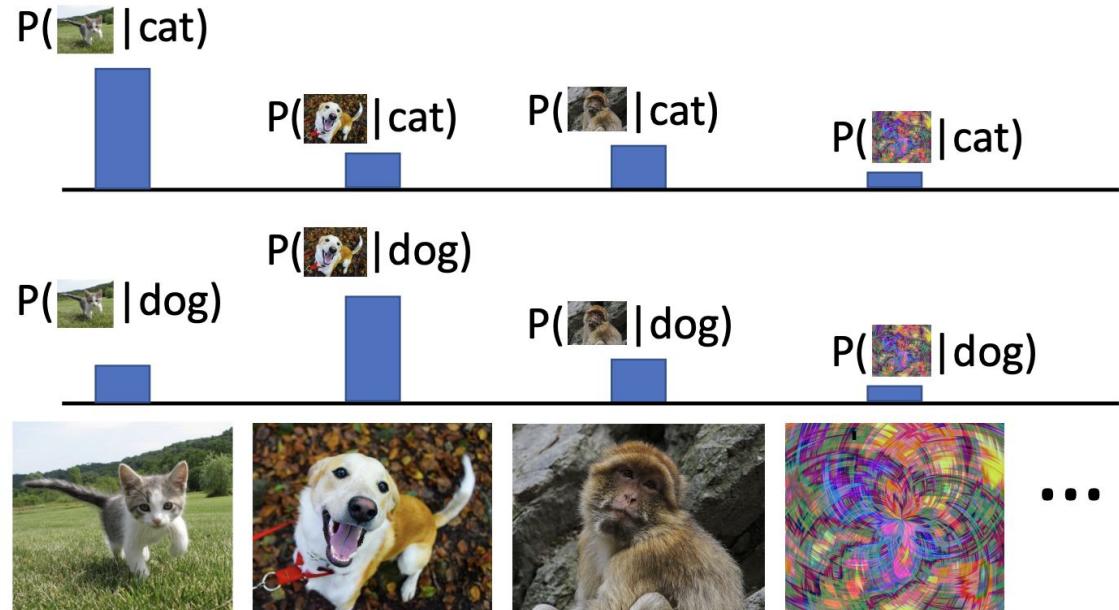
Discriminative Model:

Learn a probability distribution $P(y|x)$

Generative Model:

Learn a probability distribution $P(x)$

Conditional Generative Model: Learn a probability distribution $P(x|y)$



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 a probability distribution $P(x|y)$



These models may seem to be very distinct.
Although, they are actually **not fully** distinct.

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 a probability distribution $P(x|y)$

Recall Bayes' Rule:

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

Discriminative Model (Unconditional)
Conditional Generative Model
Generative Model
Prior over labels

We can build a conditional generative model from other components!

Discriminative vs. Generative Models

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

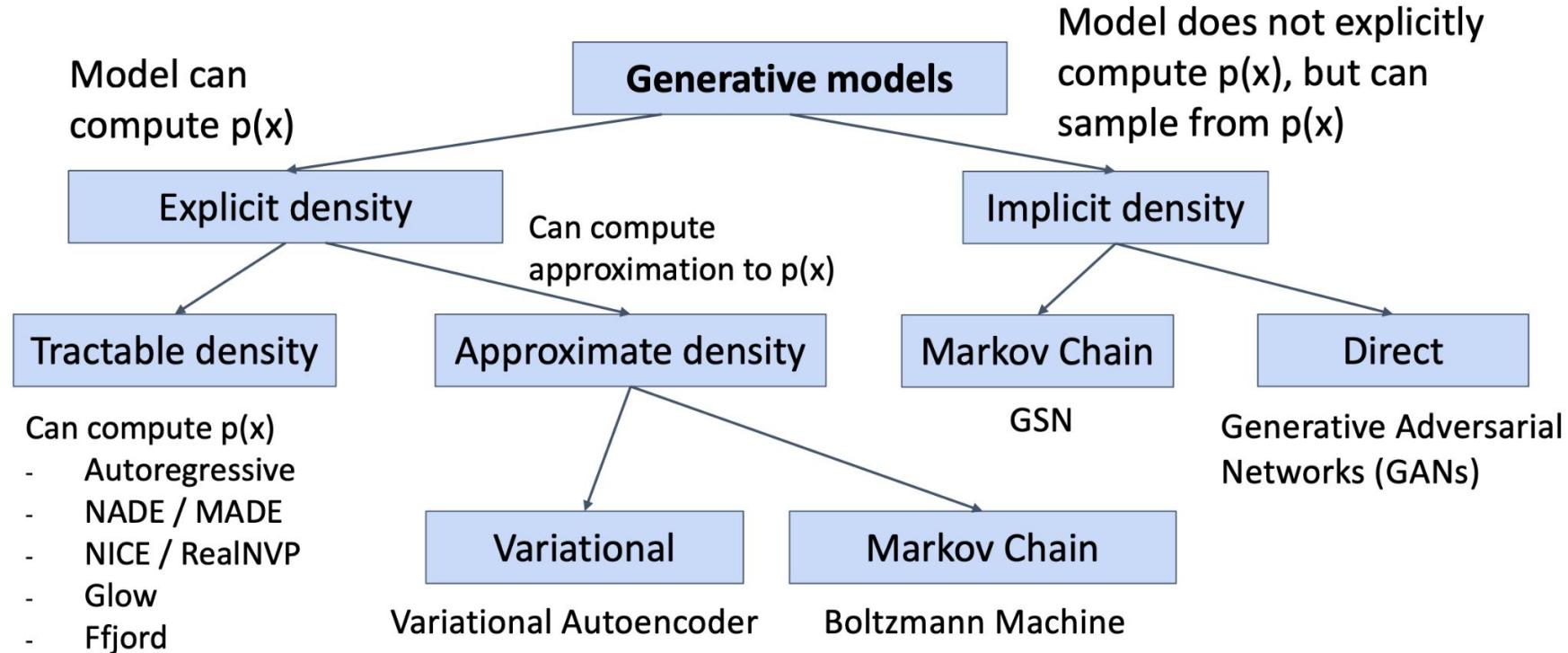
Conditional Generative

Model: Learn a probability distribution $P(x|y)$

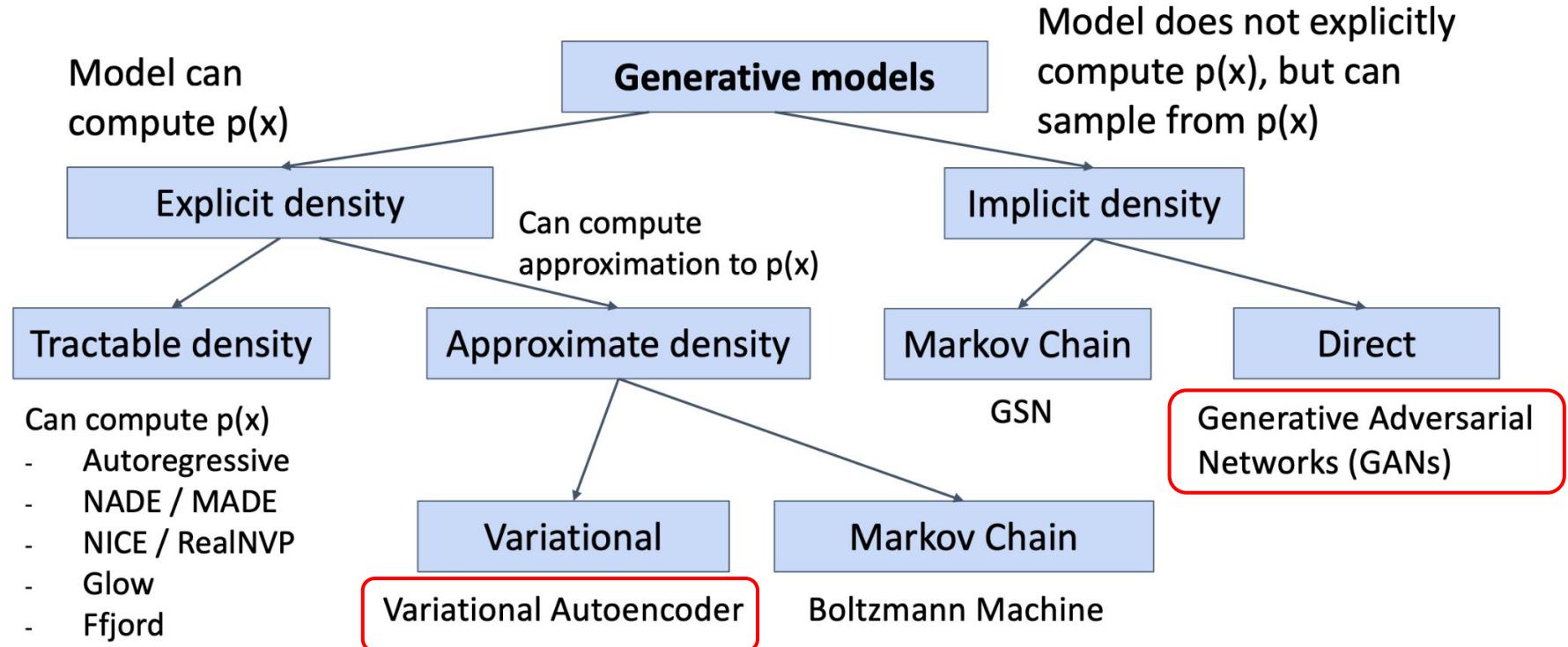


Assign labels, while rejecting outliers!
Generate new data conditioned on input labels

Taxonomy of Generative Models



Taxonomy of Generative Models



Autoencoder

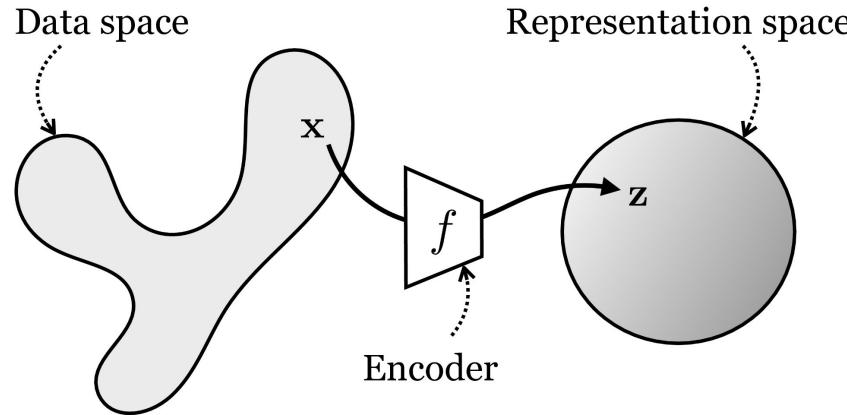
Representation Learning

The choice of representation is **critical**. Each type of representation makes some operations easy and others hard.

Q: How can we come up with the best representation possible?

Representation Learning

Goal of representation learning



The goal is to learn to map from datapoints, $x \in \mathcal{X}$, to abstract representation, $z \in \mathcal{Z}$

$$f : X \rightarrow Z$$

z is called a **vector embedding** of x

Representation Learning

What Makes for a Good Representation?

A good representation is one that makes subsequent problem solving easier.

- **Compression**

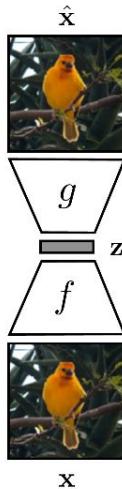
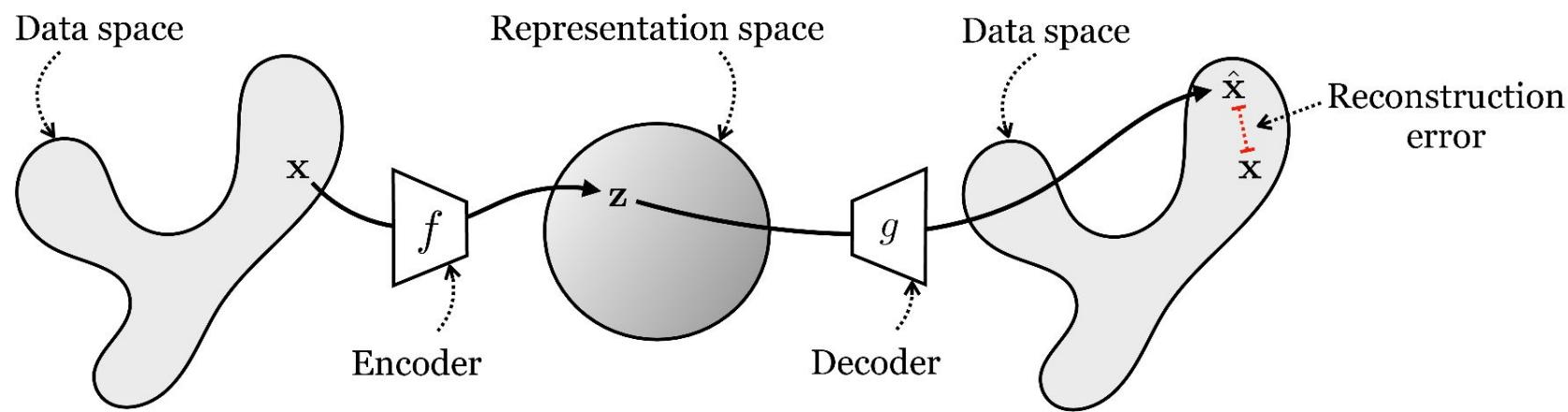
- Compressed representations require less memory to store.
- Compression is a way to achieve invariance to nuisance factors.

- **Prediction**

- Facilitating the prediction of important world properties—the future, the past, mental states, cause and effect, and so on—is perhaps the defining property of a good representation.

Autoencoder

Autoencoders are one of the **oldest and most common kinds of representation learners**.

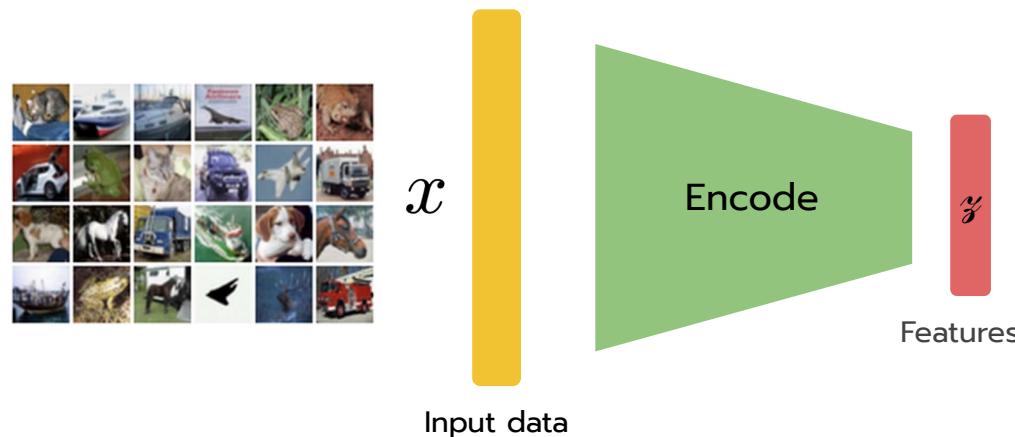


Autoencoder

Architecture

Unsupervised method for learning feature vectors from raw data x , without any labels.

“Autoencoding” = encoding itself



Features should extract useful information (maybe object identities, properties, scene type, etc) that we can use for downstream tasks

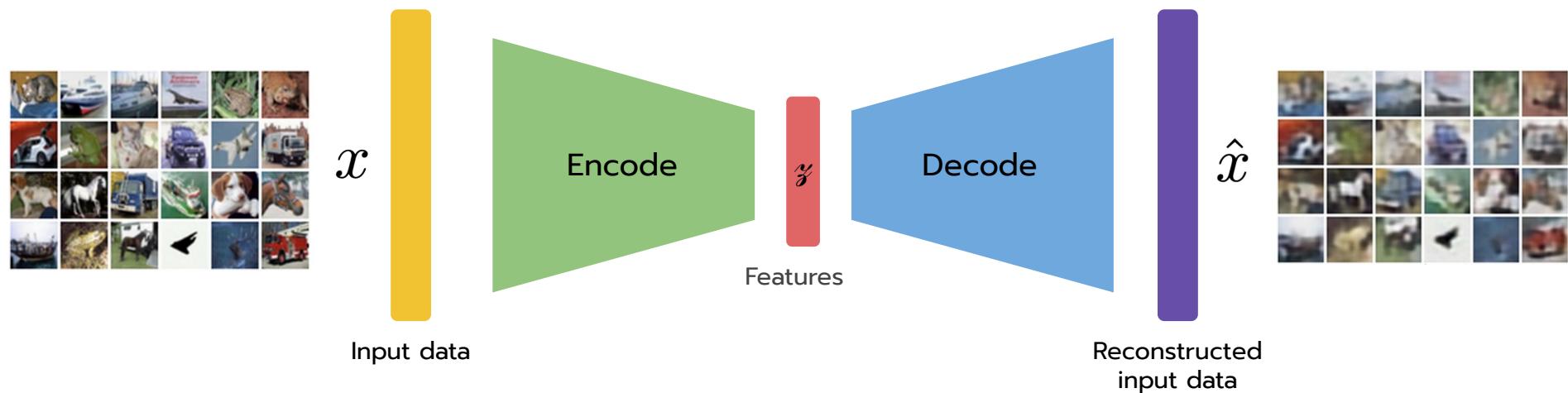
Problem: How can we learn this **feature transform** from raw data?

Autoencoder

Architecture

Idea: Use the features to **reconstruct the input data** with a decoder

“Autoencoding” = encoding itself

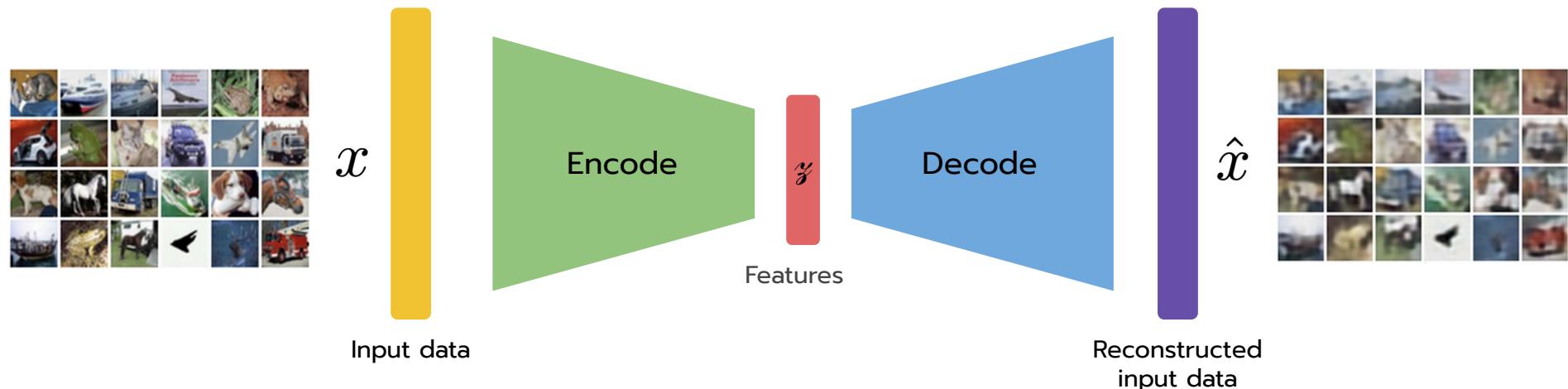


Autoencoder

Architecture

Idea: Use the features to **reconstruct the input data** with a decoder

"Autoencoding" = encoding itself



"Encoder" learns mapping from the data, x , to a low-dimensional latent space, z

"Decoder" learns mapping back from latent space, z , to reconstructed observation, \hat{x}

Autoencoder

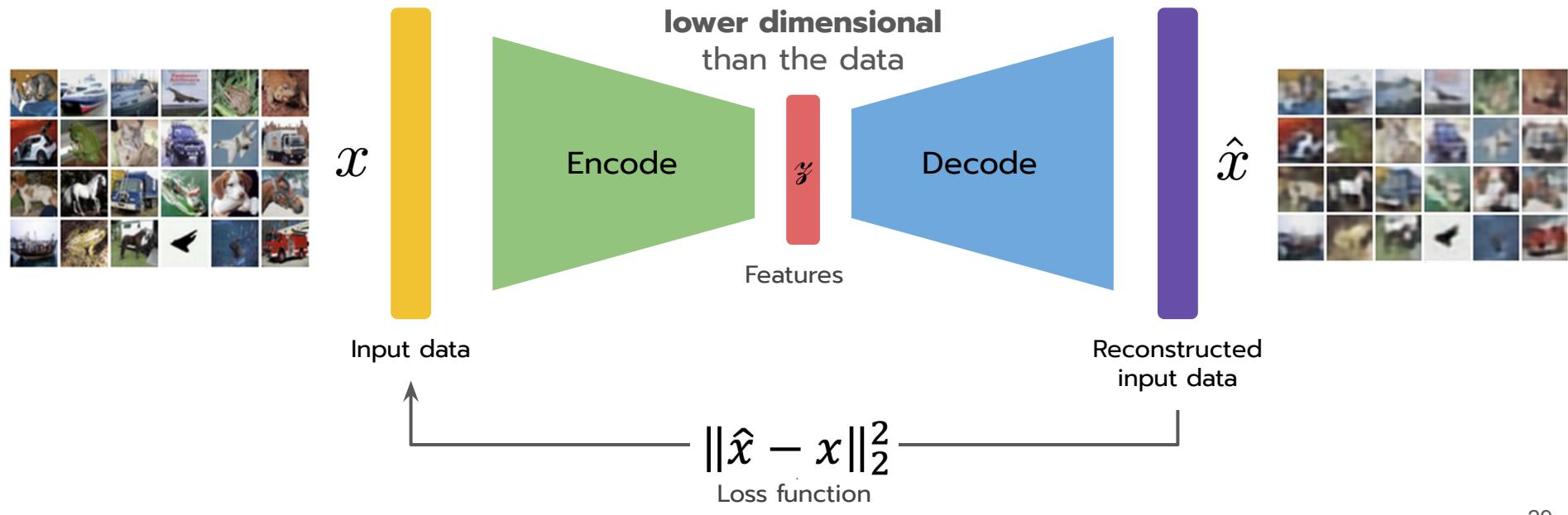
Architecture

Loss: L2 distance between input and reconstructed data.

What can we use the encoder for?
we can actually throw away for decoder (since we have autoencoder)..for classification > transfer learning & classification

Does not use any labels! Just raw data!

so that the memory of the model of the space is compressed version of the input

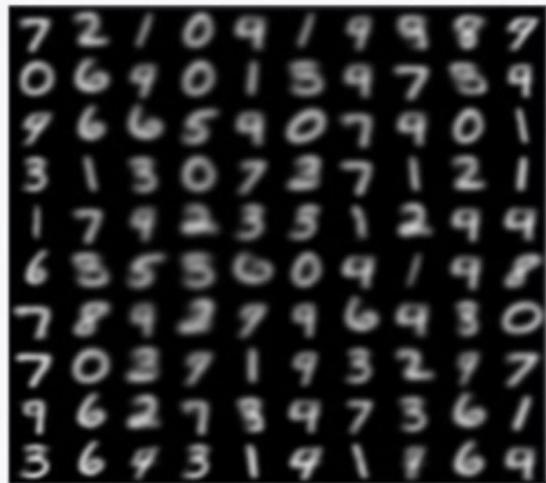


Autoencoder

Latent vector

Autoencoding is a form of compression! **Size of latent vector** affects the reconstructed data quality.

2D latent space



5D latent space



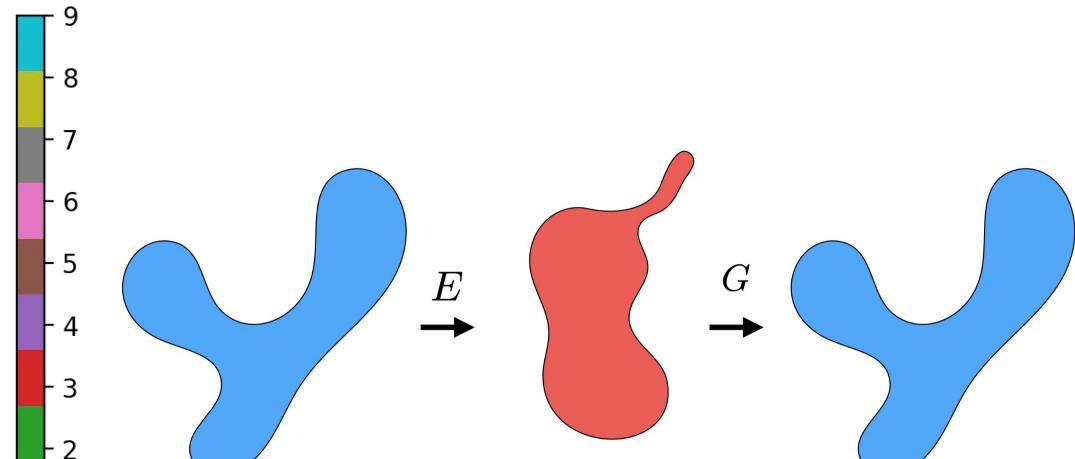
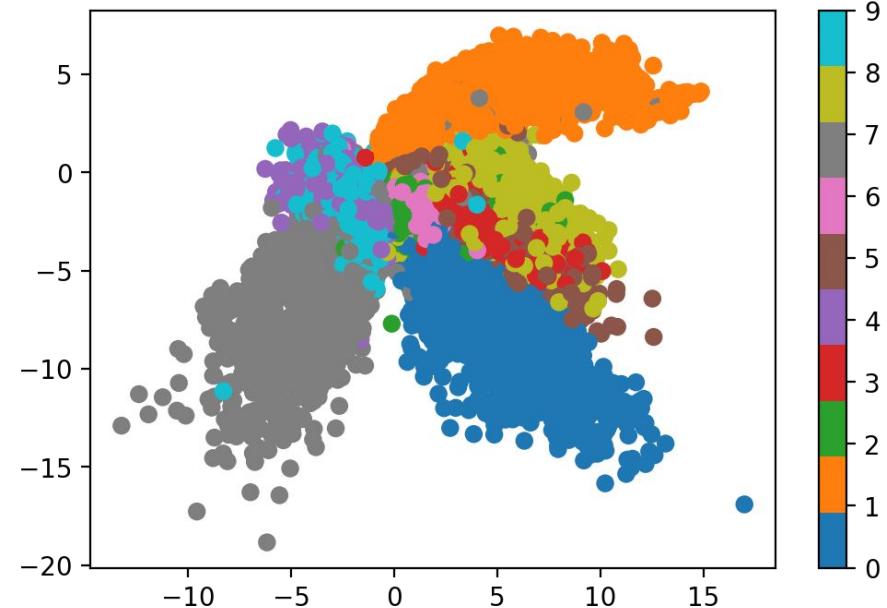
Ground truth



Autoencoder

Representation (Latent) Space

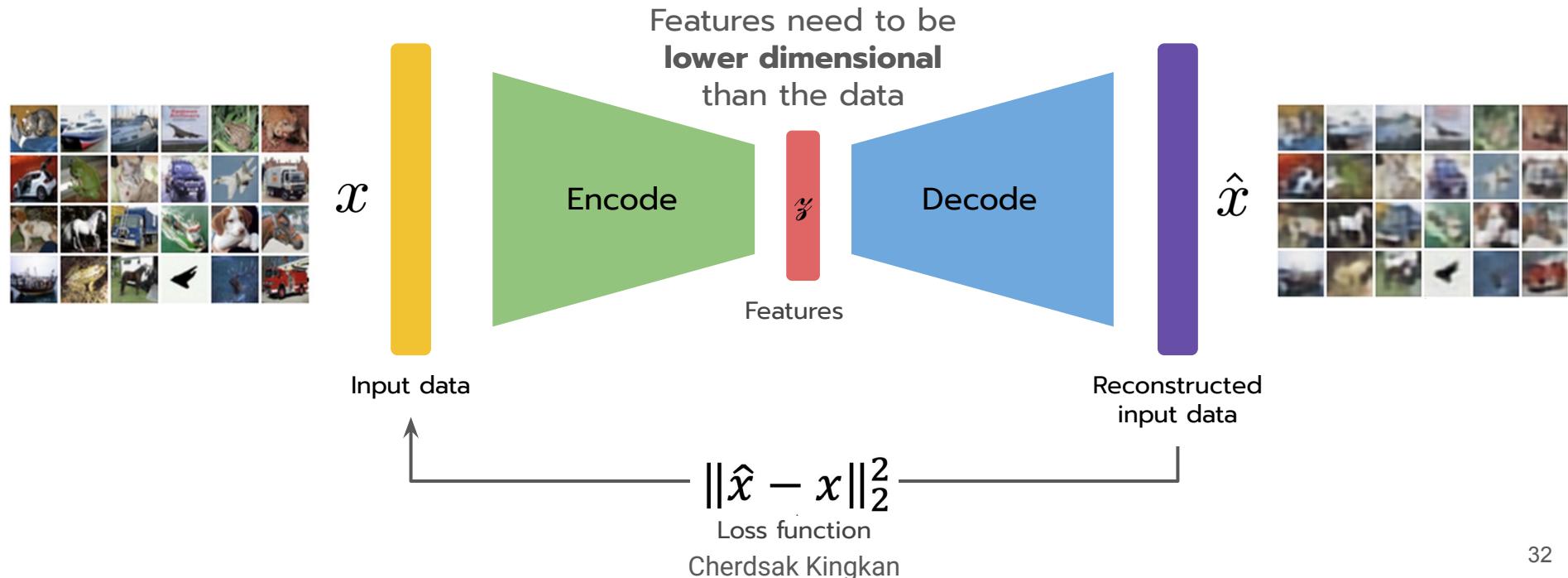
Latent Space of MNIST dataset



Autoencoder

Autoencoders (Regular, non-variational)

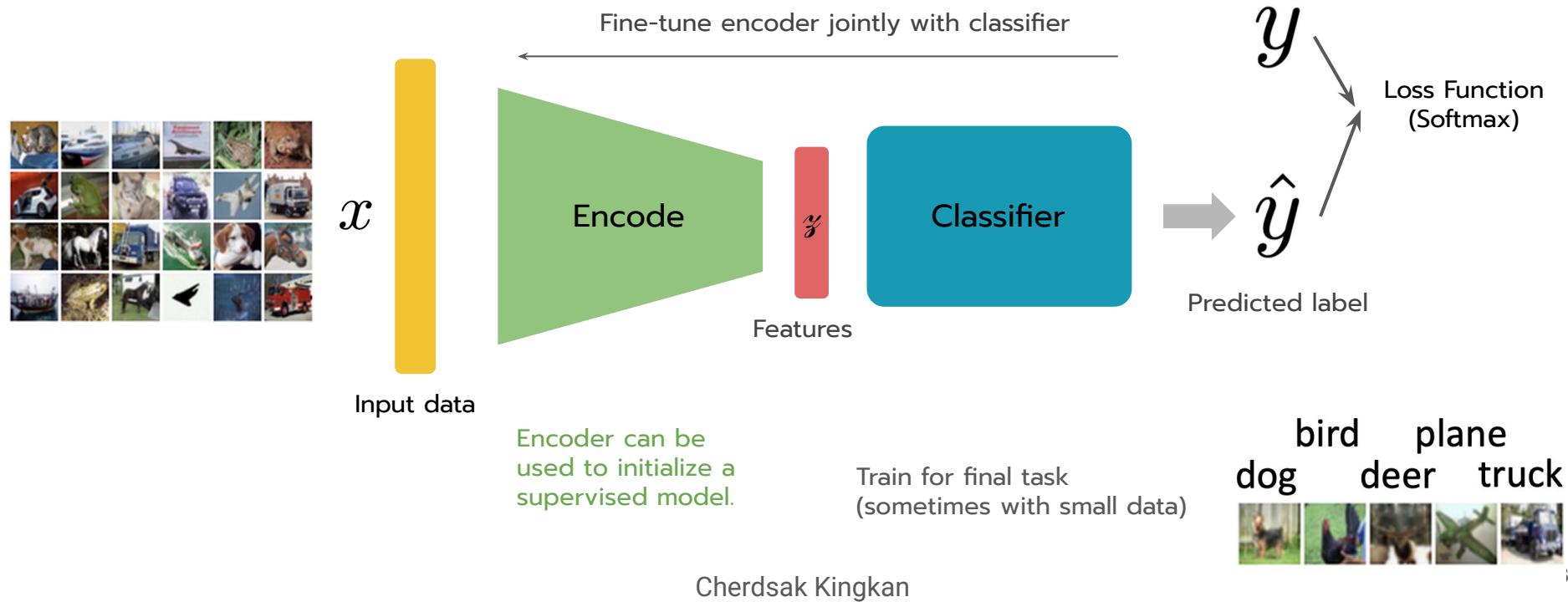
After training, **throw away decoder** and use encoder for a downstream task



Autoencoder

Autoencoders (Regular, non-variational)

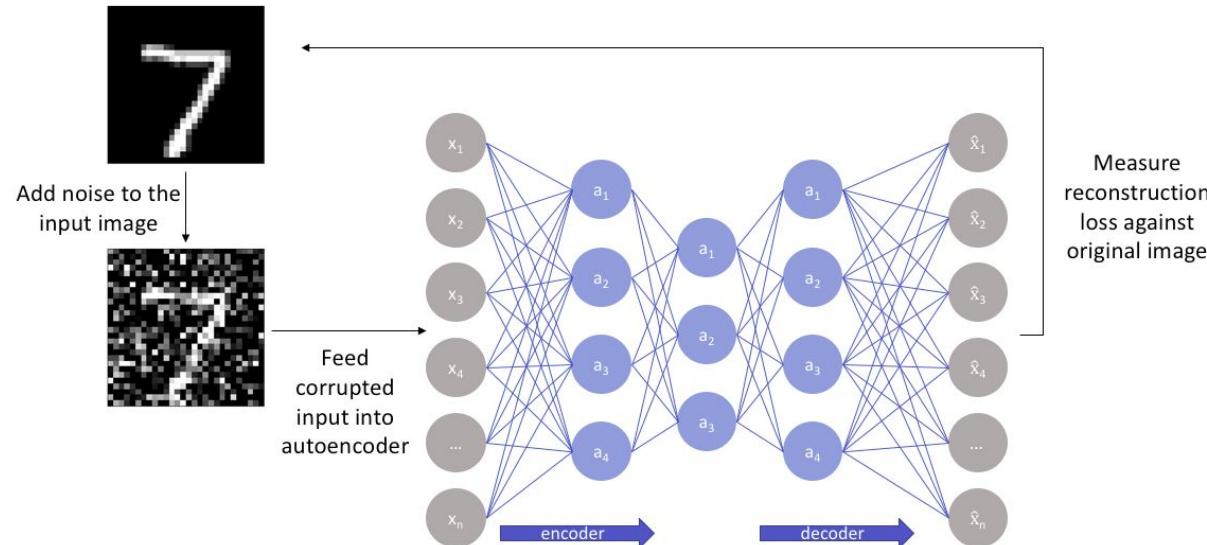
After training, **throw away decoder** and use encoder for a downstream task



Autoencoder

Applications

- Dimensionality Reduction and Feature Extraction
- Denoising Data
- Anomaly detections



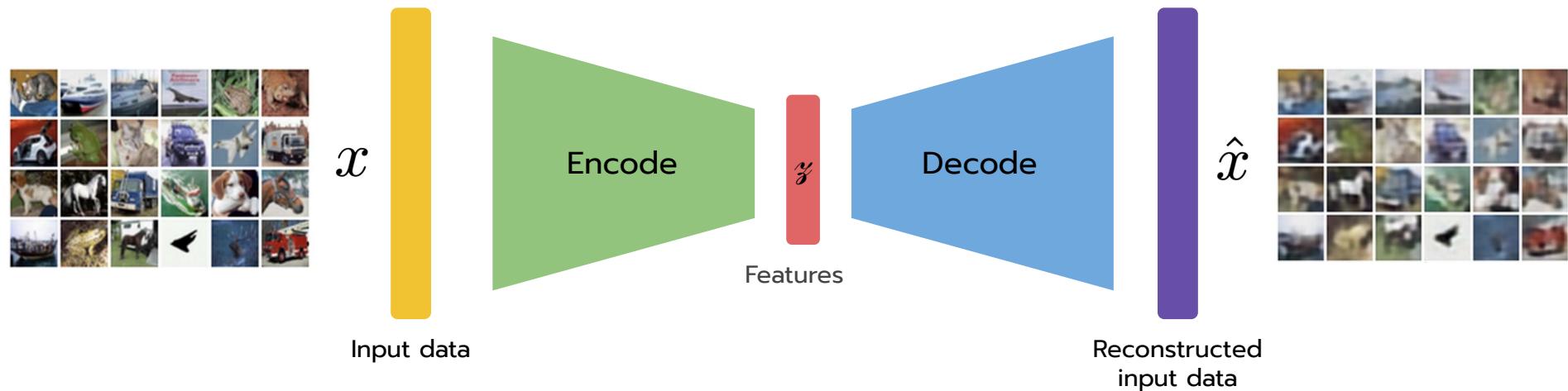
Autoencoder

Autoencoders (Regular, non-variational)

Autoencoders learn **latent features** for data without any labels!

Can use features to initialize a supervised model

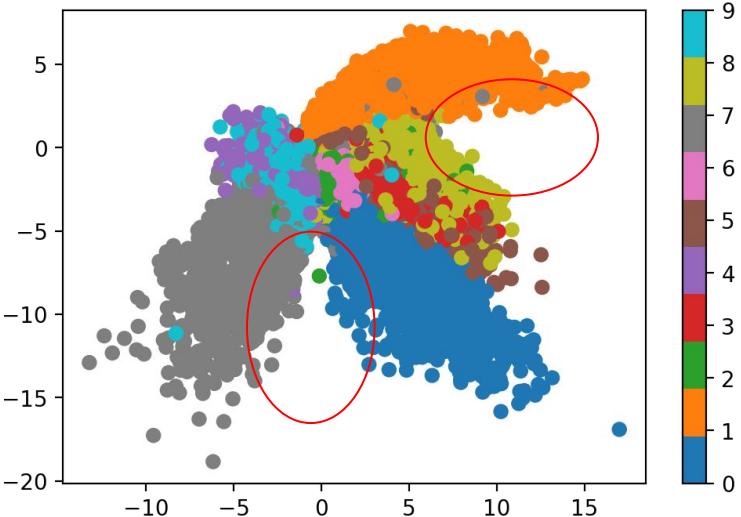
Not probabilistic: No way to sample **new data** from learned model



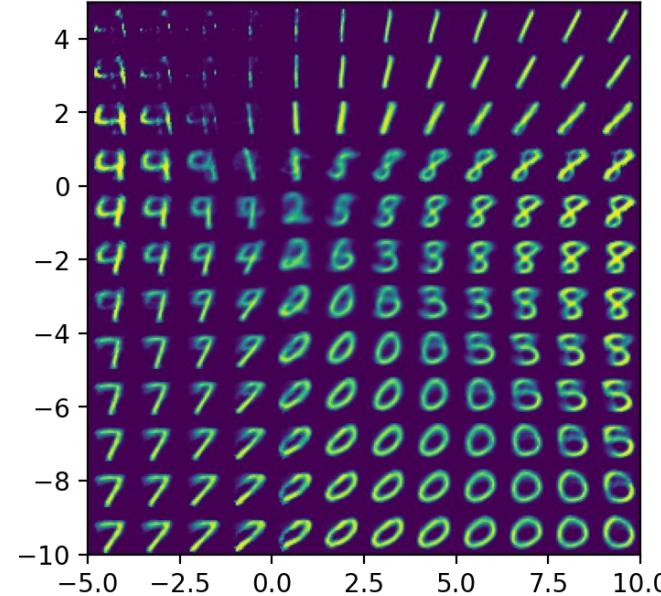
Autoencoder

Autoencoders (Regular, non-variational)

There are “gaps” in the latent space, where data is never mapped to.



If we sample a latent vector from a **region** in the **latent space that was never seen** by the decoder during training, the output might not make any sense at all.



the latent space, \mathbb{z} , can become **disjoint** and **non-continuous**.

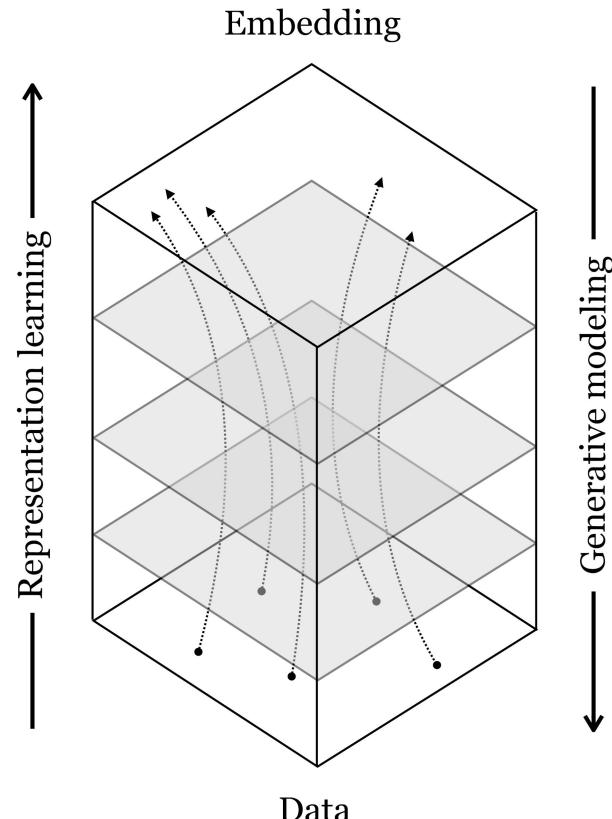
Variational Autoencoder

Generative Modeling meets Representation Learning

- Deep nets transform datapoints, layer by layer
- Each layer is a different representation of the data

Representation Learning (analysis):
maps data to abstract representations (**embeddings**).

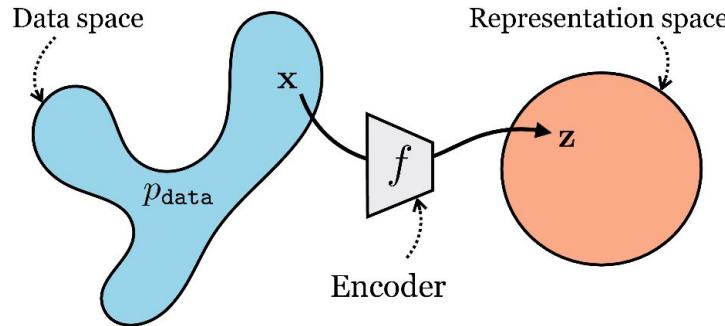
Generative Modeling (synthesis):
map abstract representations to data



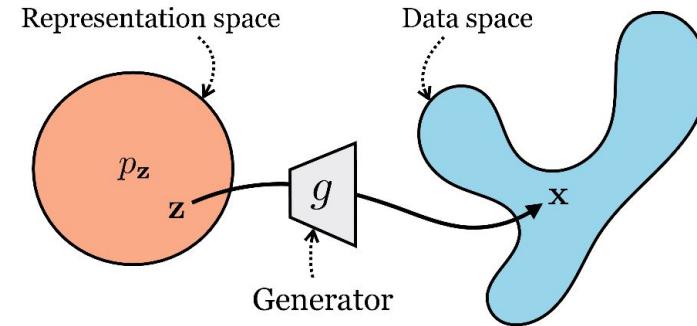
Generative Modeling meets Representation Learning

Latent Variables as Representations

Representation learning



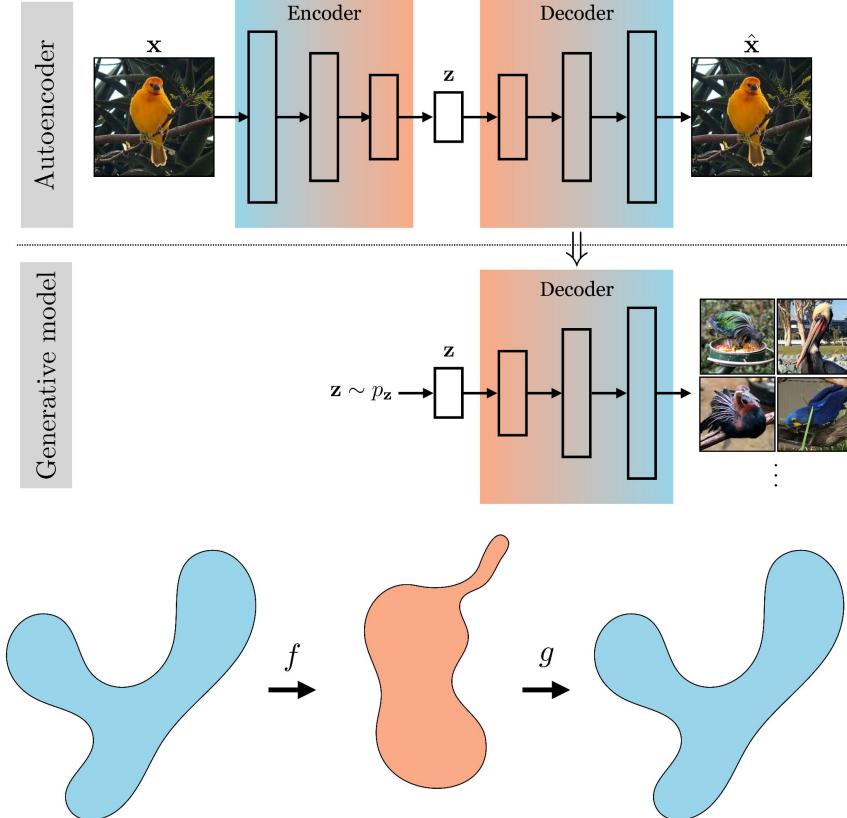
Generative modeling



Variational Autoencoder

The Decoder of an Autoencoder is a Data Generator

MORE LATENT SIZE, THE BETTER WE ARE TO



Variational autoencoders (VAEs)

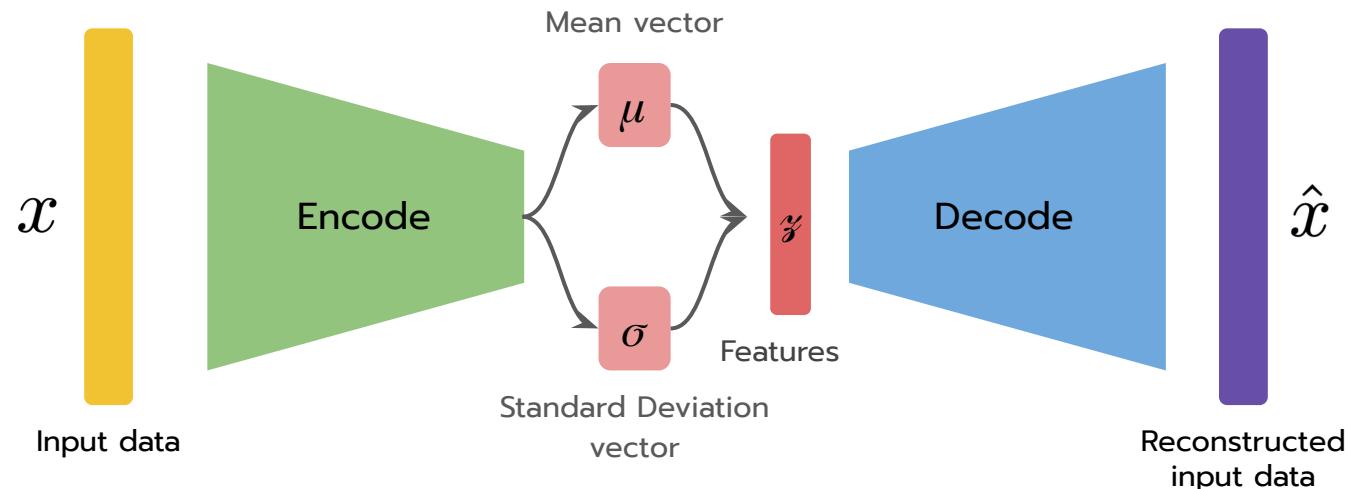
are a way to turn an autoencoder into a proper generative model, which can be sampled from and which maximizes data likelihood under a formal probabilistic model.

The trick is very simple: just take a vanilla autoencoder and (1) regularize the latent distribution to squish it into a Gaussian (or some other base distribution), (2) add noise to the output of the encoder.

Variational Autoencoder

Variational autoencoders are a probabilistic twist on autoencoders!

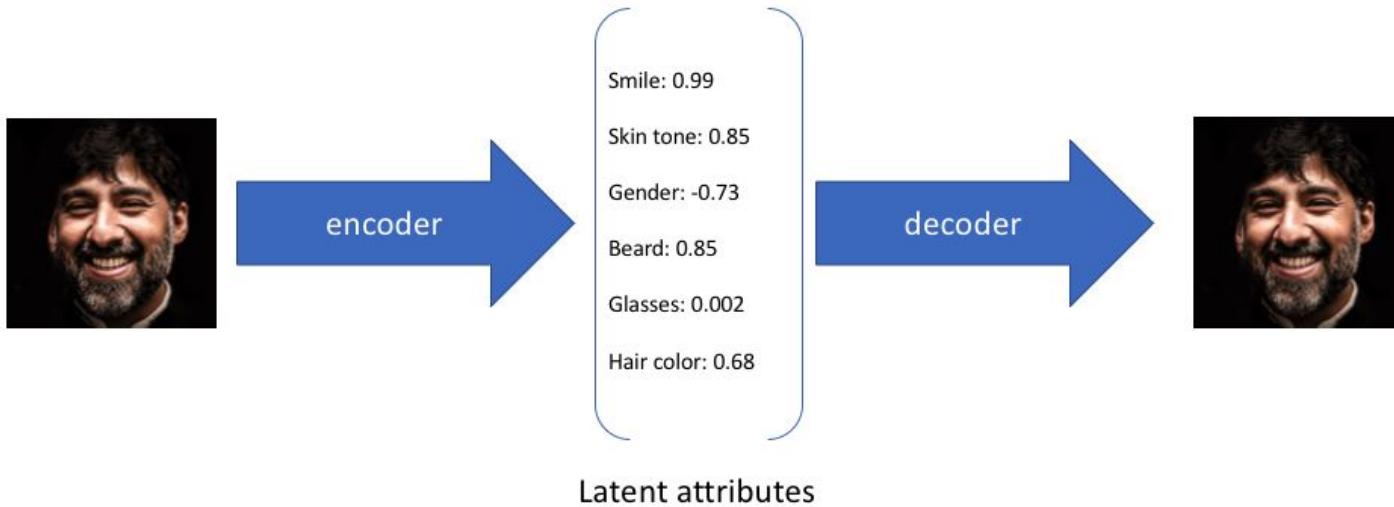
Sample from the mean and standard deviation to compute latent sample



Overview of VAE : Instead of learning an arbitrary function, the network learns the parameters of **probability distribution** modeling the data

Variational Autoencoder

Intuition



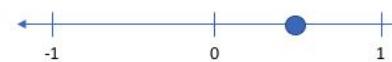
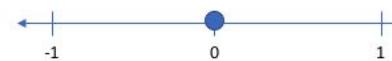
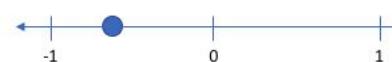
Variational Autoencoder

Intuition

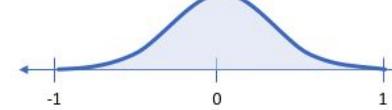
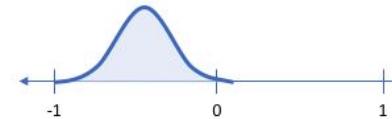
Using a variational autoencoder, we can describe latent attributes in probabilistic terms.



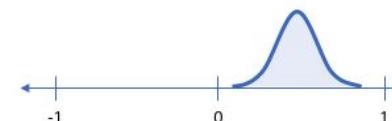
Smile (discrete value)



Smile (probability distribution)

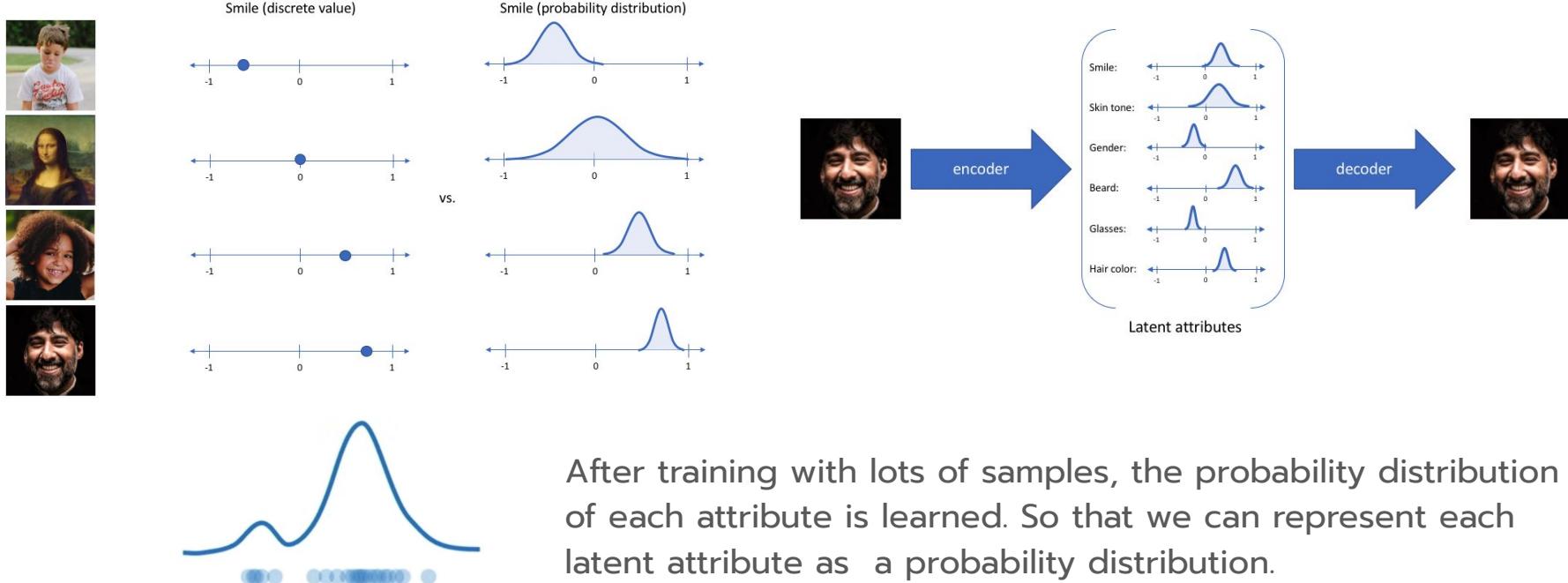


vs.



Variational Autoencoder

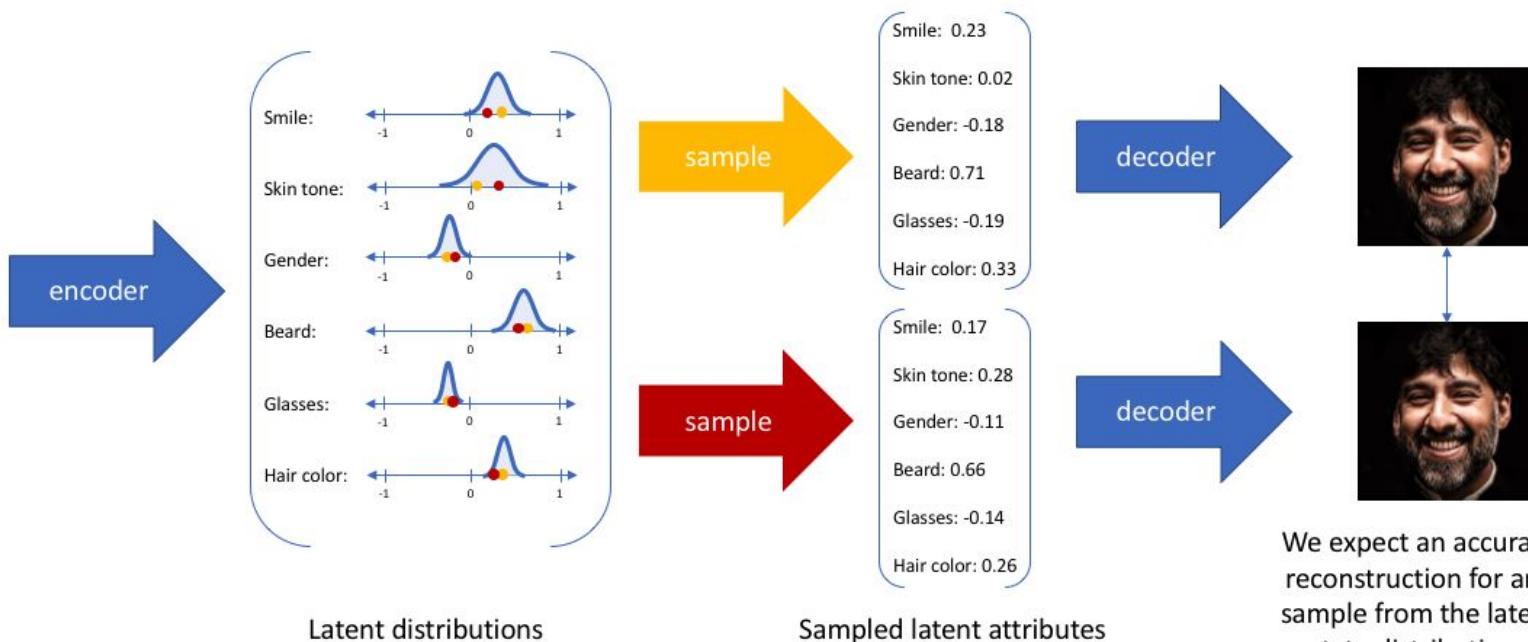
Intuition



Variational Autoencoder

Intuition

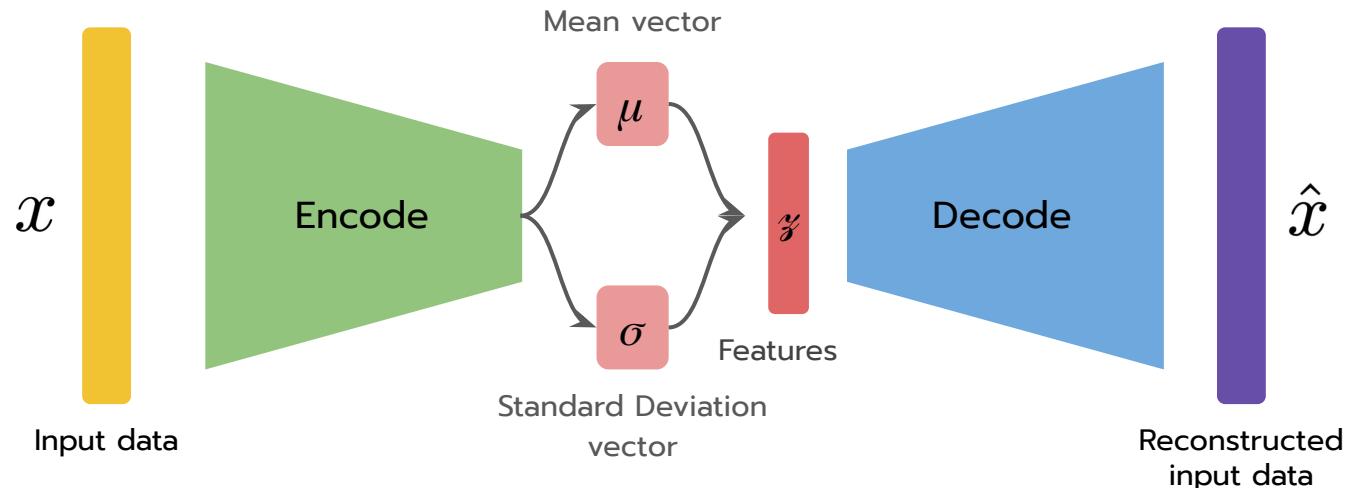
Values which are nearby to one another in latent space should correspond with very similar reconstructions.



Variational Autoencoder

Probabilistic spin on autoencoders: we want to do

1. Learn **latent features** \hat{z} from raw data
2. Sample from the model to **generate new data**



Variational Autoencoder

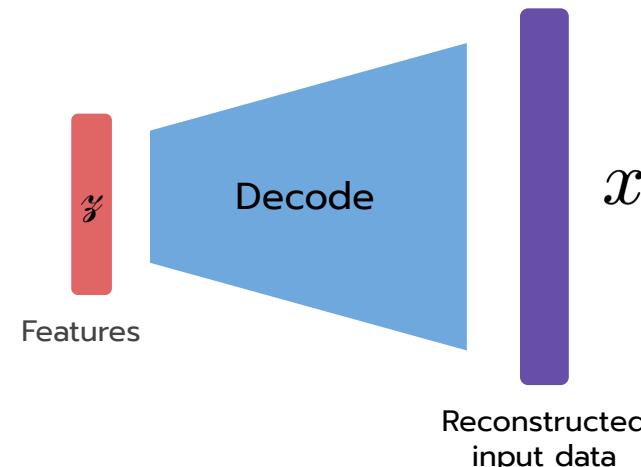
Probabilistic spin on autoencoders: we want to do

1. Learn latent features z from raw data
2. Sample from the model to **generate new data**

We can only see x , but we want to infer the attributes of z , i.e., we want to compute

$$p(z|x) = \frac{p(x|z)p(z)}{p(x)}$$

After training, we want to **sample** a hidden variable z which **generate** an observation x



Variational Autoencoder

Probabilistic spin on autoencoders: we want to do

1. Learn latent features z from raw data
2. Sample from the model to **generate new data**

We can only see x , but we want to infer the attributes of z , i.e., we want to compute

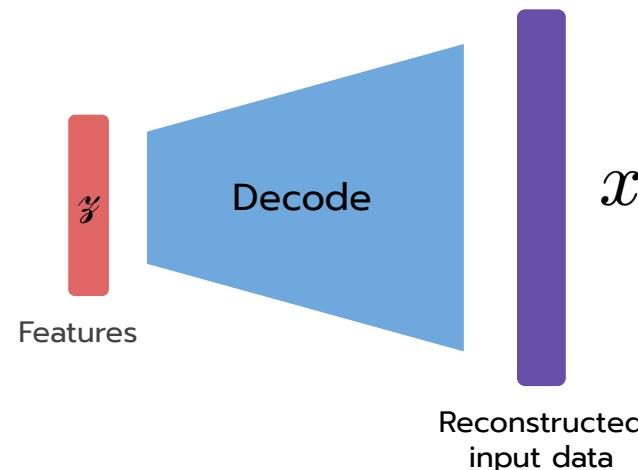
compute with
Decoder network

we assumed
Gaussian prior

$$p(z|x) = \frac{p(x|z)p(z)}{p(x)}$$

Very difficult to compute (technically possible,
but practical impossible)

After training, we want to **sample** a hidden variable z which **generate** an observation x



Variational Autoencoder

Probabilistic spin on autoencoders: we want to do

1. Learn latent features z from raw data
2. Sample from the model to **generate new data**

We can only see x , but we want to infer the attributes of z , i.e., we want to compute

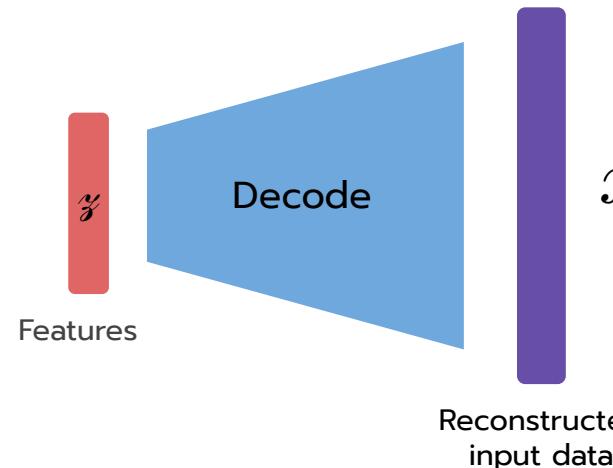
compute with
Decoder network

we assumed
Gaussian prior

$$p(z|x) = \frac{p(x|z)p(z)}{p(x)}$$

$$p(x) = \int p(x|z)p(z)dz$$

After training, we want to **sample** a hidden variable z which **generate** an observation x

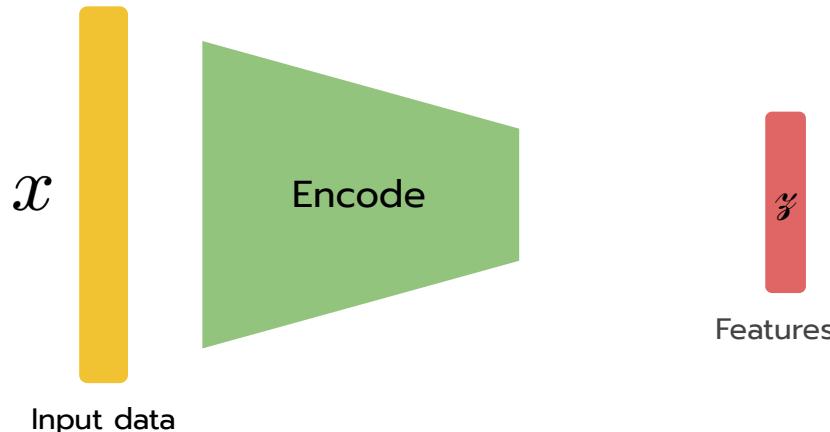


Impossible to integrate over all z (intractable)

Variational Autoencoder

Probabilistic spin on autoencoders: we want to do

1. Learn **latent features** \mathbf{z} from raw data
2. Sample from the model to generate new data



We can approximate $p(z|x)$ by another distribution $q(z|x)$, i.e.

$$q(z|x) \approx p(z|x)$$

the KL divergence is a measure of difference between two probability distributions. Thus, if we wanted to ensure that $q(z|x)$ was similar to $p(z|x)$, we could minimize the KL divergence between the two distributions.

$$\min D_{KL}(q(z|x) || p(z|x))$$

Variational Autoencoder

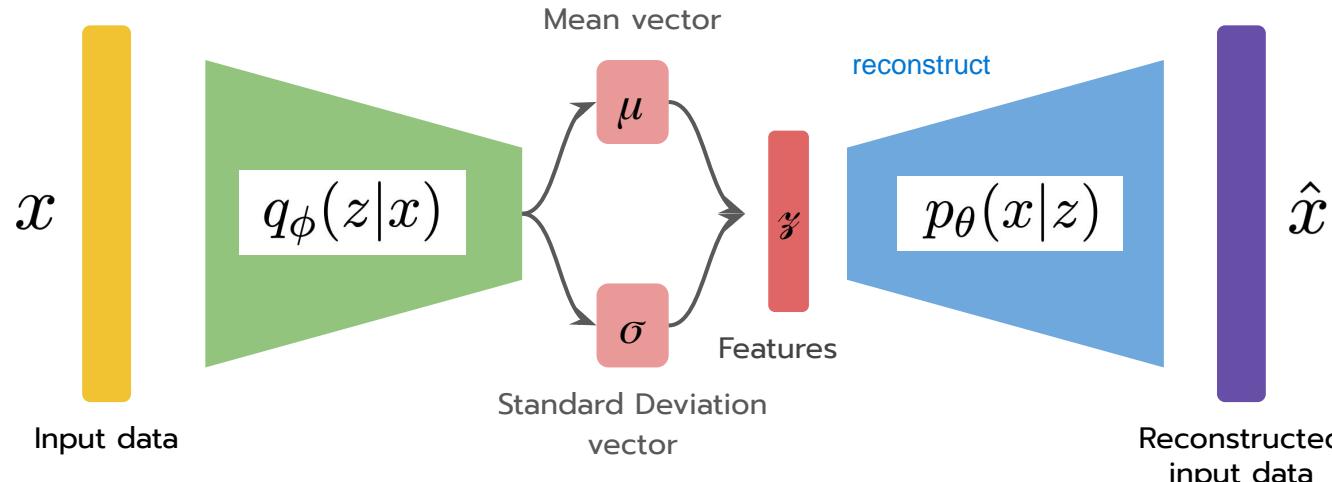
Common choice of prior - Gaussian Distribution

Probabilistic spin on autoencoders: we want to do

1. Learn **latent features** \mathbf{z} from raw data
2. Sample from the model to **generate new data**

$$p(z) = \mathcal{N}(\mu, \sigma)$$

mean / standard deviation



The reconstruction likelihood

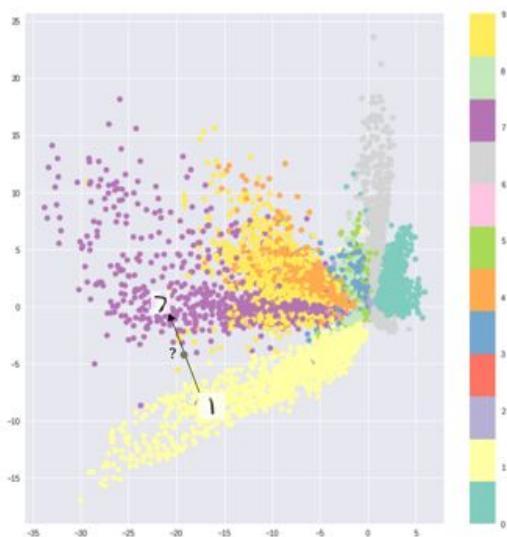
$$E_{q(z|x)} \log p(x|z) - D_{KL}(q(z|x)||p(z))$$

Regularization Loss: ensures that our **learned distribution** is similar to the true **prior distribution**

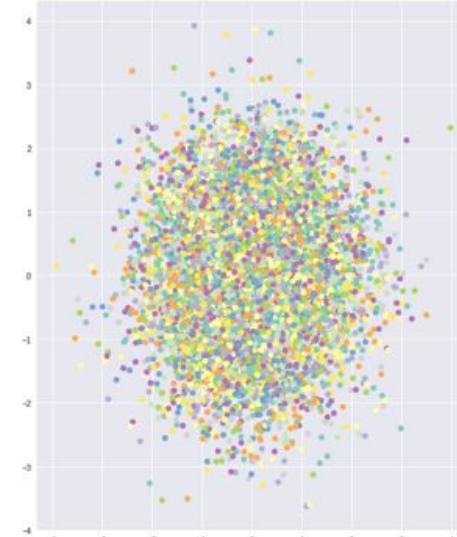
Variational Autoencoder

Why regularization term is needed

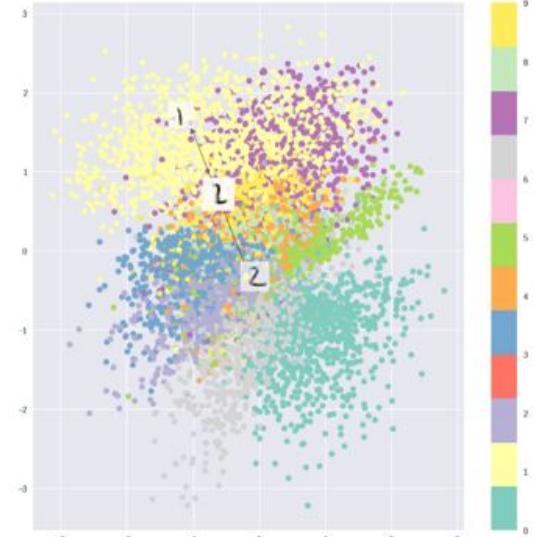
Only reconstruction loss



Only KL divergence



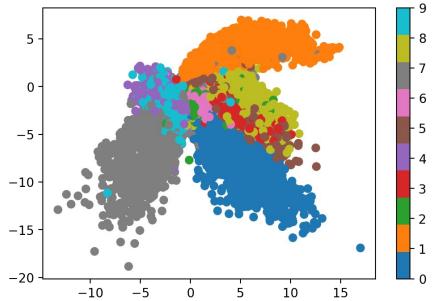
Combination



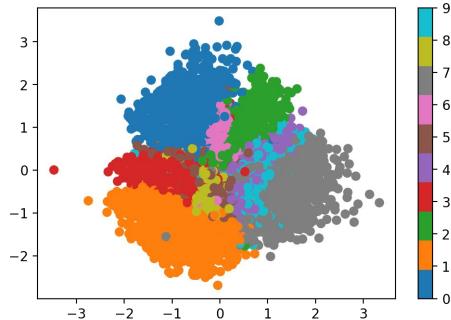
Variational Autoencoder

Why regularization term is needed

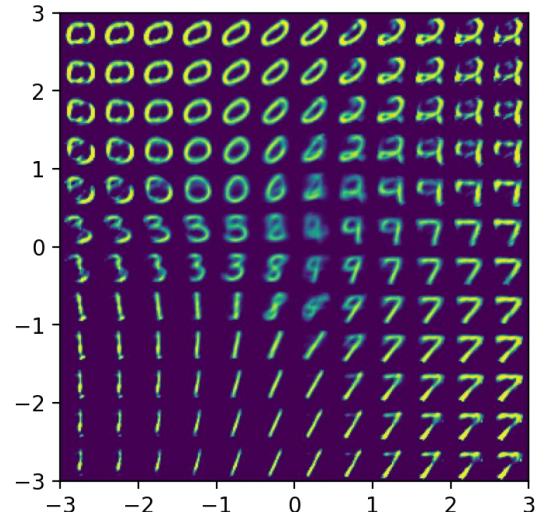
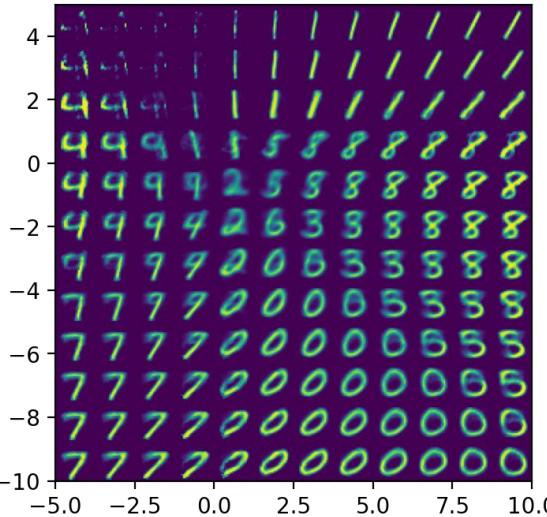
Autoencoder



Variational
Autoencoder

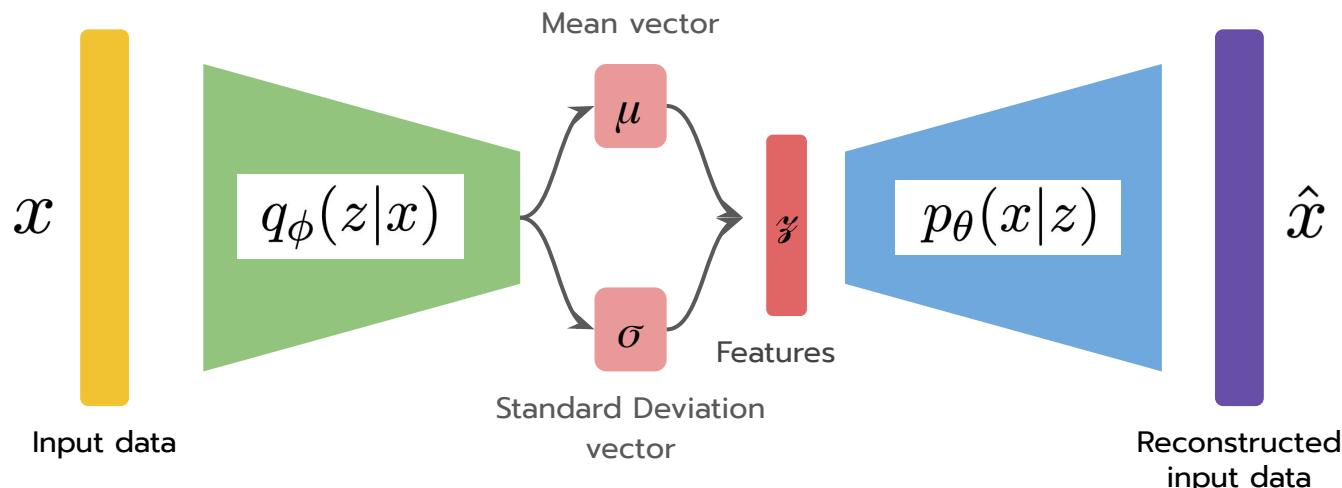


Cherdsak Kingkan



Variational Autoencoder

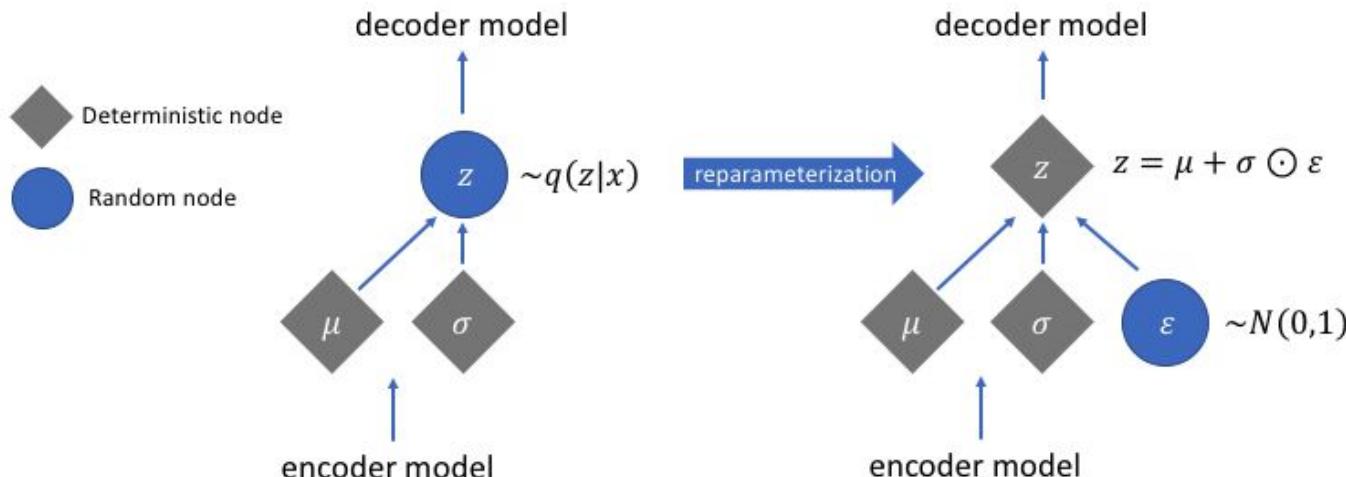
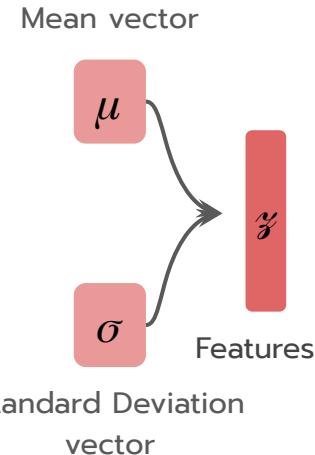
How to perform backpropagation when training training the model



$$E_{q(z|x)} \log p(x|z) - D_{KL}(q(z|x)||p(z))$$

Variational Autoencoder

Reparameterization trick



Variational Autoencoder

Latent perturbation

Slowly increase or decrease a **single latent variable**. Keep all other variables fixed.

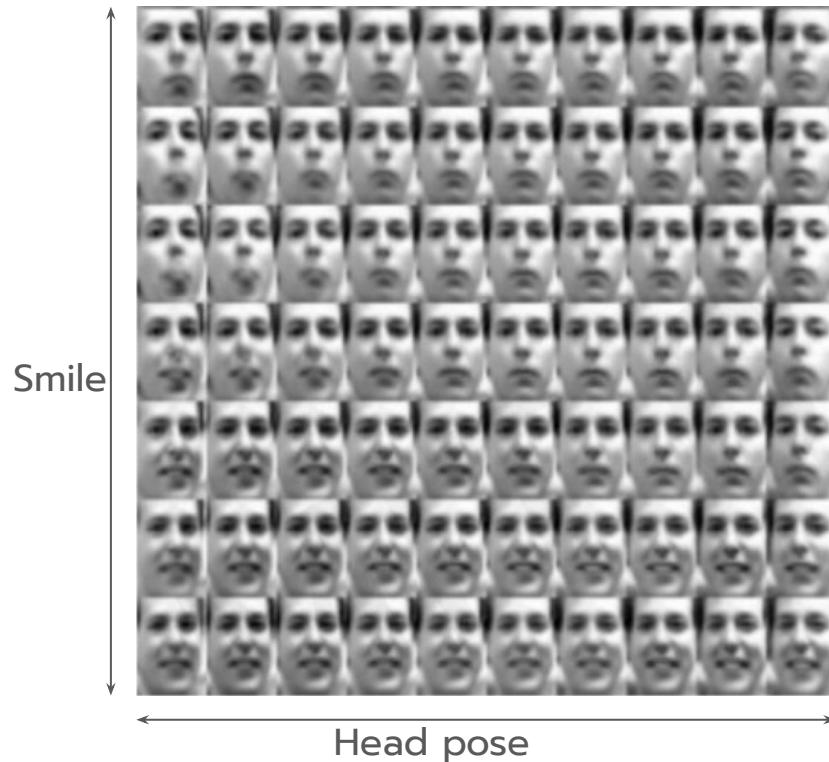


Head pose

Different dimension of γ encodes **different interpretable latent features**

Variational Autoencoder

Latent perturbation



Ideally, we want latent variables that are uncorrelated with each other.

Enforce diagonal prior on the latent variables to encourage independence

Disentanglement

Generative Adversarial Network

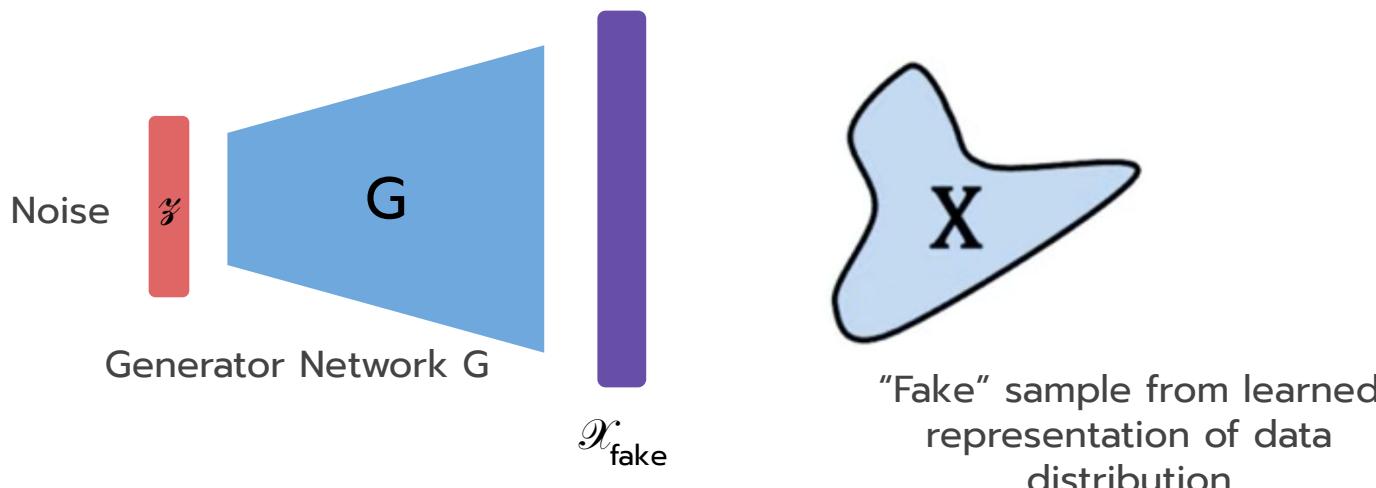
Generative Adversarial Network (GAN)

What if we just want to sample

Idea: do not explicitly model density, and instead just sample to generate new instances.

Problem: want to sample from complex distribution - cannot do this directly

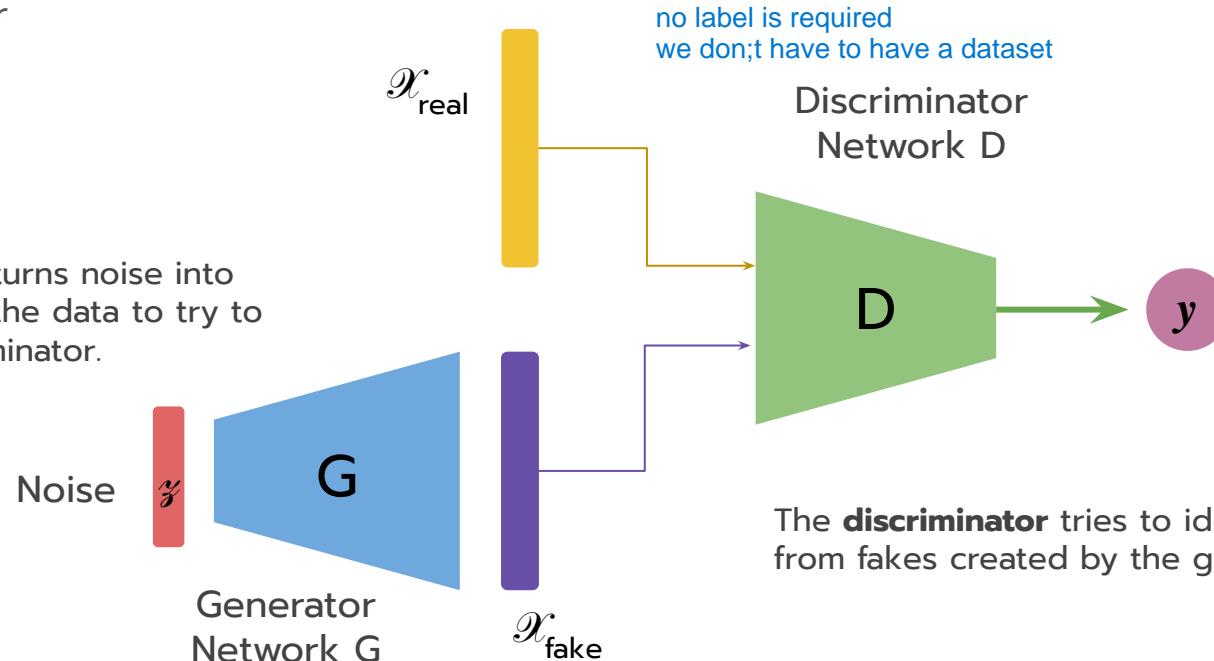
Solution: sample from something simple (e.g., noise), learn a transformation to the data distribution.



Generative Adversarial Networks (GANs)

GANs are a way to make a generative model by having two neural networks compete with each other

The **generator** turns noise into an imitation of the data to try to trick the discriminator.



The **discriminator** tries to identify real data from fakes created by the generator.

Generative Adversarial Network (GAN)

Intuition behind GANs

Discriminator looks at both real and fake data created by the generator

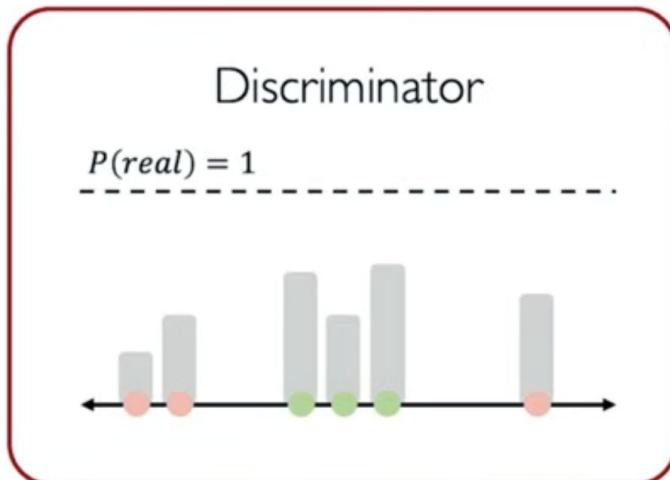
Generator starts from noise to try to create an imitation of the data



Generative Adversarial Network (GAN)

Intuition behind GANs

Discriminator tries to predict what's real and what's fake



Real data

Generator starts from noise to try to create an imitation of the data

Generator

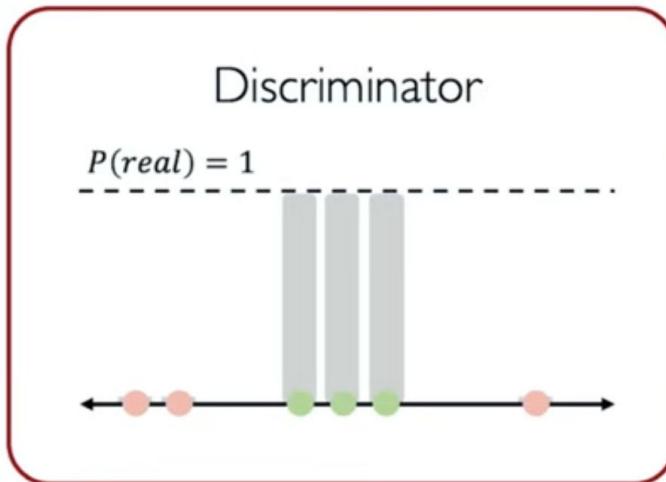


Fake data

Generative Adversarial Network (GAN)

Intuition behind GANs

Discriminator tries to predict what's real and what's fake



Generator starts from noise to try to create an imitation of the data

Generator



Real data



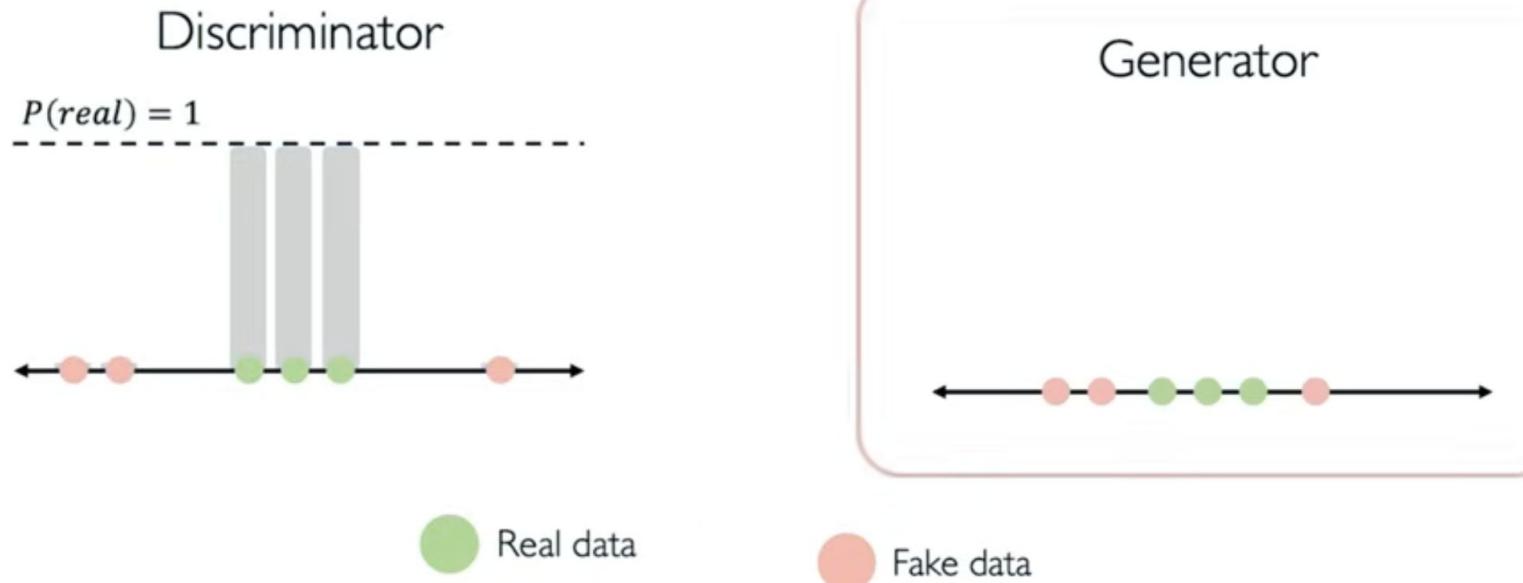
Fake data

Generative Adversarial Network (GAN)

Intuition behind GANs

Discriminator tries to predict what's real and what's fake

Generator tries to improve its imitation of the data

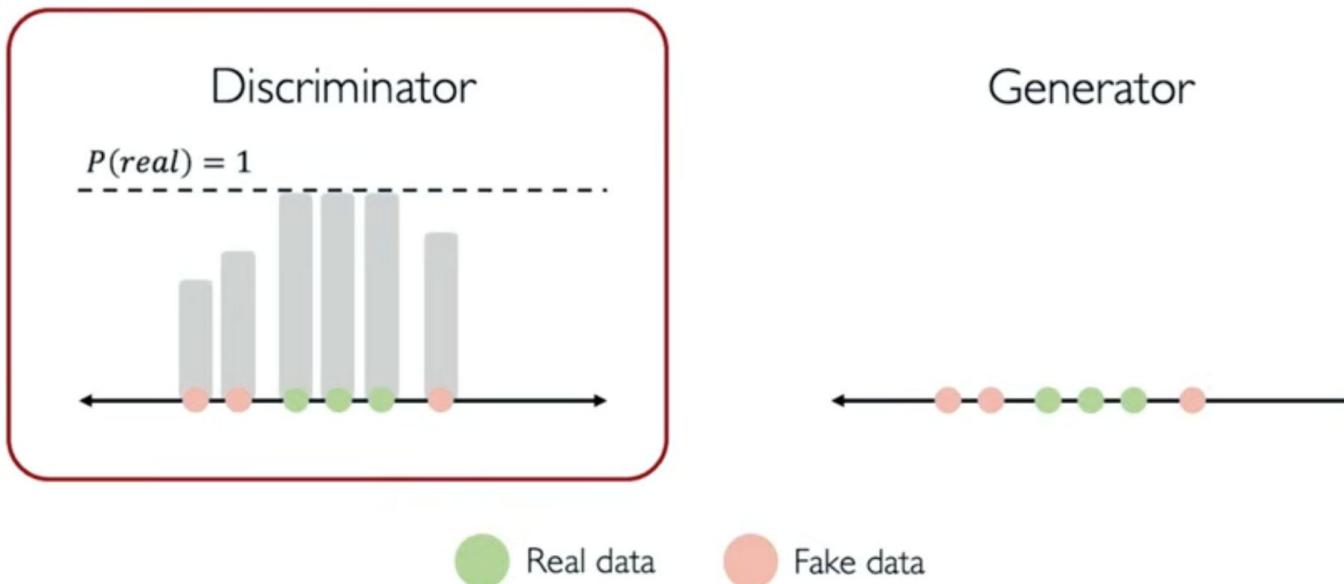


Generative Adversarial Network (GAN)

Intuition behind GANs

Discriminator tries to predict what's real and what's fake

Generator tries to improve its imitation of the data

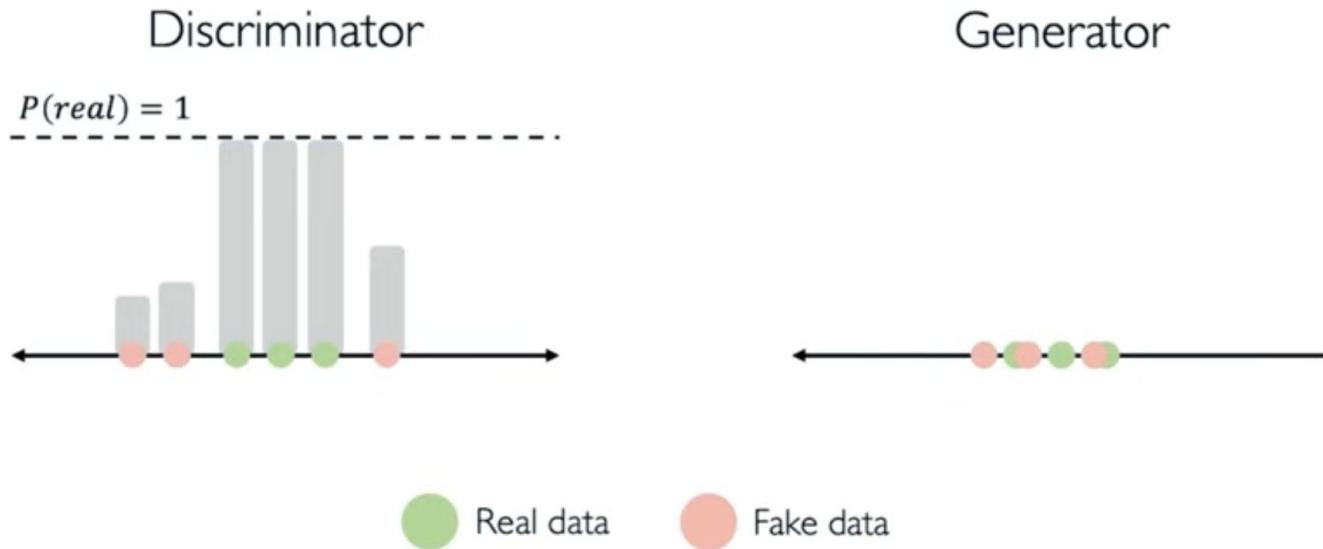


Generative Adversarial Network (GAN)

Intuition behind GANs

Discriminator tries to identify real data from fakes created by the generator.

Generator tries to create imitations of data to trick the discriminator.

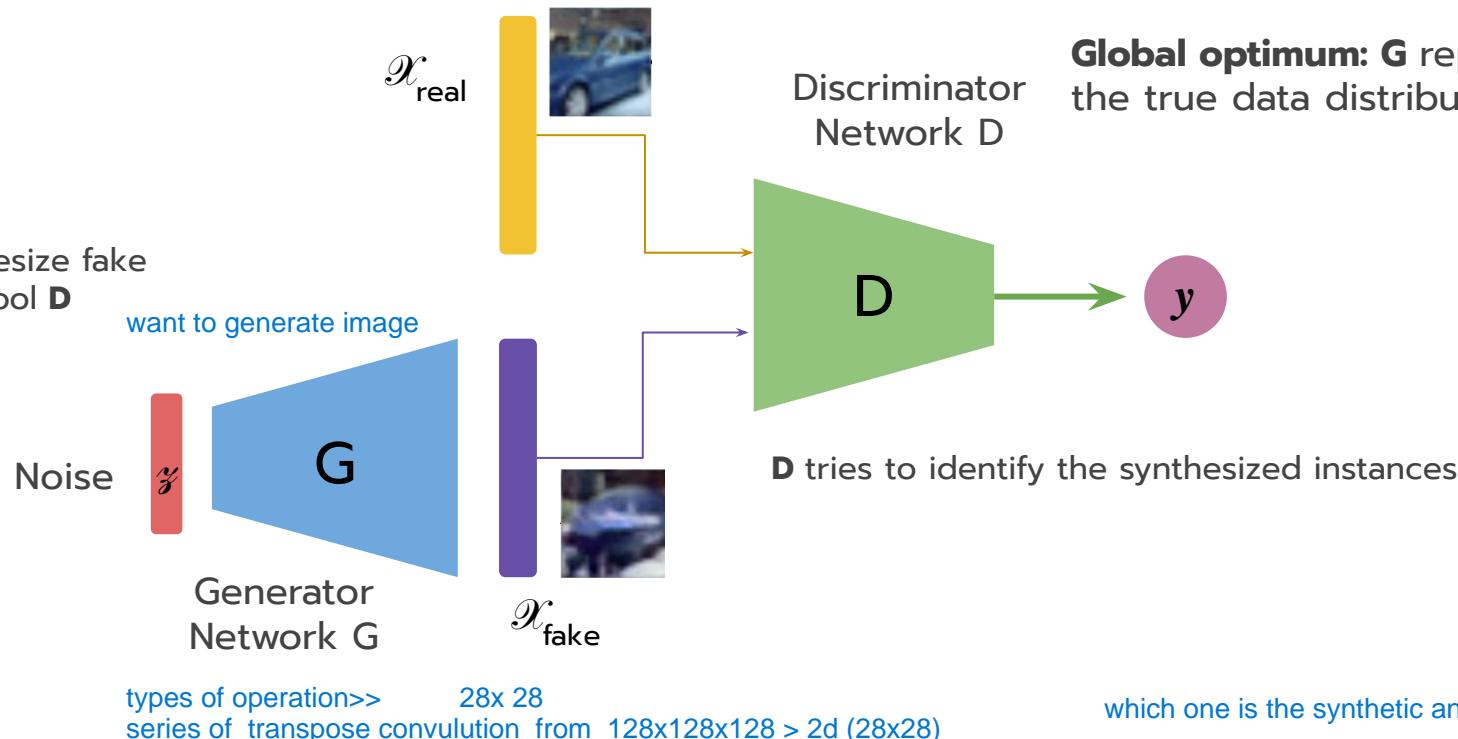


Generative Adversarial Networks (GANs)

Training GANs

Training: adversarial objective for **D** and **G**

Global optimum: **G** reproduces the true data distribution.



Generative Adversarial Networks (GANs)

Training GANs

Jointly train generator G and discriminator D with a **minimax game**

$$\min_G \max_D \left(E_{x \sim p_{data}} [\log D(x)] + E_{z \sim p(z)} [\log (1 - D(G(z)))] \right)$$

Discriminator wants $D(x) = 1$ for real data

Discriminator wants $D(x) = 0$ for fake data

Generator wants $D(x) = 1$ for fake data

generate as real as possible
discriminator cannot distinguish between real and synthetic image

Generative Adversarial Networks (GANs)

Training GANs

Jointly train generator G and discriminator D with a **minimax game**

Training G and D using alternating gradient updates

$$\min_{\mathbf{G}} \max_{\mathbf{D}} \left(E_{x \sim p_{data}} [\log \mathbf{D}(x)] + E_{\mathbf{z} \sim p(\mathbf{z})} \left[\log \left(1 - \mathbf{D}(\mathbf{G}(\mathbf{z})) \right) \right] \right)$$

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

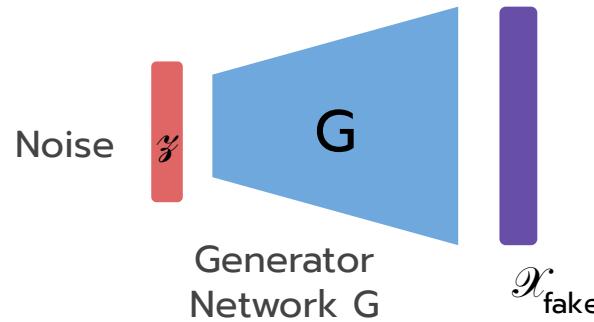
For t in 1, ... T:

1. (Update D) $\mathbf{D} = \mathbf{D} + \alpha_{\mathbf{D}} \frac{\partial \mathbf{V}}{\partial \mathbf{D}}$
2. (Update G) $\mathbf{G} = \mathbf{G} - \alpha_{\mathbf{G}} \frac{\partial \mathbf{V}}{\partial \mathbf{G}}$

Generative Adversarial Network (GAN)

Generating new data with GANs

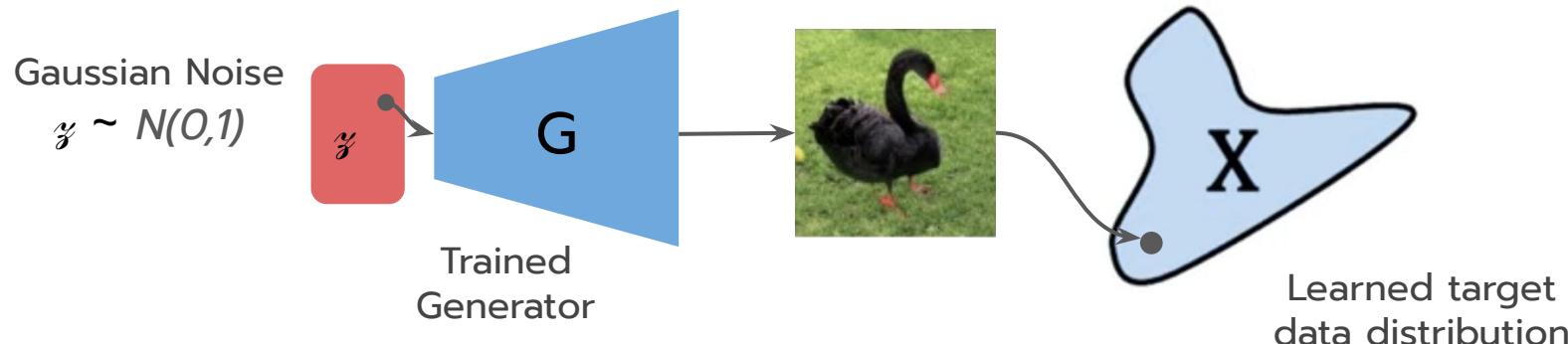
After training, use generator network to create **new data** that's never been seen before.



Generative Adversarial Network (GAN)

GANs are distribution transformers

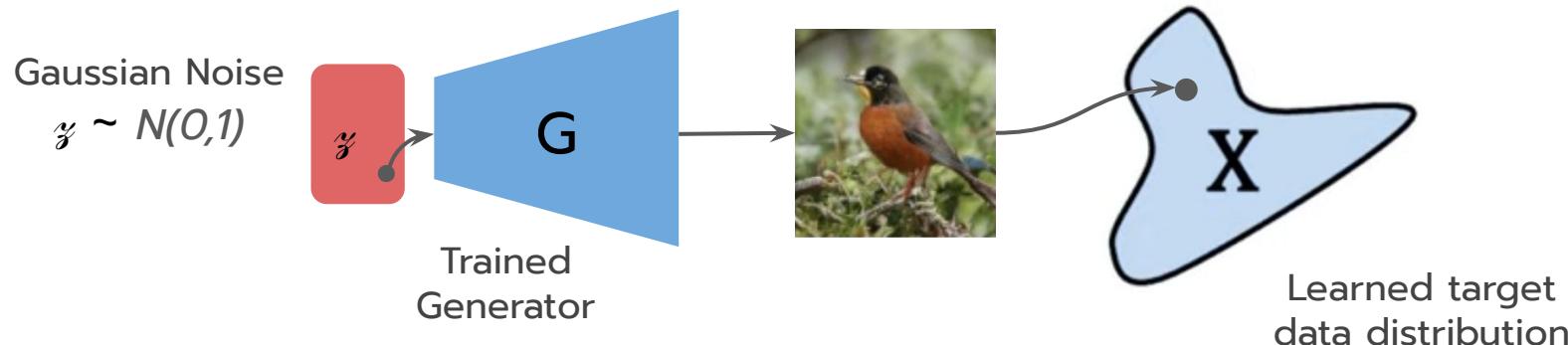
After training, use generator network to create **new data** that's never been seen before.



Generative Adversarial Network (GAN)

GANs are distribution transformers

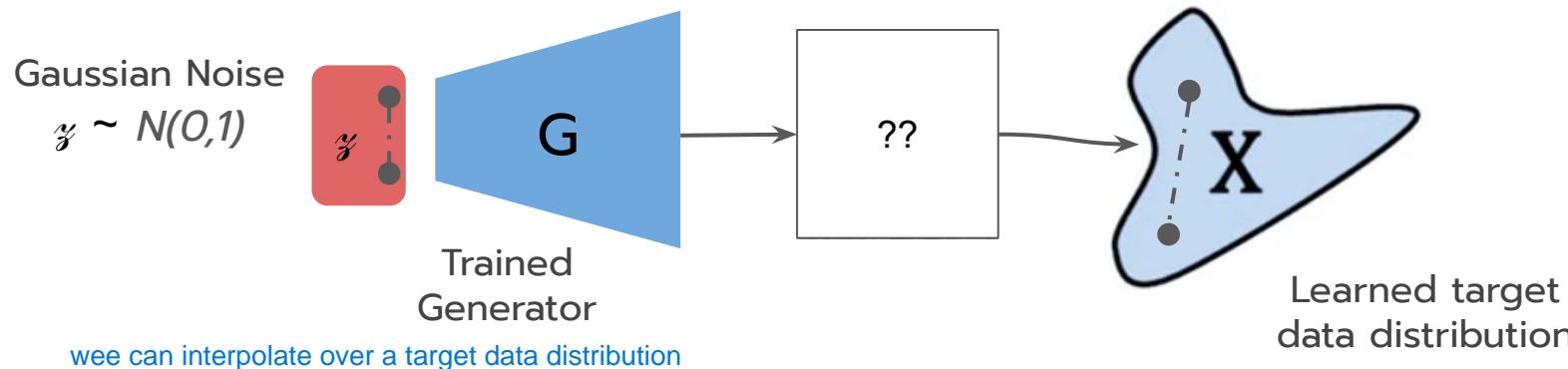
After training, use generator network to create **new data** that's never been seen before.



Generative Adversarial Network (GAN)

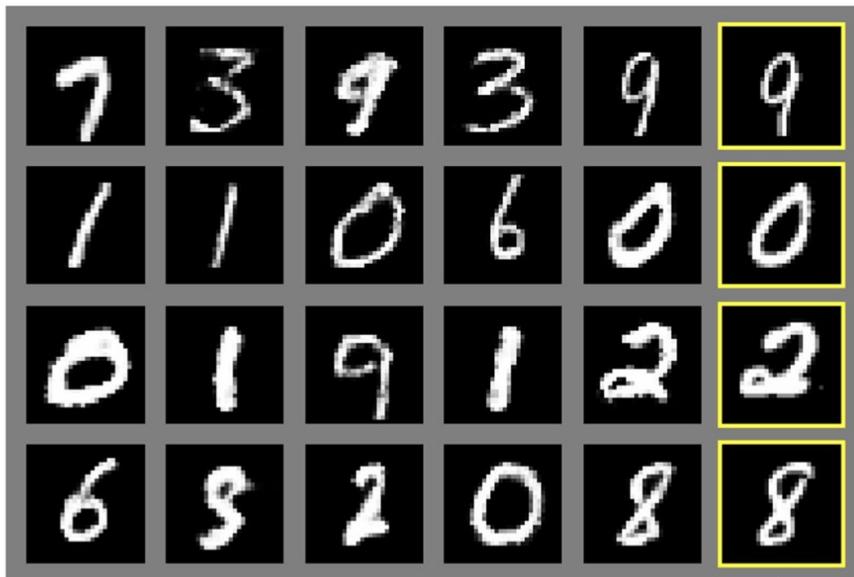
GANs are distribution transformers

After training, use generator network to create **new data** that's never been seen before.



Generative Adversarial Network (GAN)

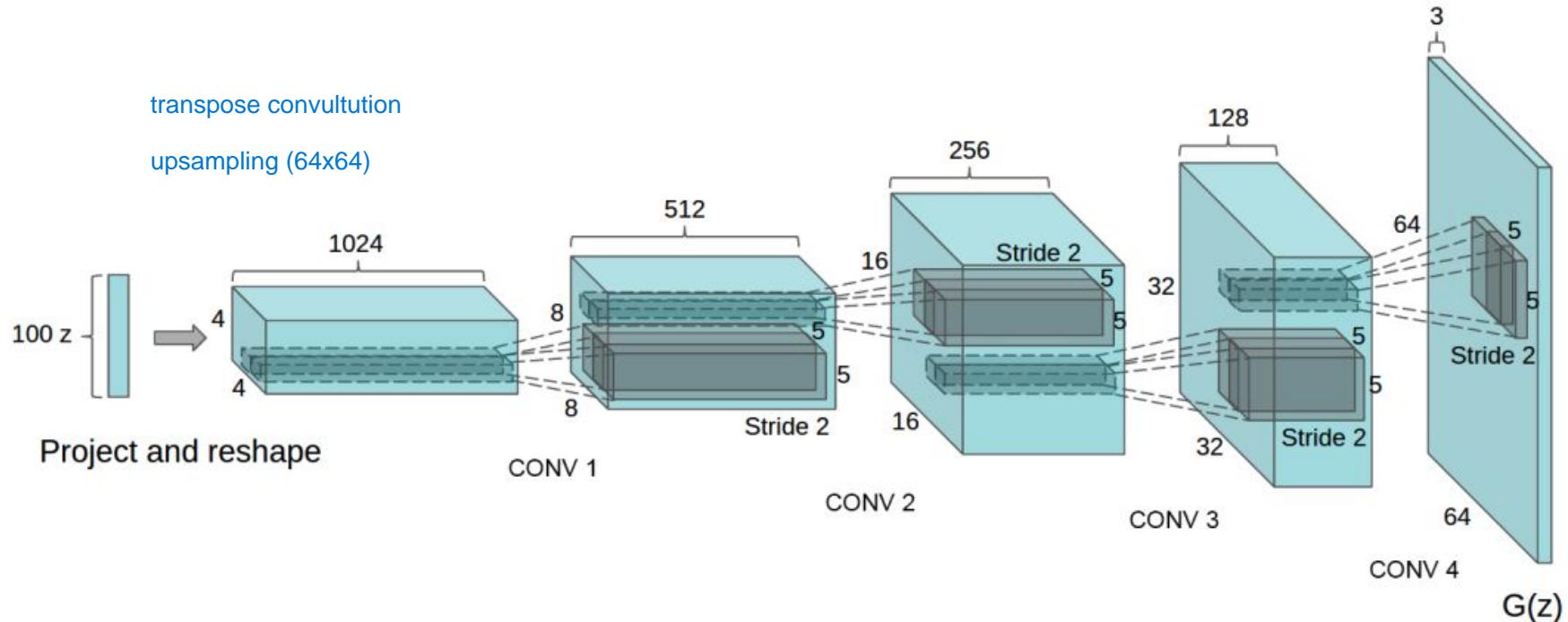
Results



Nearest neighbor from training set

Generative Adversarial Network (GAN)

Variations - DC GAN



Generative Adversarial Network (GAN)

Variations - DC GAN : Results



Generative Adversarial Network (GAN)

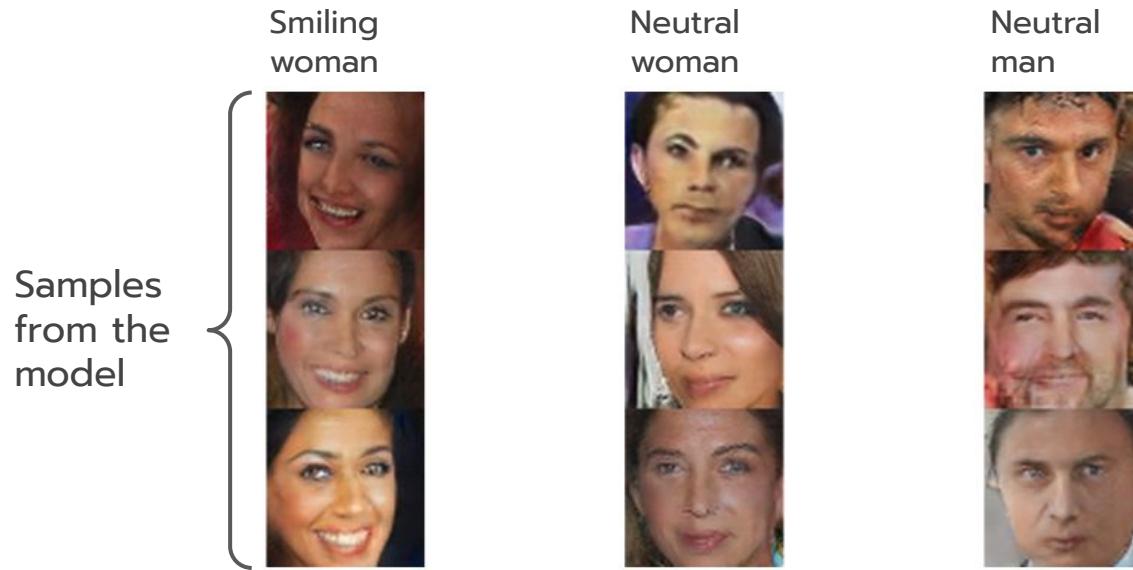
Variations - DC GAN : Results

Interpolating
between points in
latent z space,
show that the
space learned has
smooth
transitions



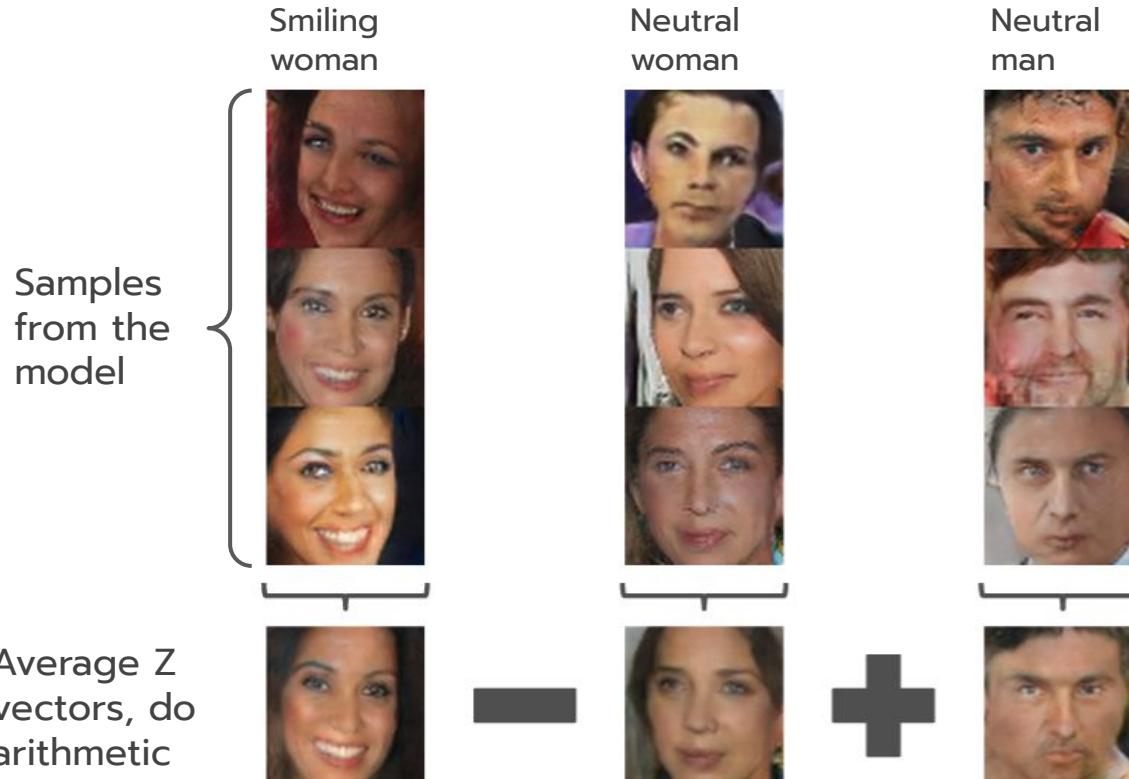
Generative Adversarial Network (GAN)

Variations - DC GAN : Vector Math



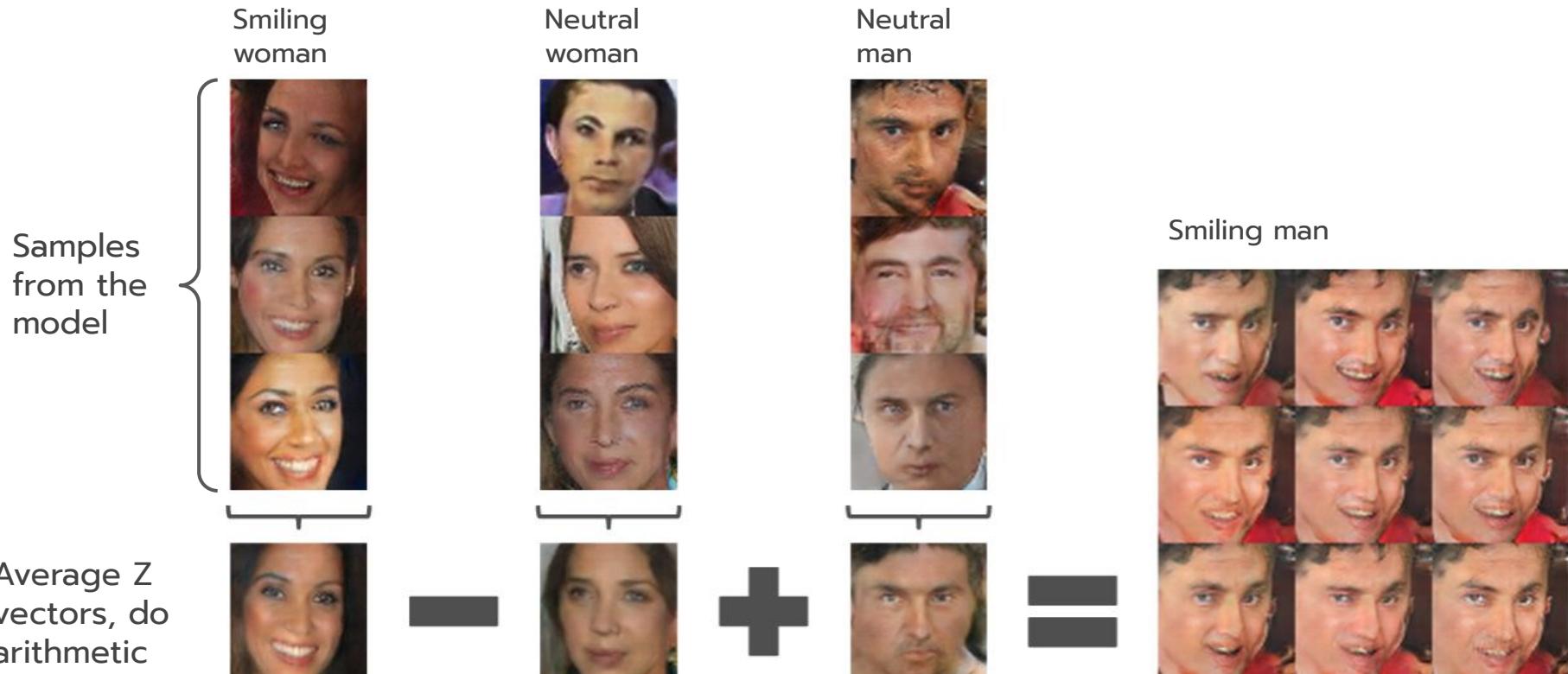
Generative Adversarial Network (GAN)

Variations - DC GAN : Vector Math



Generative Adversarial Network (GAN)

Variations - DC GAN : Vector Math



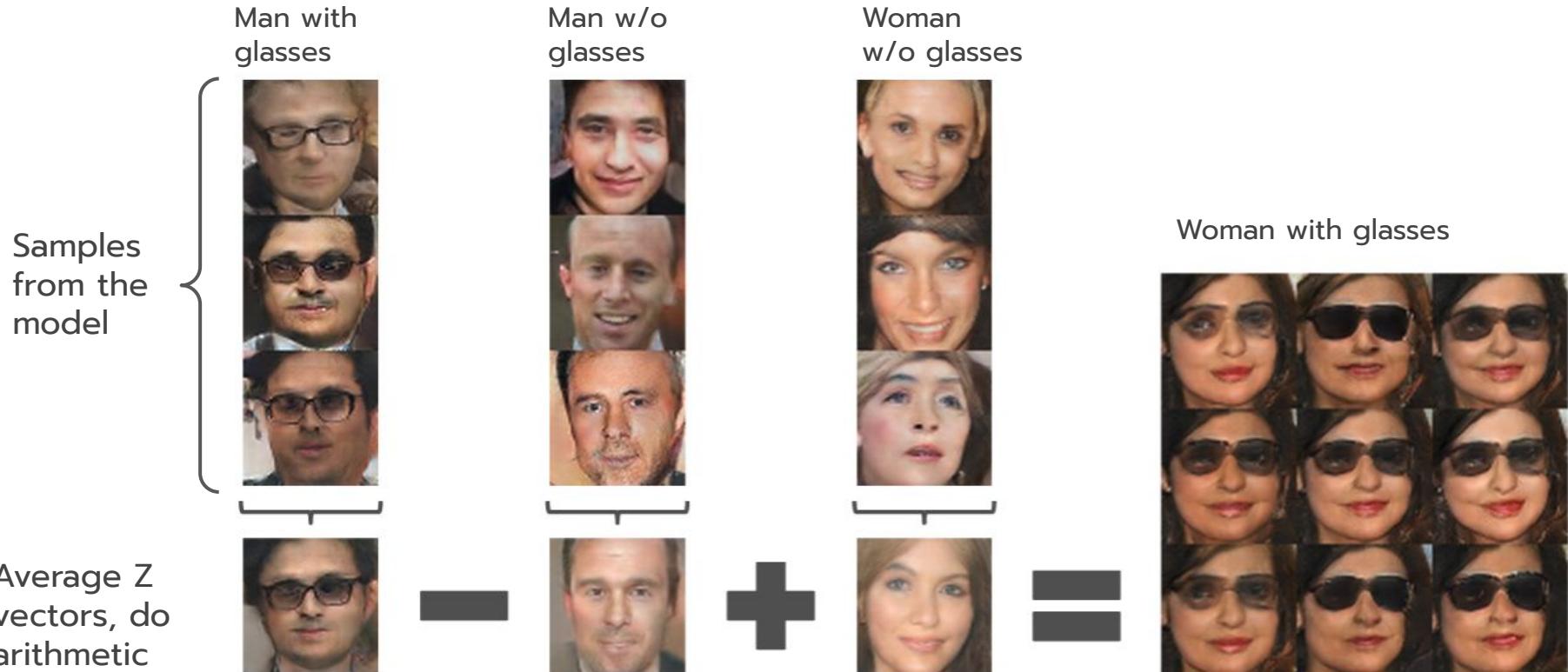
Generative Adversarial Network (GAN)

Variations - DC GAN : Vector Math



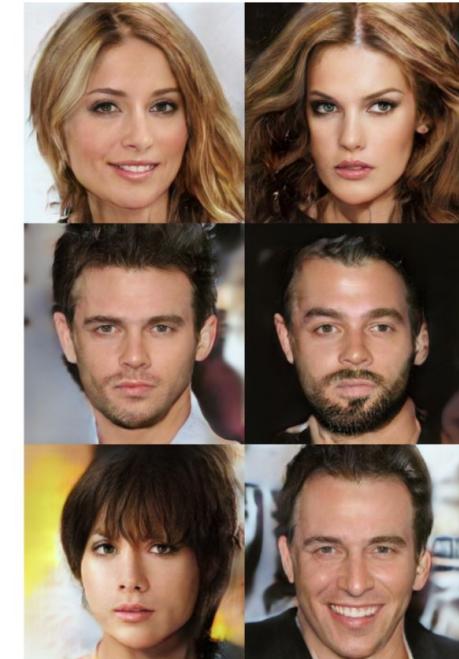
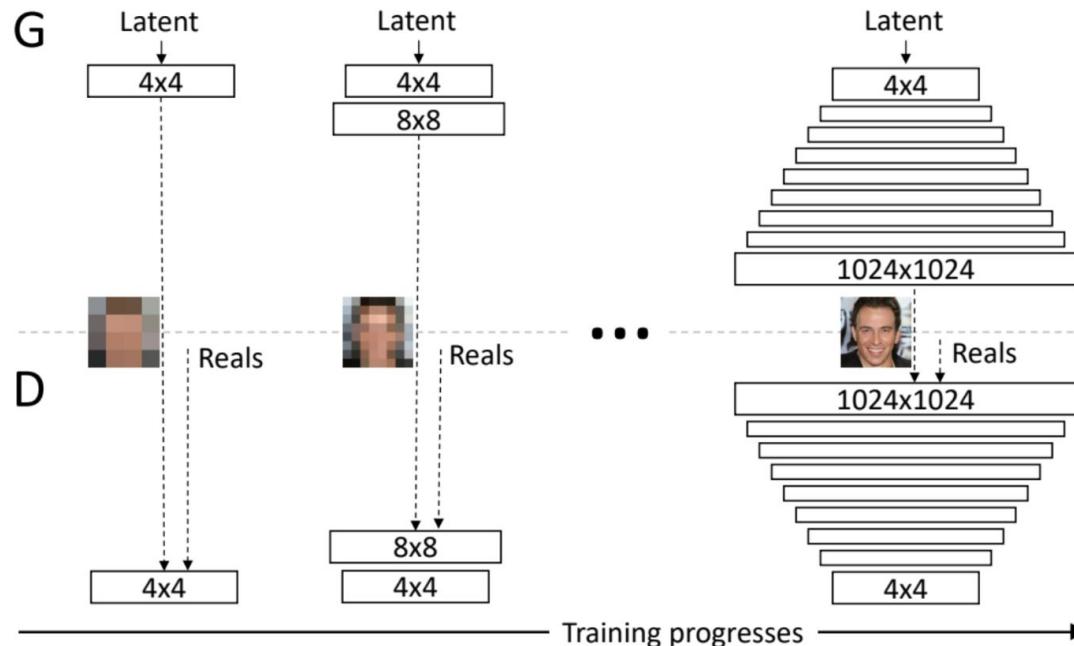
Generative Adversarial Network (GAN)

Variations - DC GAN : Vector Math



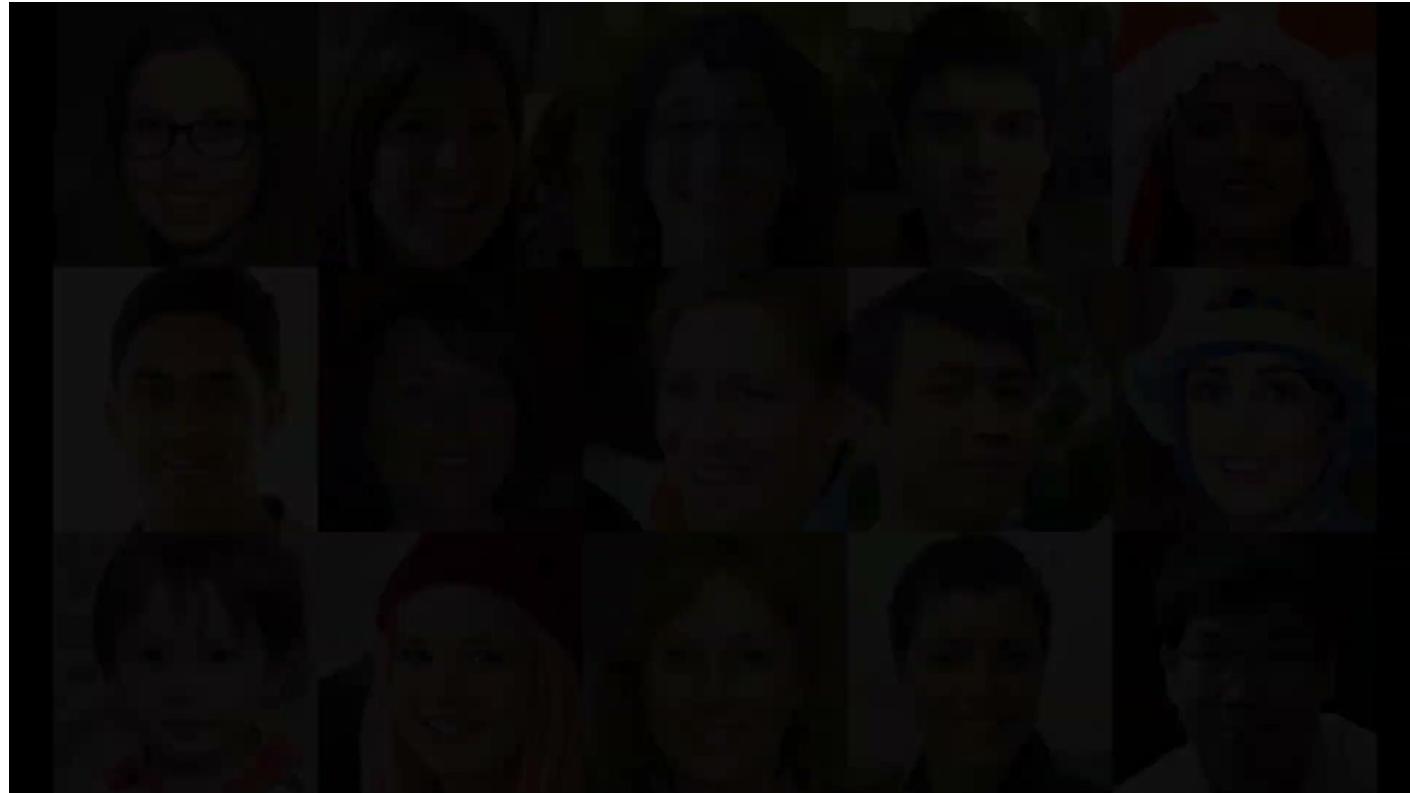
Generative Adversarial Network (GAN)

Variations - Progressively Growing GAN



Generative Adversarial Network (GAN)

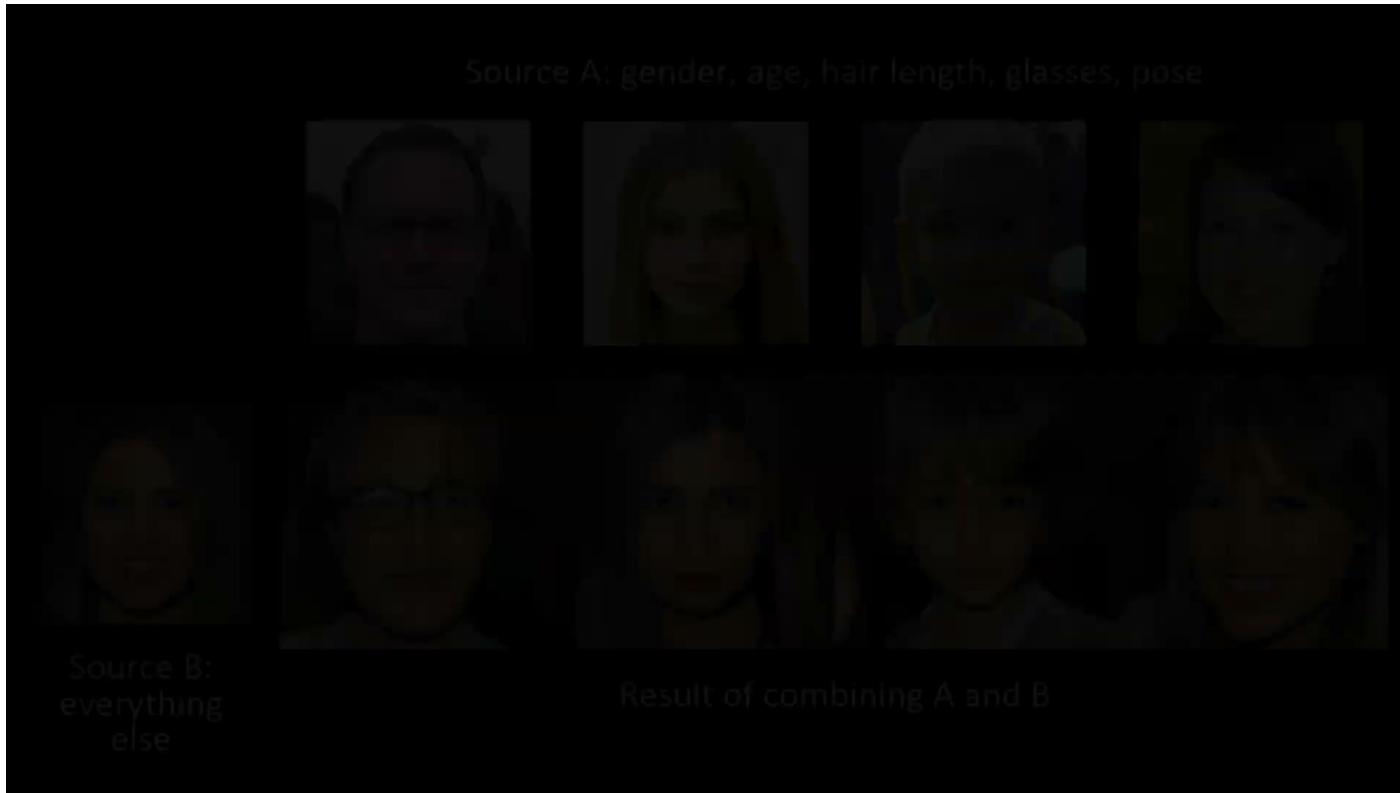
Variations - StyleGAN



Karras et al, "A Style-Based Generator Architecture for Generative Adversarial Networks", CVPR 2019
Video is licensed under CC BY-NC 4.0.
Source: https://drive.google.com/drive/folders/1NFO7_vH0t98J13ckJYFd7kuaTkyeRJ86

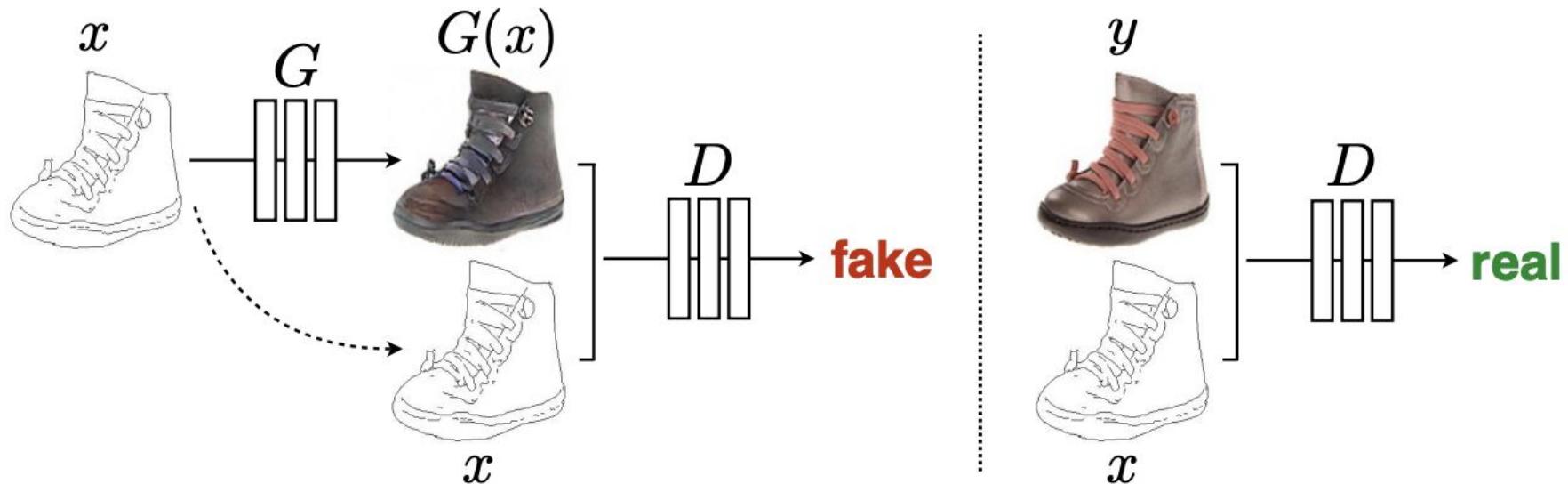
Generative Adversarial Network (GAN)

Variations - StyleGAN



Generative Adversarial Network (GAN)

Variations - Conditional GAN



Generative Adversarial Network (GAN)

Variations - Image-to-Image Translation: Pix2Pix

Labels to Street Scene

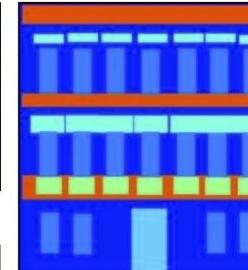


input



output

Labels to Facade



input



output

BW to Color



input

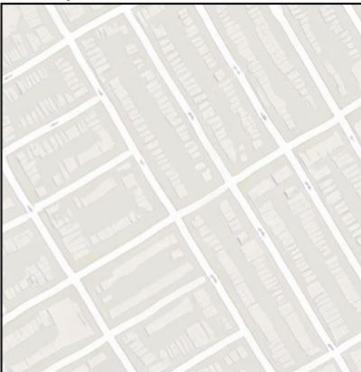


output

Aerial to Map



input



output

Day to Night



input



output

Edges to Photo



input



output

Generative Adversarial Network (GAN)

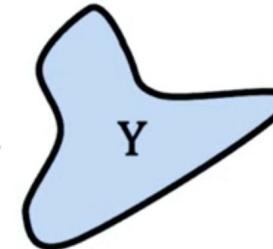
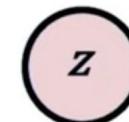
Variations - Cycle GAN

Distribution transformations.

GANs:

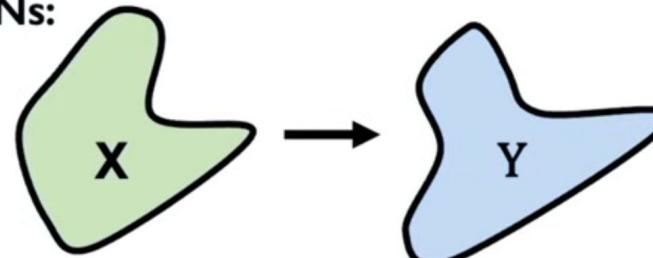
Gaussian noise

$$z \sim N(0,1)$$



Gaussian noise → target data manifold

CycleGANs:

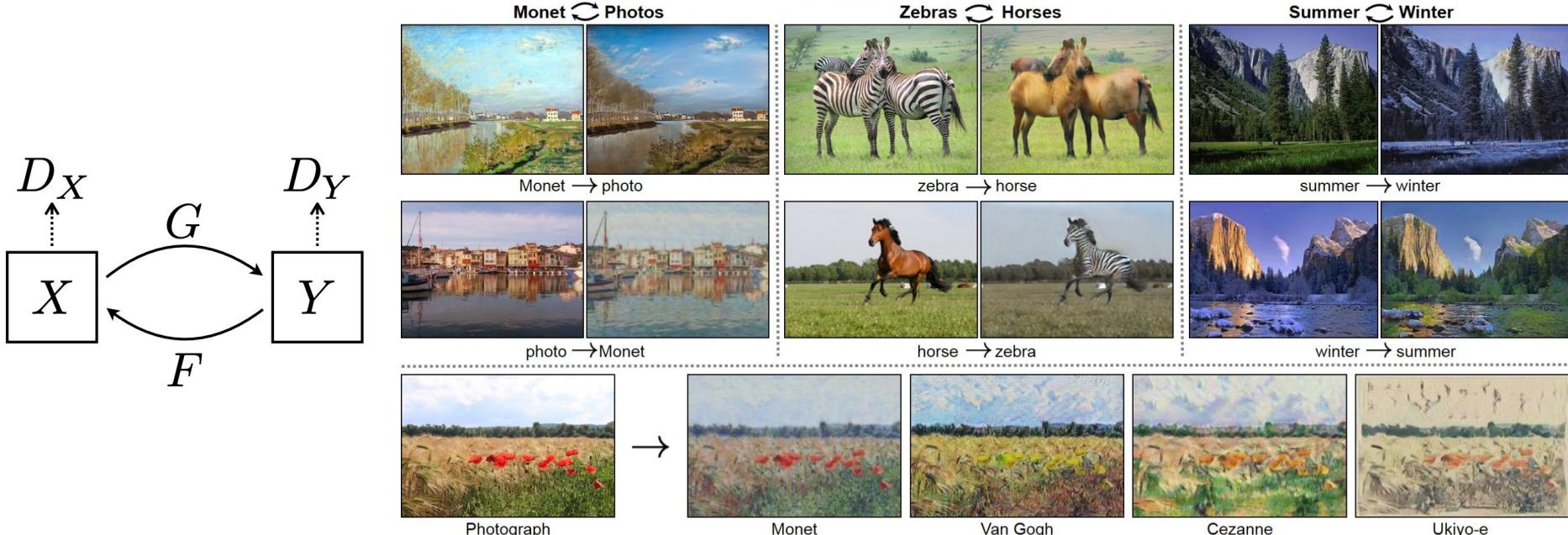


data manifold X → data manifold Y

Generative Adversarial Network (GAN)

Variations - Cycle GAN

Domain transformation: CycleGAN learns transformations across domains with unpaired data.



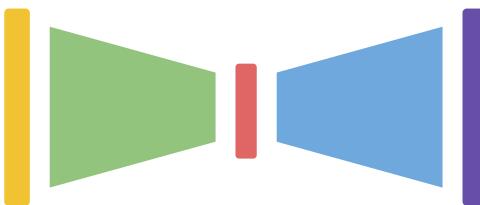
Generative Adversarial Network (GAN)

Variations - Cycle GAN

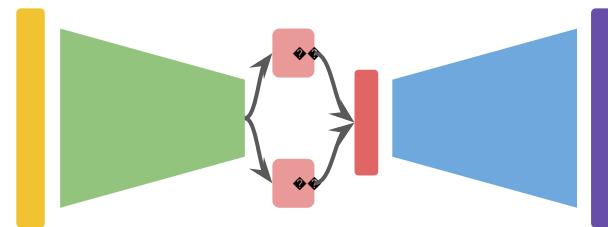


Summarize

Autoencoder (AE)



Variational Autoencoder (VAE)



Generative Adversarial Network (GAN)

