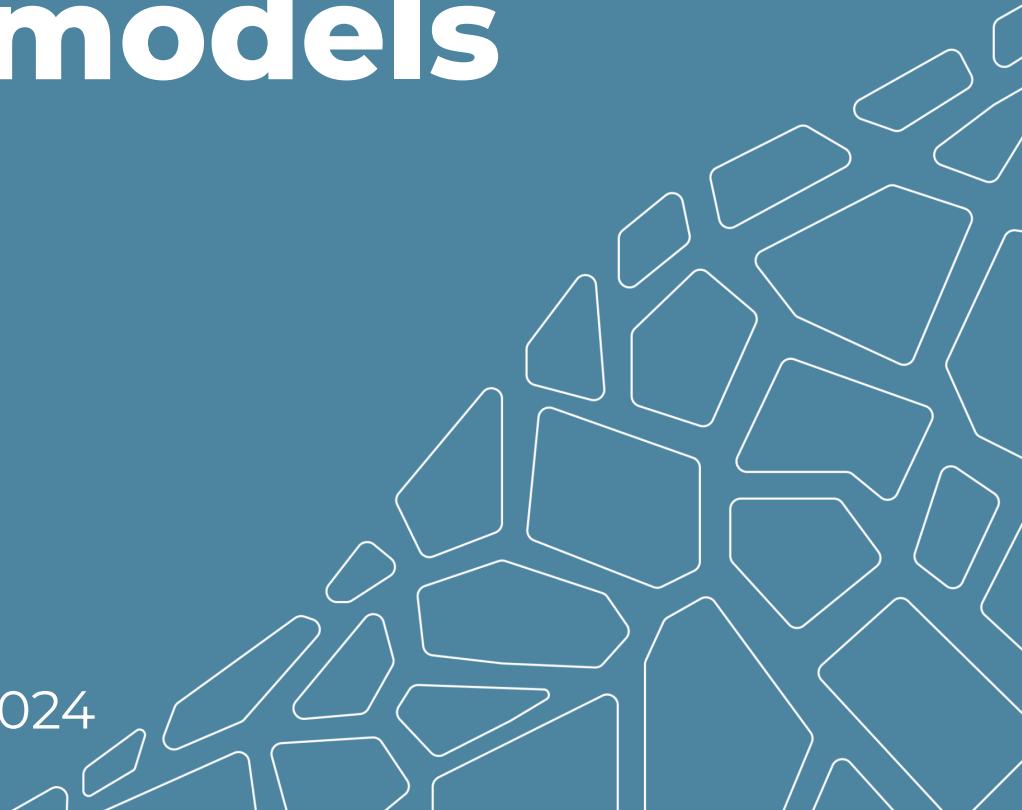


Generative models

Vladislav Goncharenko



MIPT, spring 2024



Outline

- Problem statement: probabilistic view
- Generative models overview
- Autoencoders (not generative models!)
- Variational Autoencoders (VAE)
 - Conditional VAE
- Generative Adversarial Networks
 - Conditional 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.

Classification

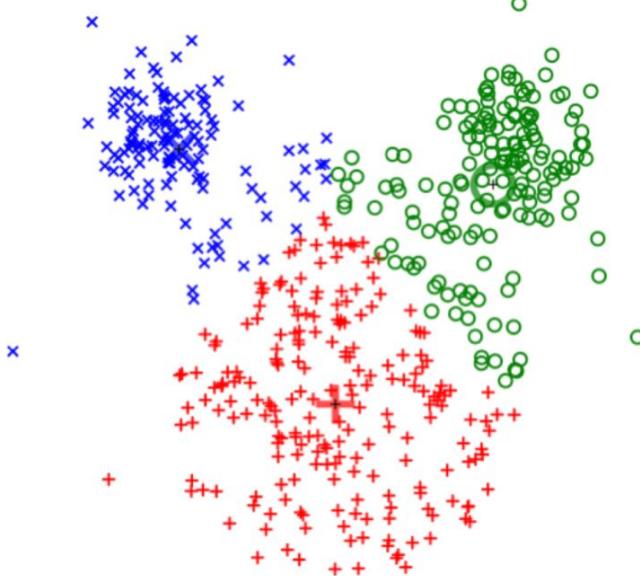


Cat

[This image](#) is CC0 public domain

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.

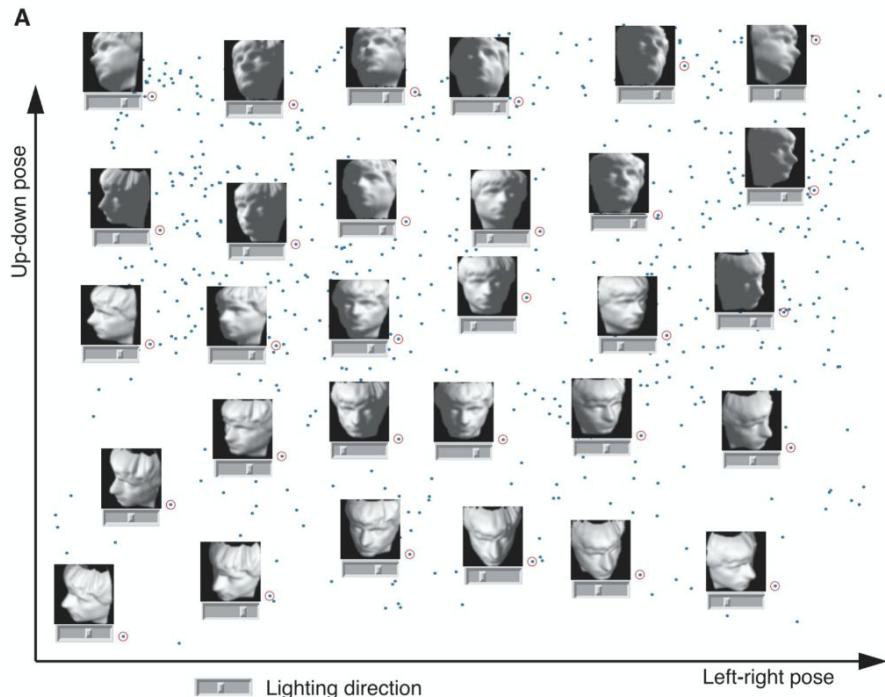


Manifold assumption

The data lie approximately on a surface (called manifold) of usually much lower dimension than the input space

So problem dimensionality could be (non-)linearly reduced or other tasks solved

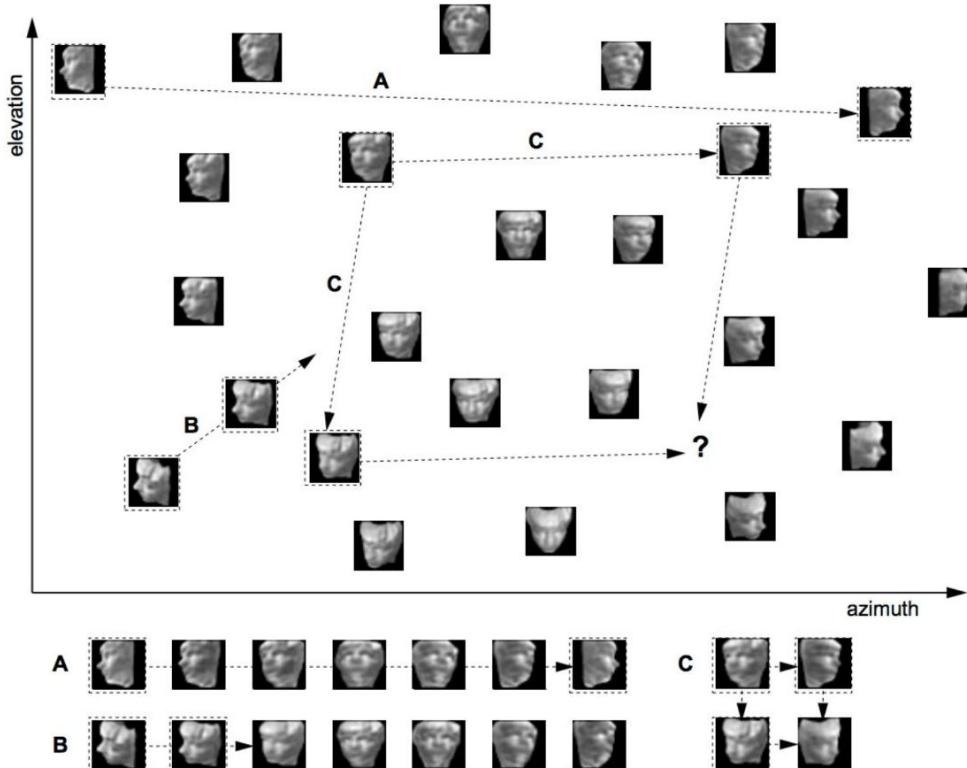
Sometimes dimensionality of manifold is referred as [intrinsic dimension](#) (see [this article](#))



[Tenenbaum, de Silva, Langford](#)
[A Global Geometric Framework for Nonlinear Dimensionality Reduction](#)



Latent space



Latent (embedding) space describes data in coordinates more relevant to humans' reason and often allows useful linear operations:

- Interpolation (A)
- Extrapolation (B)
- Analogy (C)

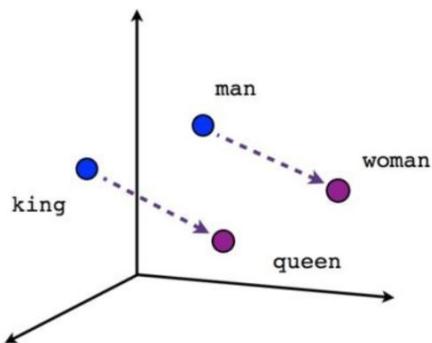
This process is also called embedding space 'walking'



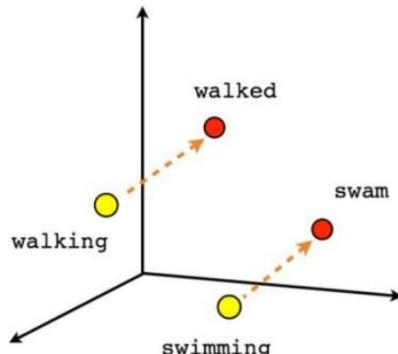
Latent space example

Word2vec is a method to embed words from text corpus into linear space

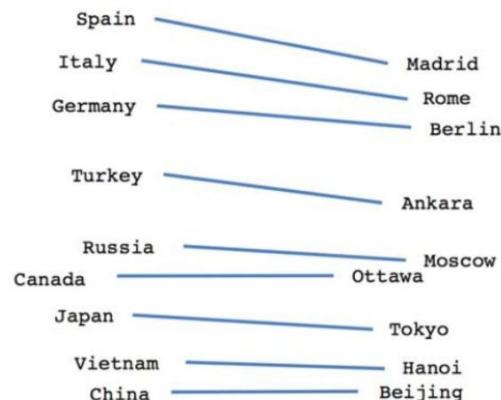
Read more: [manifold assumption](#), [assessing assumption](#)



Male-Female



Verb tense



Country-Capital

Discriminative vs Generative Models

Discriminative Model:

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

Data: x



Generative Model:

Learn a probability distribution $p(x)$

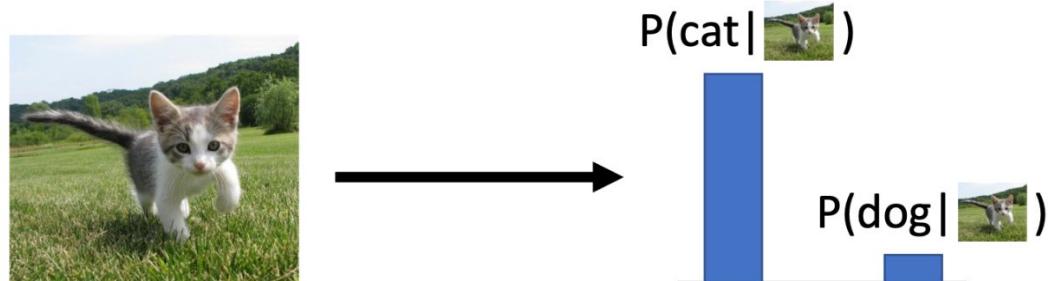
Conditional Generative Model: Learn $p(x|y)$

Label: y

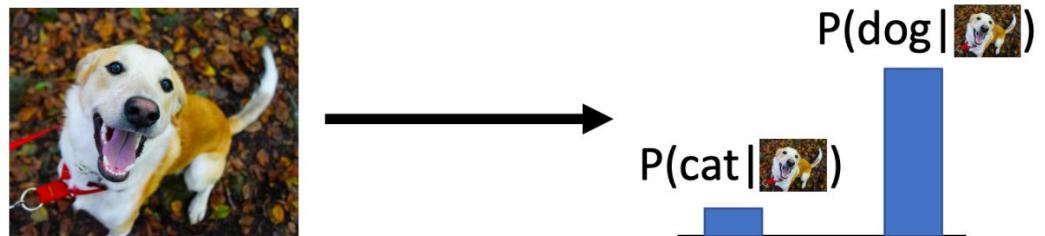
Cat

Discriminative vs Generative Models

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



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

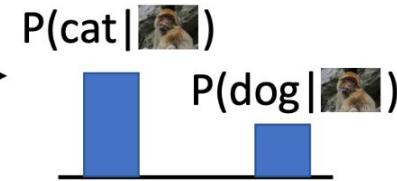


Conditional Generative Model: Learn $p(x|y)$

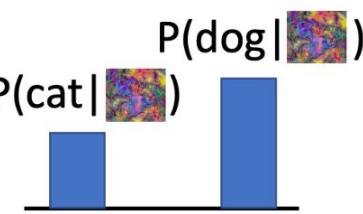
Discriminative model: the possible labels for each input "compete" for probability mass.
But no competition between **images**

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 $p(x|y)$

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

Monkey image is CC0 Public Domain
Abstract image is free to use under the [Pixabay license](#)

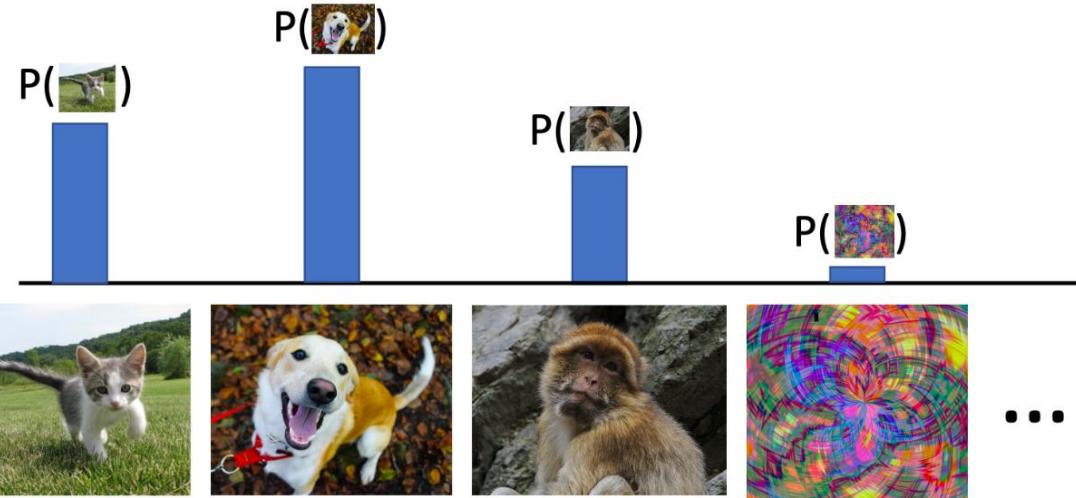
Discriminative vs Generative Models

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

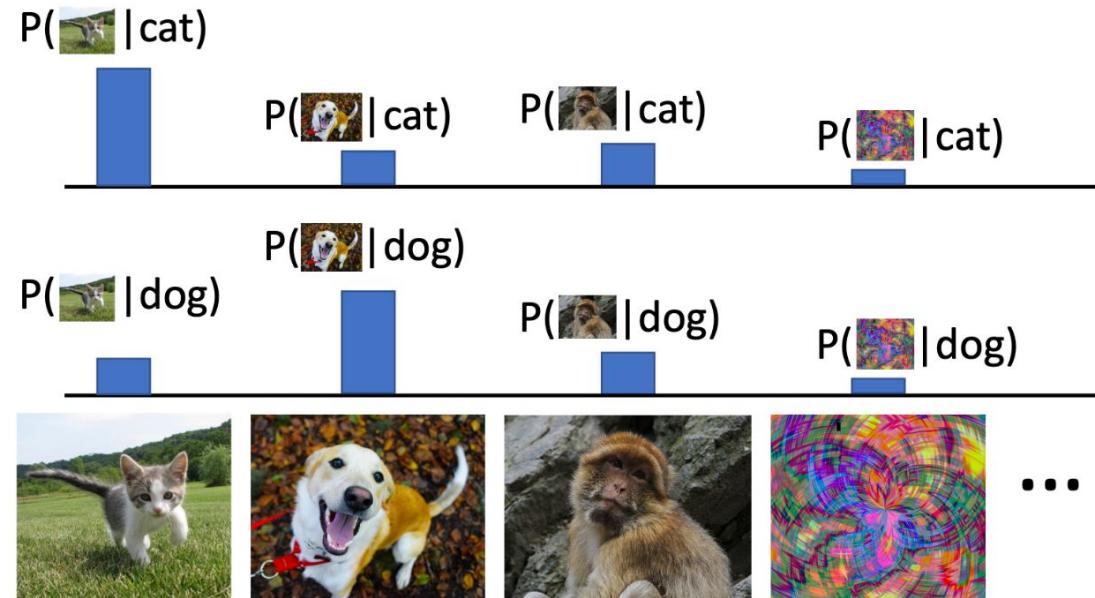
Model can “reject” unreasonable inputs by assigning them small values

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:

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

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

We can build a conditional generative model from other components!

What can we do with a generative model?

Discriminative Model:

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



Assign labels to data
Feature learning (supervised)

Generative Model:

Learn a probability distribution $p(x)$



Detect outliers
Feature learning (unsupervised)
Sample to **generate** new data

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



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

Taxonomy of Generative Models

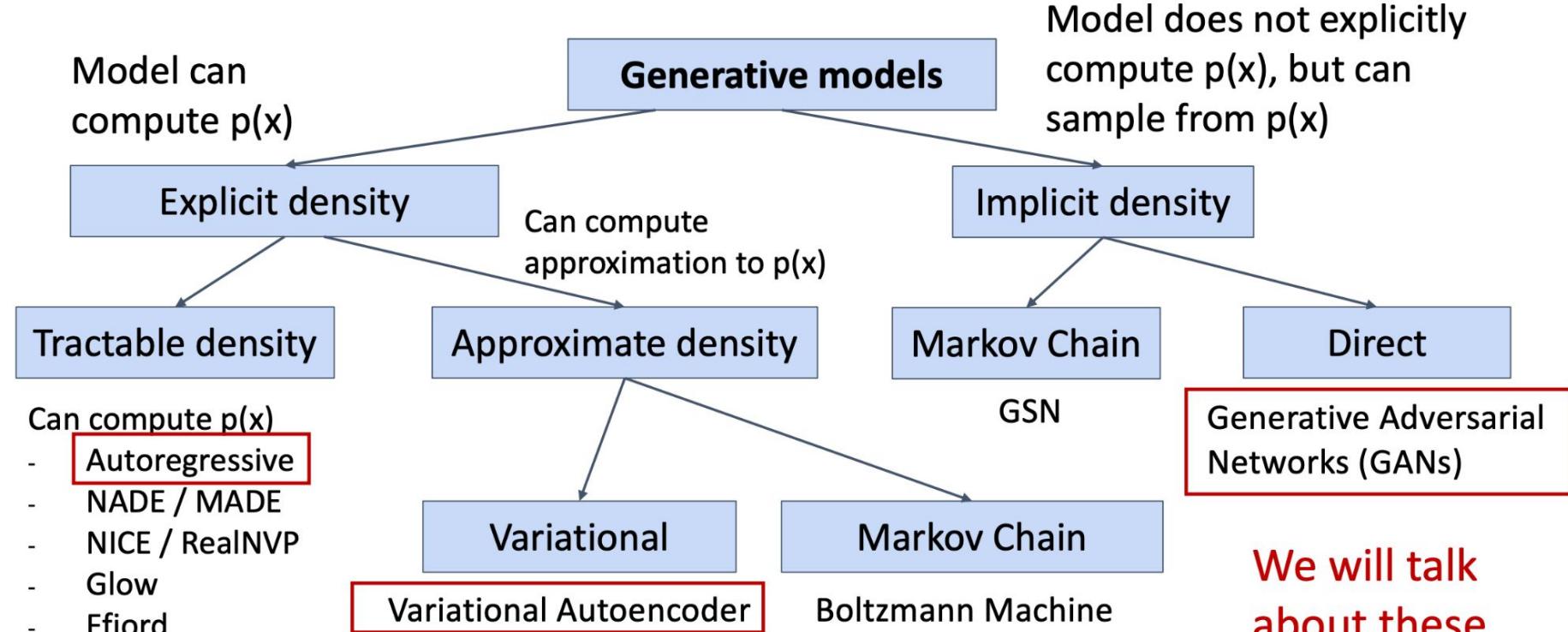


Figure adapted from Ian Goodfellow, Tutorial on Generative Adversarial Networks, 2017.

Autoencoders

Denote \mathbf{z} as encoded with encoder E input \mathbf{x}

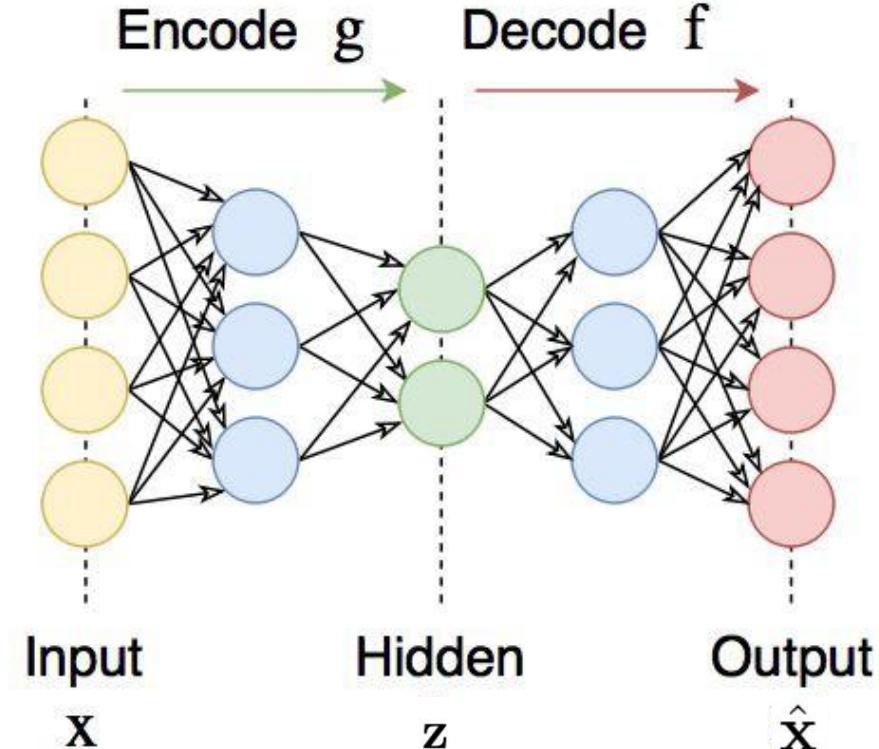
$$\mathbf{z} = E(\mathbf{x}, \theta_E)$$

Decoder D recovers \mathbf{x} from latent representation

$$\hat{\mathbf{x}} = D(\mathbf{z}, \theta_D)$$

Optimal parameters learned w.r.t. loss function L

$$[\theta_E, \theta_D] = \arg \min L(\hat{\mathbf{x}}, \mathbf{x})$$



Autoencoders

Denote \mathbf{z} as encoded with encoder E input \mathbf{x}

$$\mathbf{z} = E(\mathbf{x}, \theta_E)$$

Decoder D recovers \mathbf{x} from latent representation

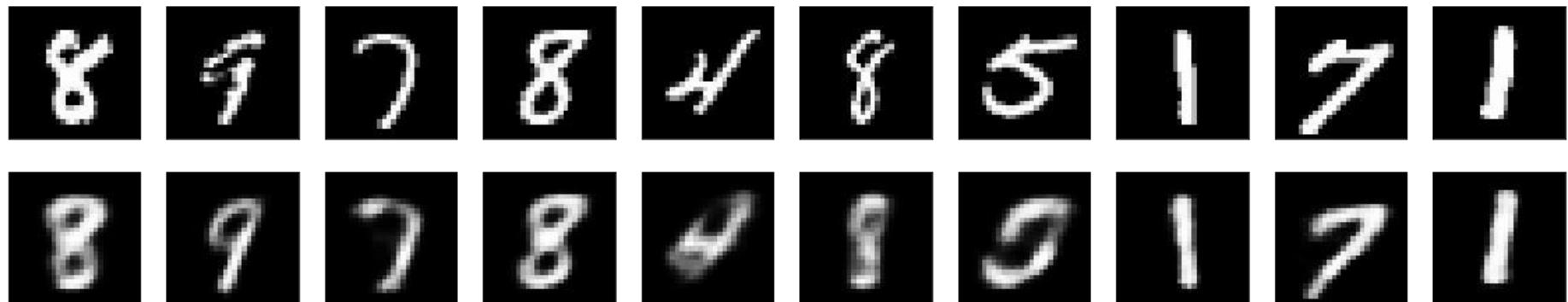
$$\hat{\mathbf{x}} = D(\mathbf{z}, \theta_D)$$

Simple example: PCA

Optimal parameters learned w.r.t. loss function L

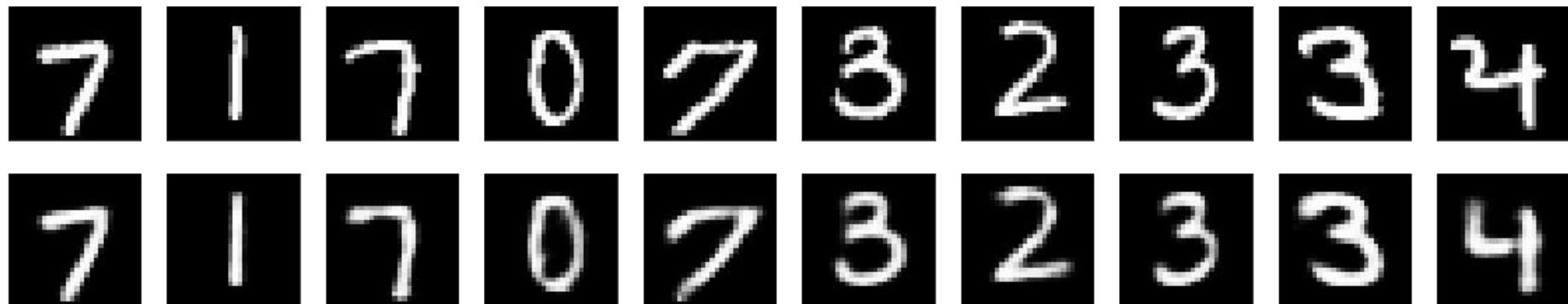
$$[\theta_E, \theta_D] = \arg \min L(\hat{\mathbf{x}}, \mathbf{x})$$

PCA performance on MNIST



16 components

Convolutional performance on MNIST

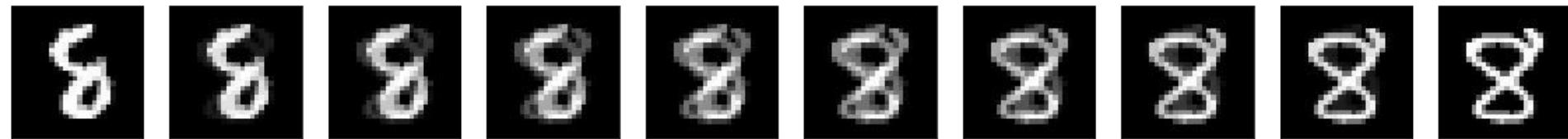


7 x 7 latent space

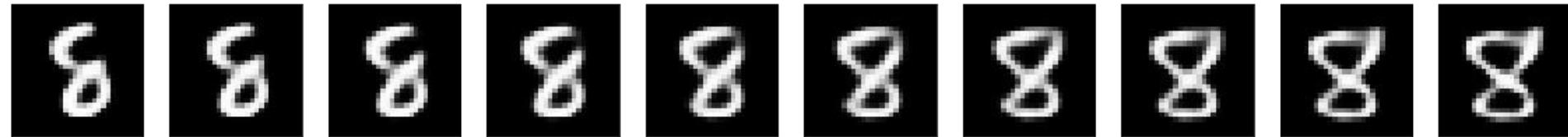
Homotopy between samples

10 steps between samples

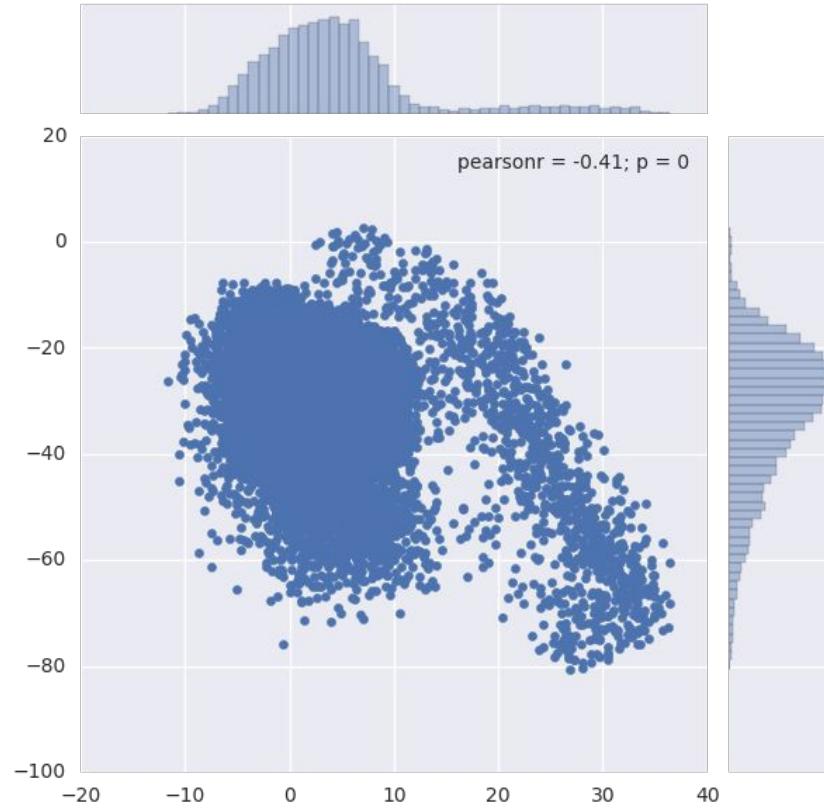
- In original feature space (28 x 28):



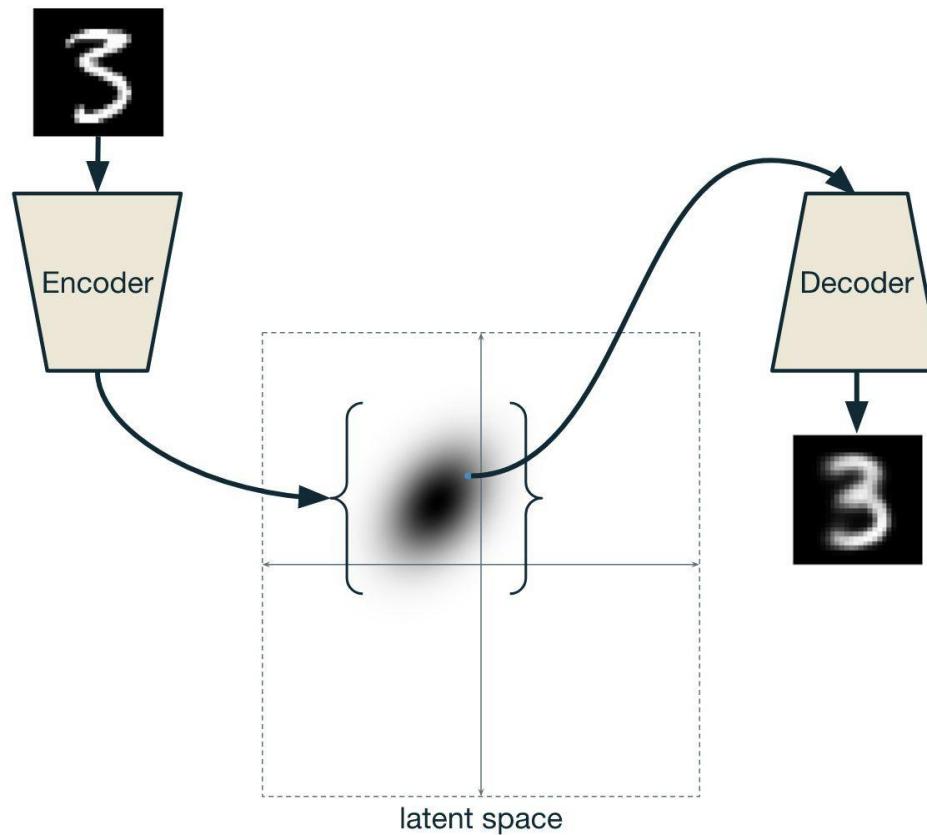
- In latent space (7 x 7):



Latent space structure



VAE intuition



KL divergence

Denote distributions $Q(z)$ and $P(z|X)$.

Kullback–Leibler divergence is defined as

$$\mathcal{D} [Q(z) \| P(z|X)] = E_{z \sim Q} [\log Q(z) - \log P(z|X)]$$

KL divergence

$$\mathcal{D} [Q(z) \| P(z|X)] = E_{z \sim Q} [\log Q(z) - \log P(z|X)]$$

Applying the Bayes rule:

$$\mathcal{D} [Q(z) \| P(z|X)] = E_{z \sim Q} [\log Q(z) - \log P(X|z) - \log P(z)] + \log P(X)$$

KL divergence

$$\mathcal{D} [Q(z) \| P(z|X)] = E_{z \sim Q} [\log Q(z) - \log P(z|X)]$$

Applying the Bayes rule:

$$\mathcal{D} [Q(z) \| P(z|X)] = E_{z \sim Q} [\log Q(z) - \log P(X|z) - \log P(z)] + \log P(X)$$

$$\log P(X) - \mathcal{D} [Q(z) \| P(z|X)] = E_{z \sim Q} [\log P(X|z)] - \mathcal{D} [Q(z) \| P(z)]$$

KL divergence

$$\mathcal{D} [Q(z) \| P(z|X)] = E_{z \sim Q} [\log Q(z) - \log P(z|X)]$$

Applying the Bayes rule:

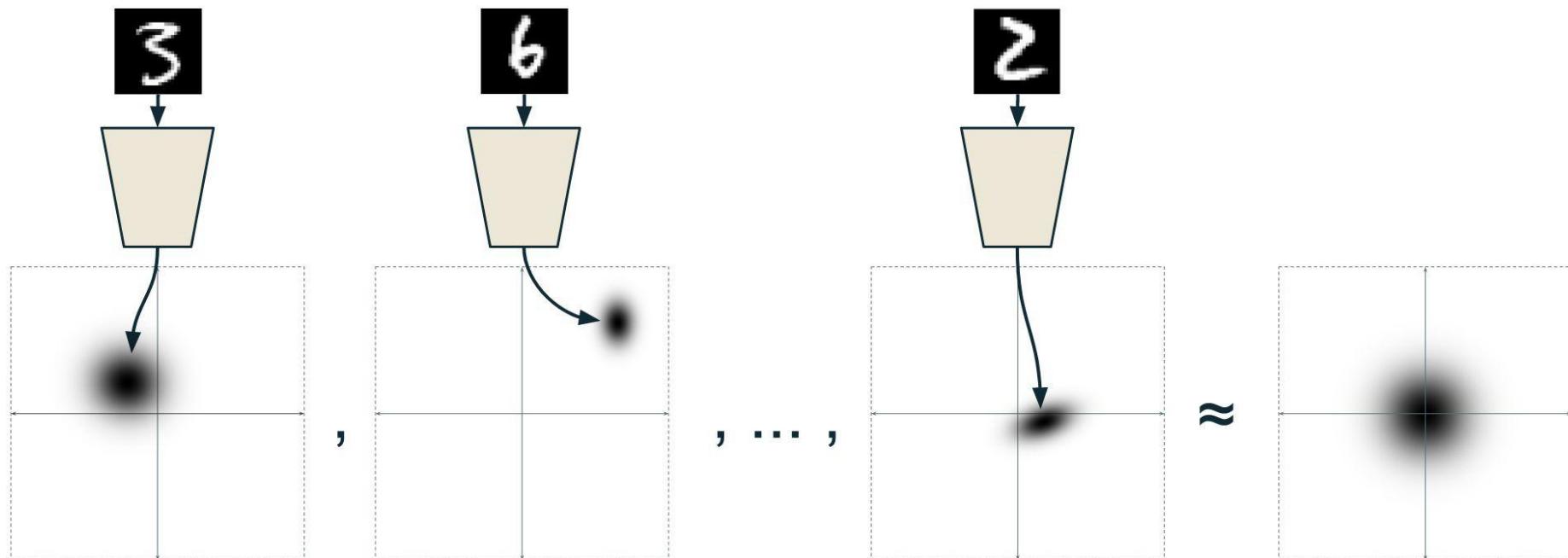
$$\mathcal{D} [Q(z) \| P(z|X)] = E_{z \sim Q} [\log Q(z) - \log P(X|z) - \log P(z)] + \log P(X)$$

$$\log P(X) - \mathcal{D} [Q(z) \| P(z|X)] = E_{z \sim Q} [\log P(X|z)] - \mathcal{D} [Q(z) \| P(z)]$$

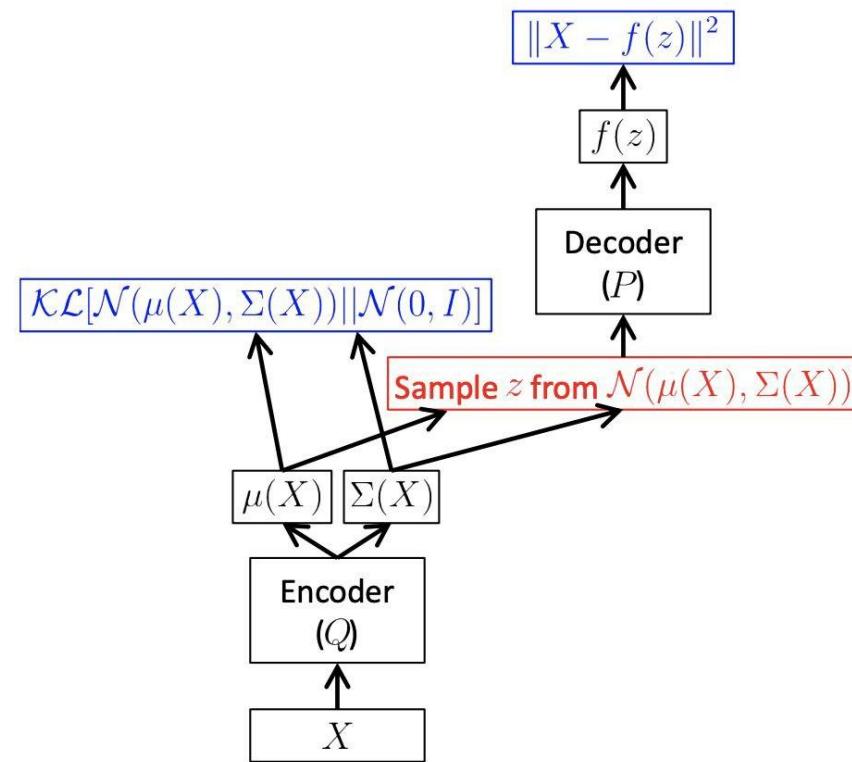
$$\boxed{\log P(X) - \mathcal{D} [Q(z|X) \| P(z|X)] = E_{z \sim Q} [\log P(X|z)] - \mathcal{D} [Q(z|X) \| P(z)]}$$

This equation is the core of Variational Autoencoders

Structure of the latent space



VAE so far

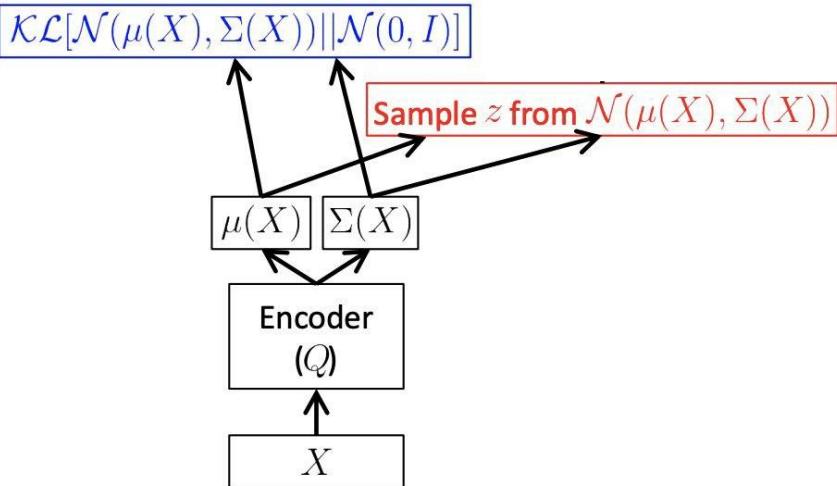


VAE so far

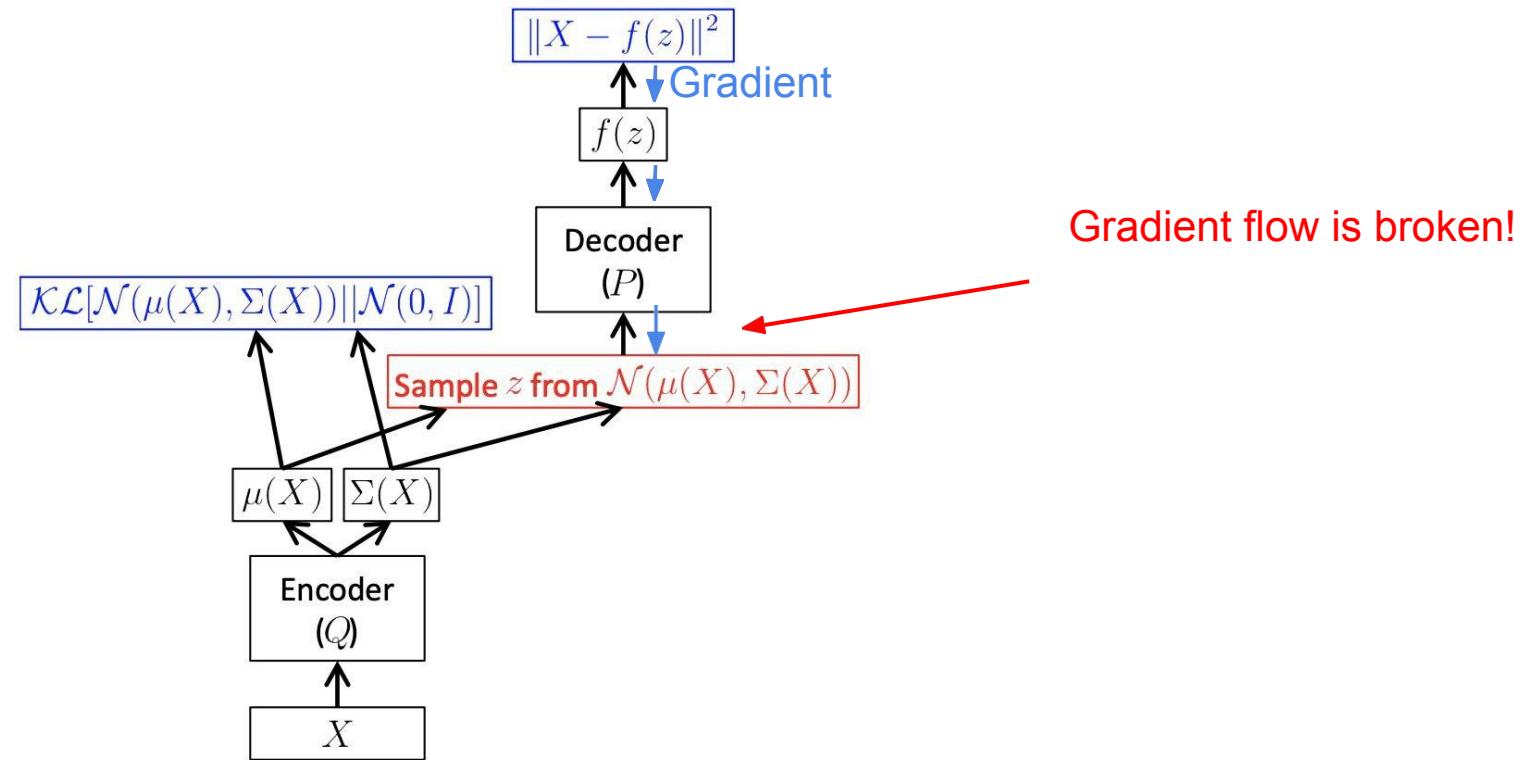
$$\mathcal{D}[\mathcal{N}(\mu(X), \Sigma(X)) || \mathcal{N}(0, I)] =$$

$$\frac{1}{2} \left(\text{tr}(\Sigma(X)) + (\mu(X))^\top (\mu(X)) - k - \log \det(\Sigma(X)) \right)$$

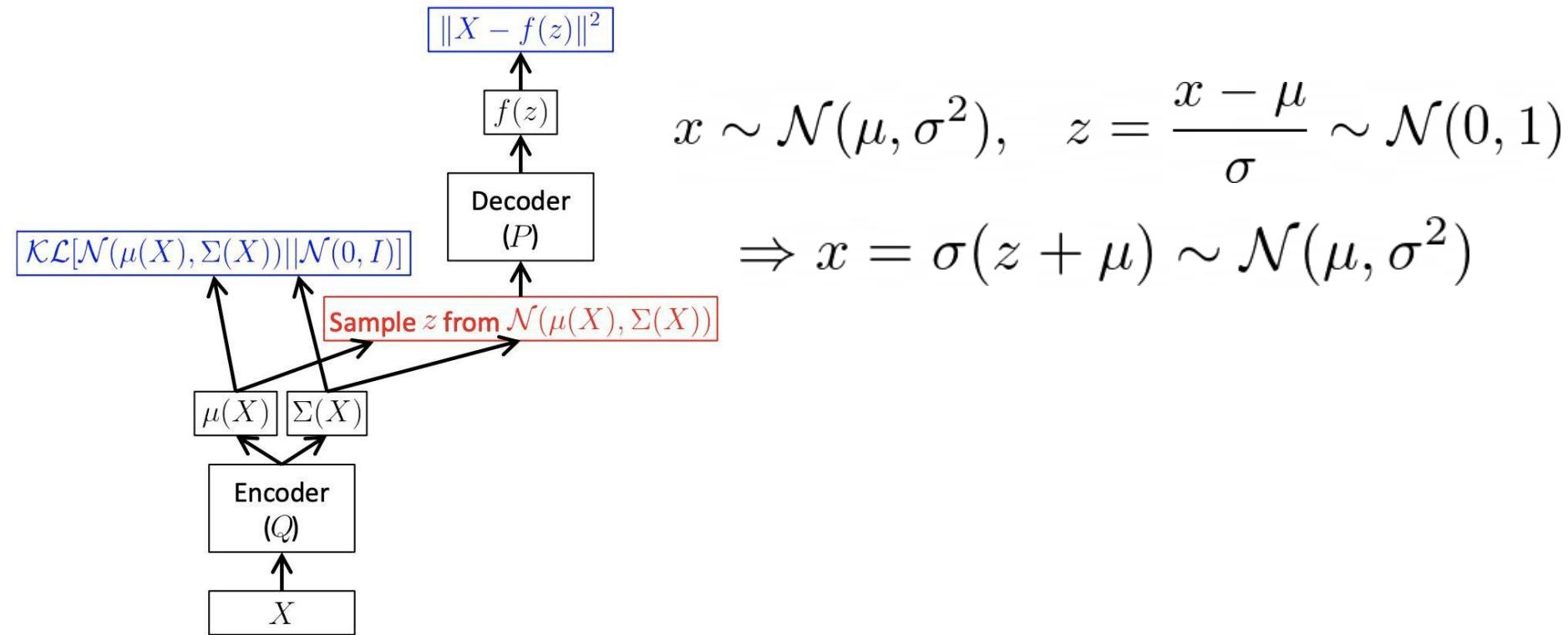
Try to derive it by yourself



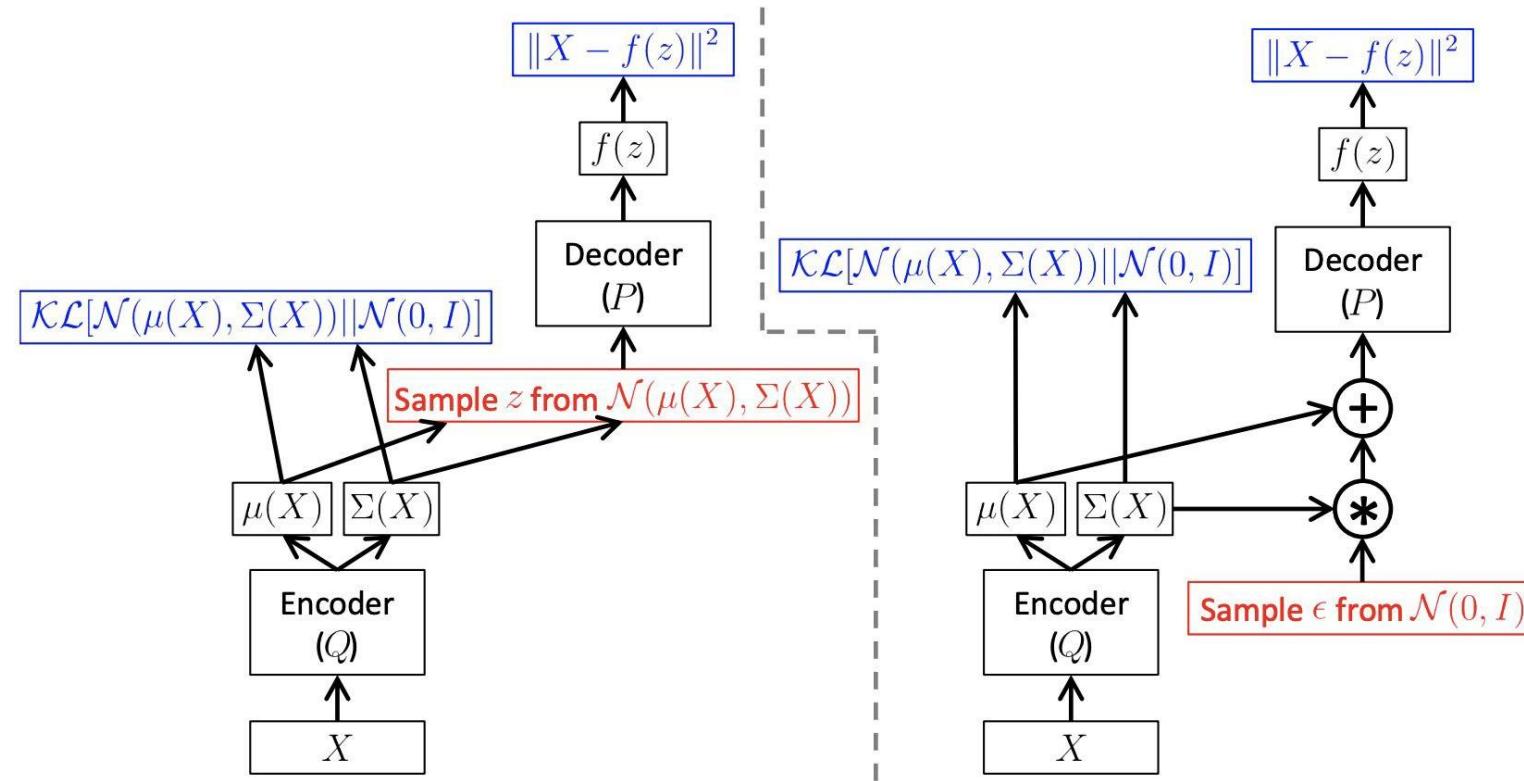
VAE so far



Reparametrization trick

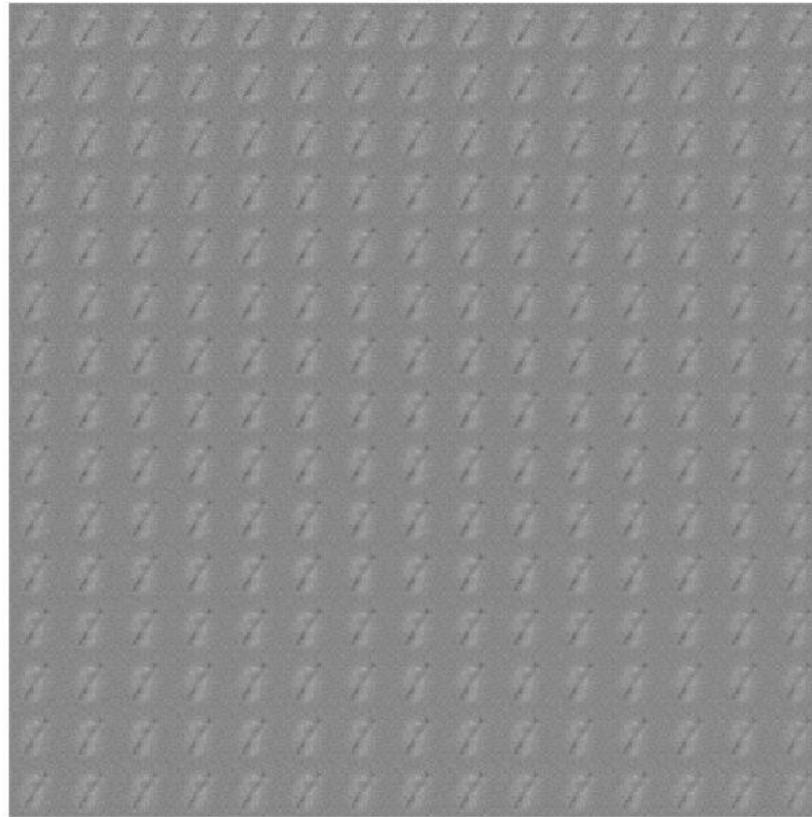


Reparametrization trick

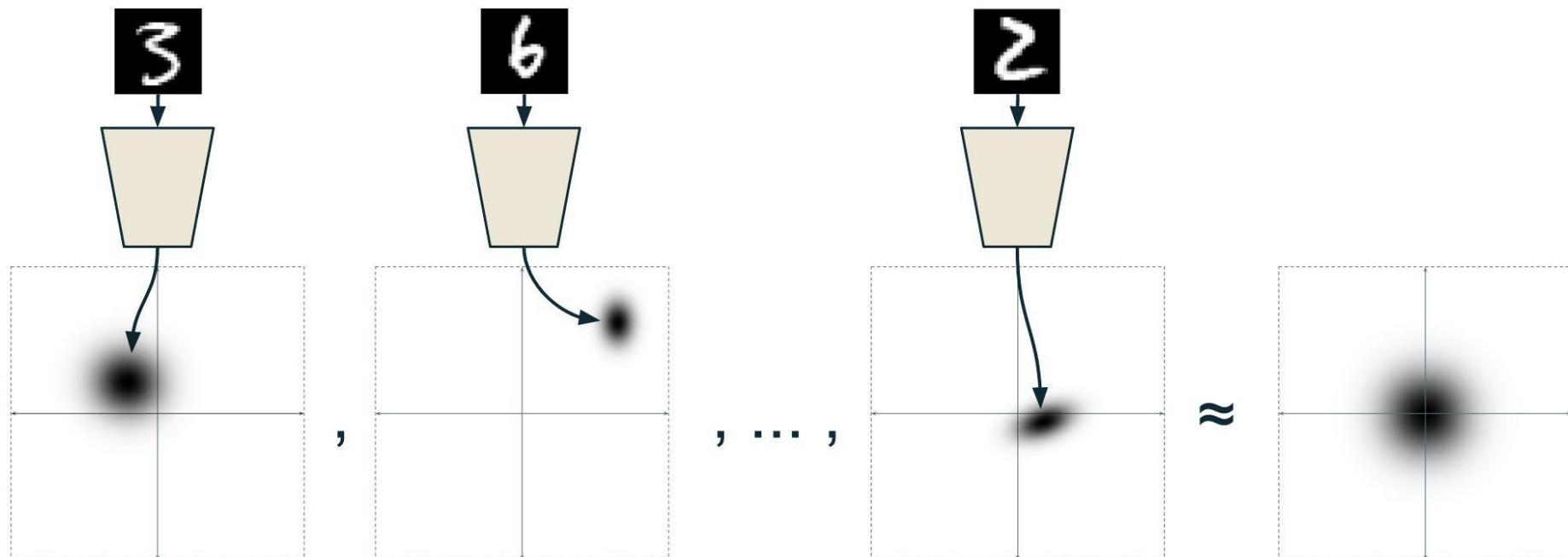


VAE manifold

Epoch: 0

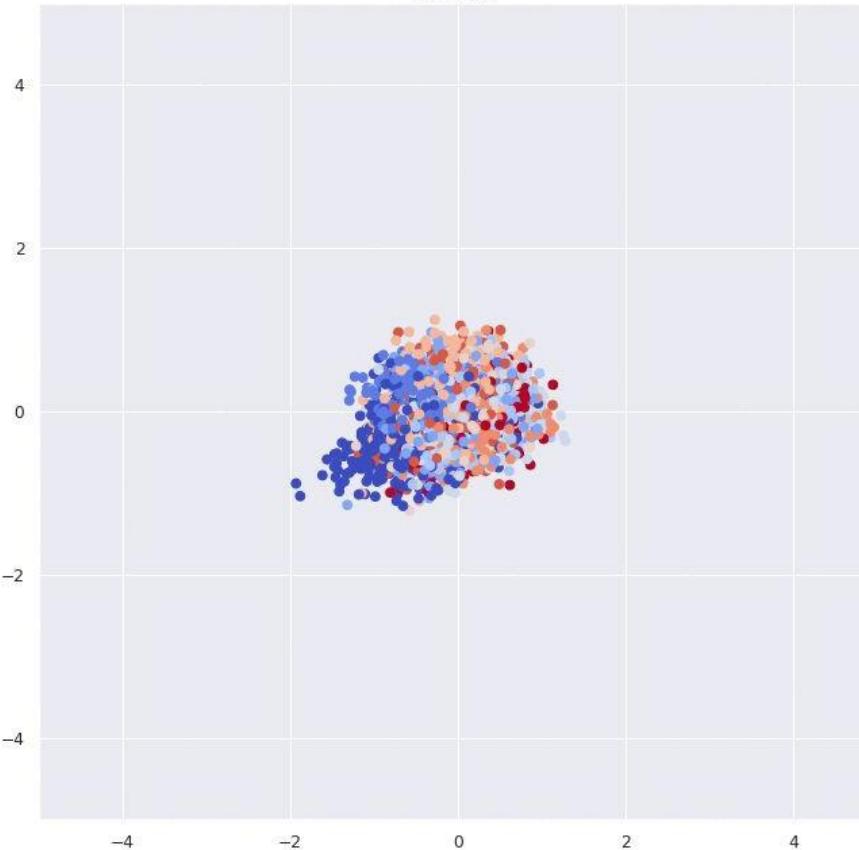


Structure of the latent space

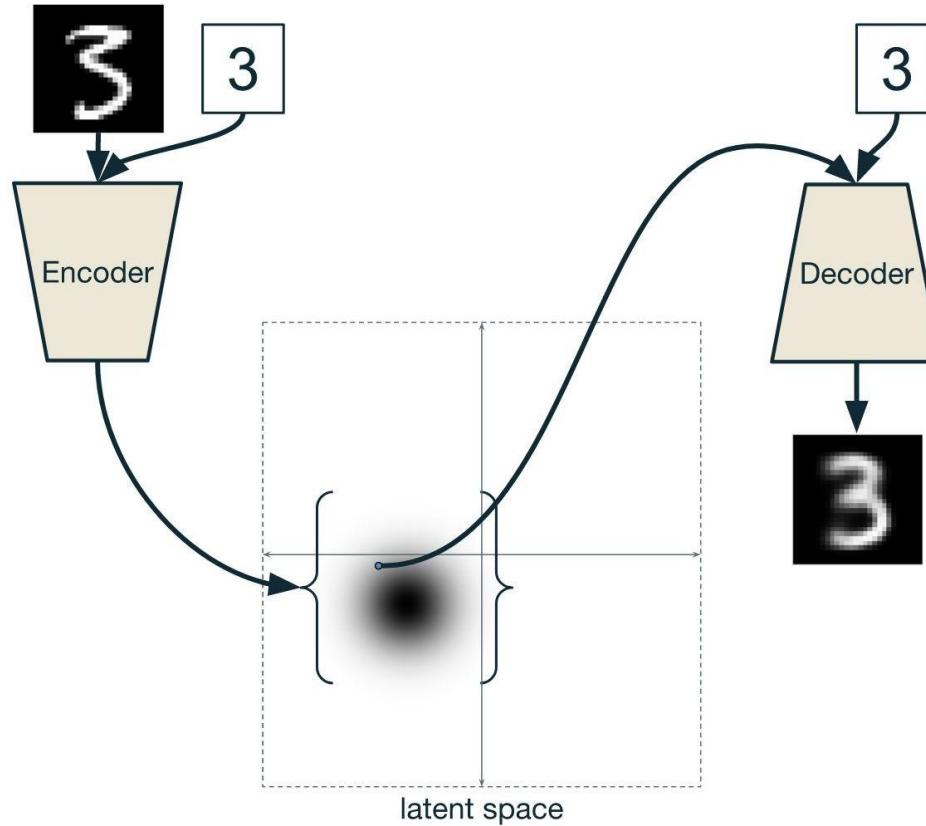


VAE latent space distribution

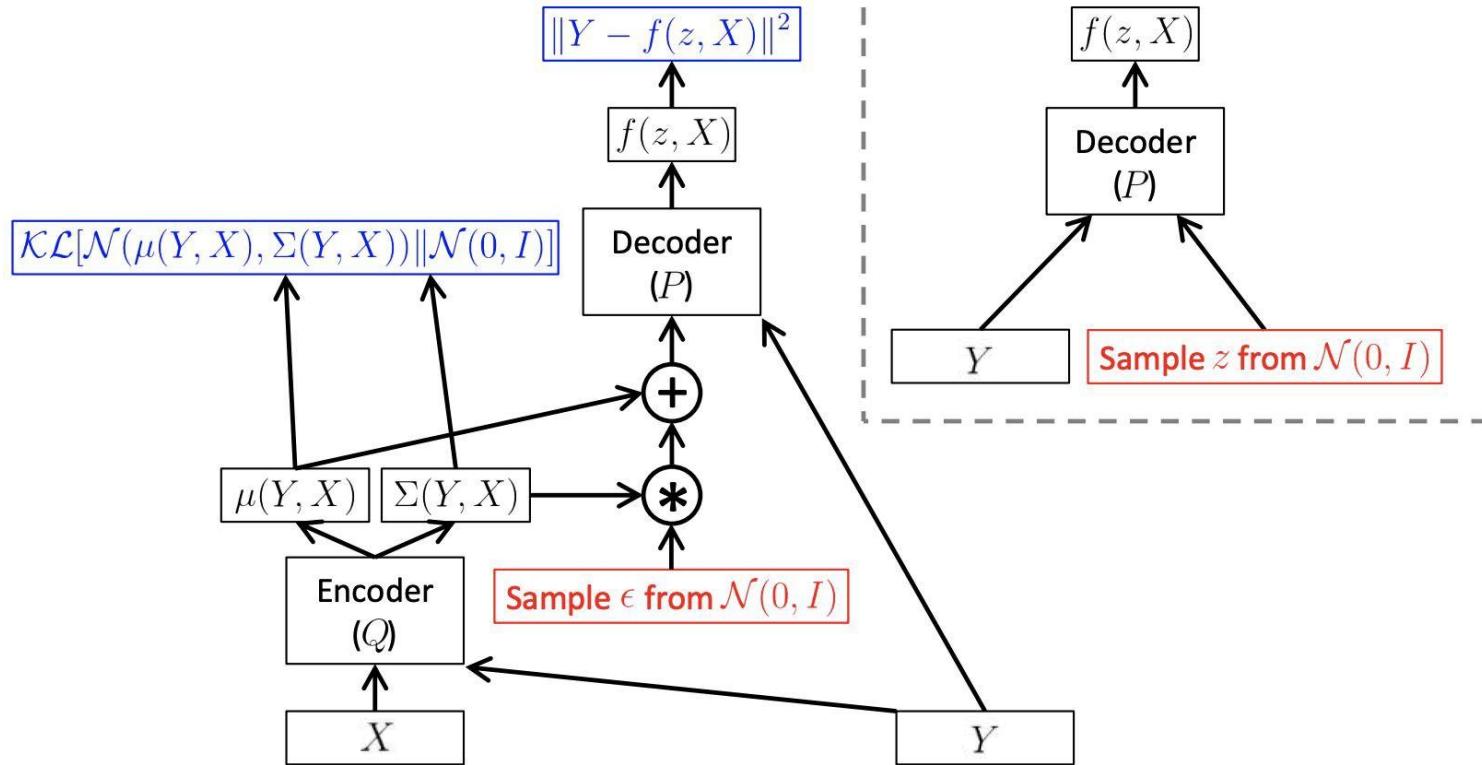
Epoch: 0



Conditional VAE intuition

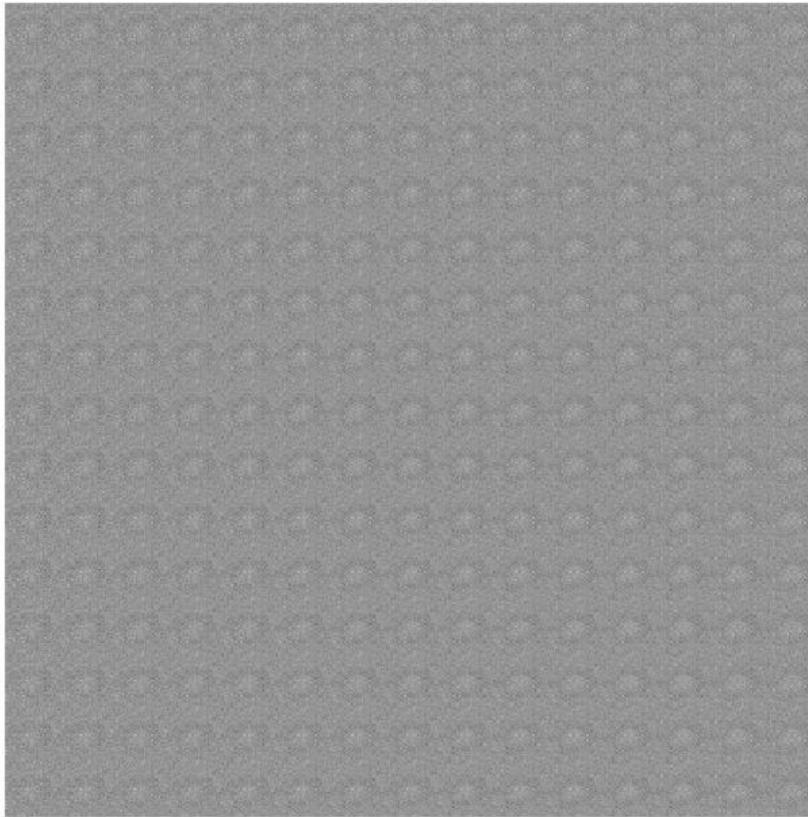


Conditional VAE



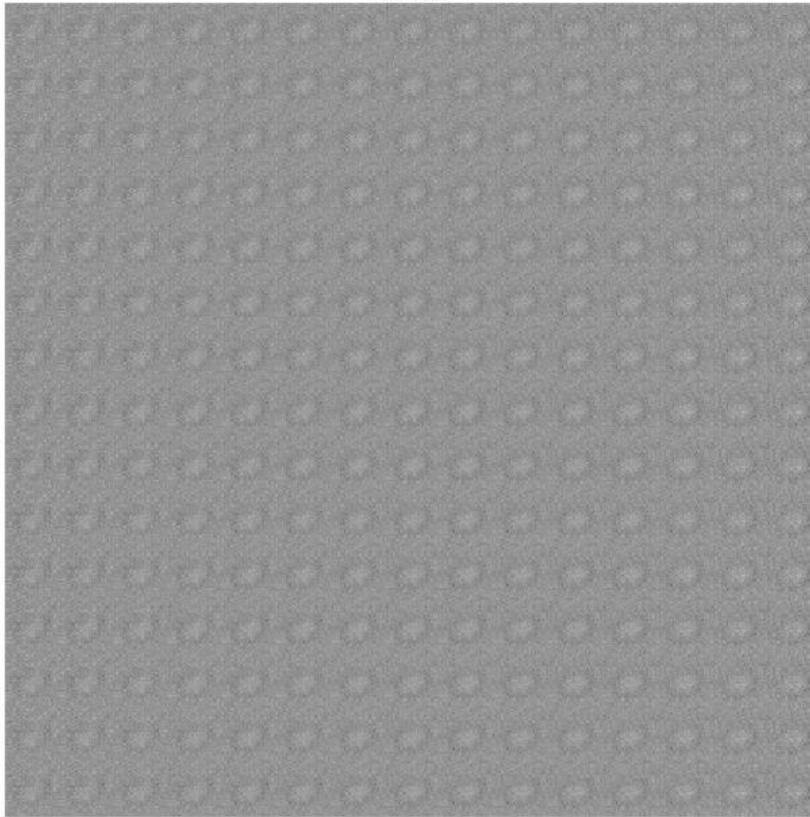
cVAE manifold

Epoch: 0



cVAE manifold

Epoch: 0

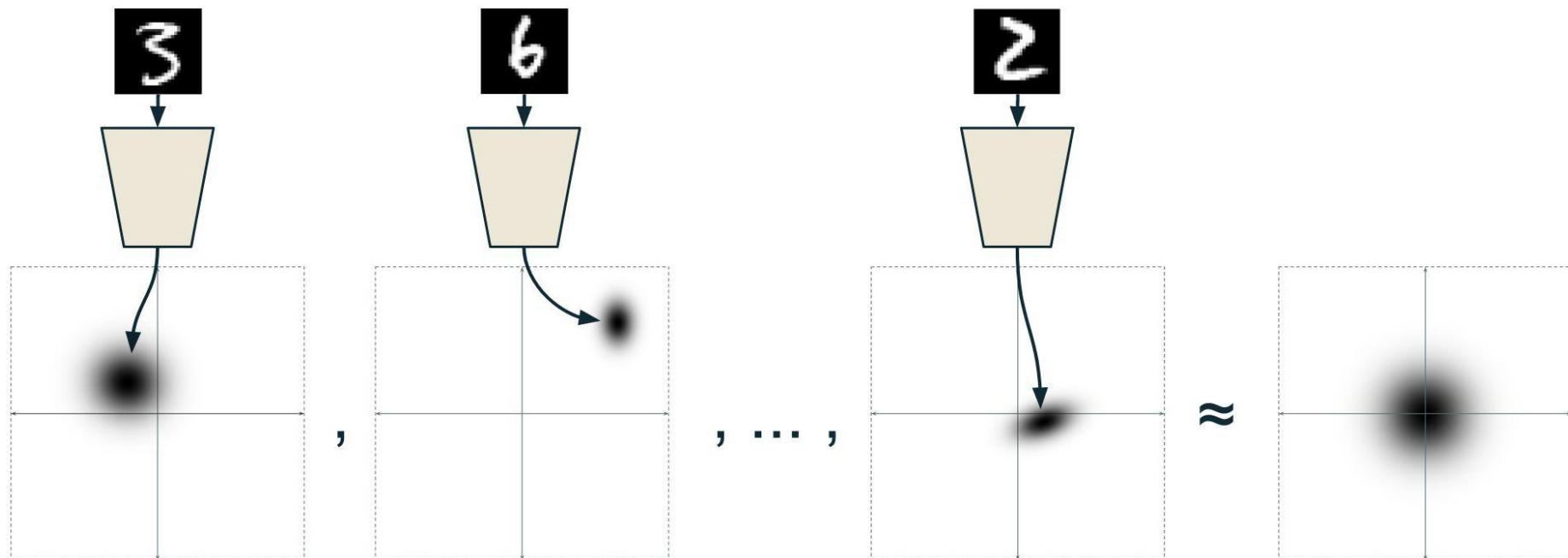


Transferring style with

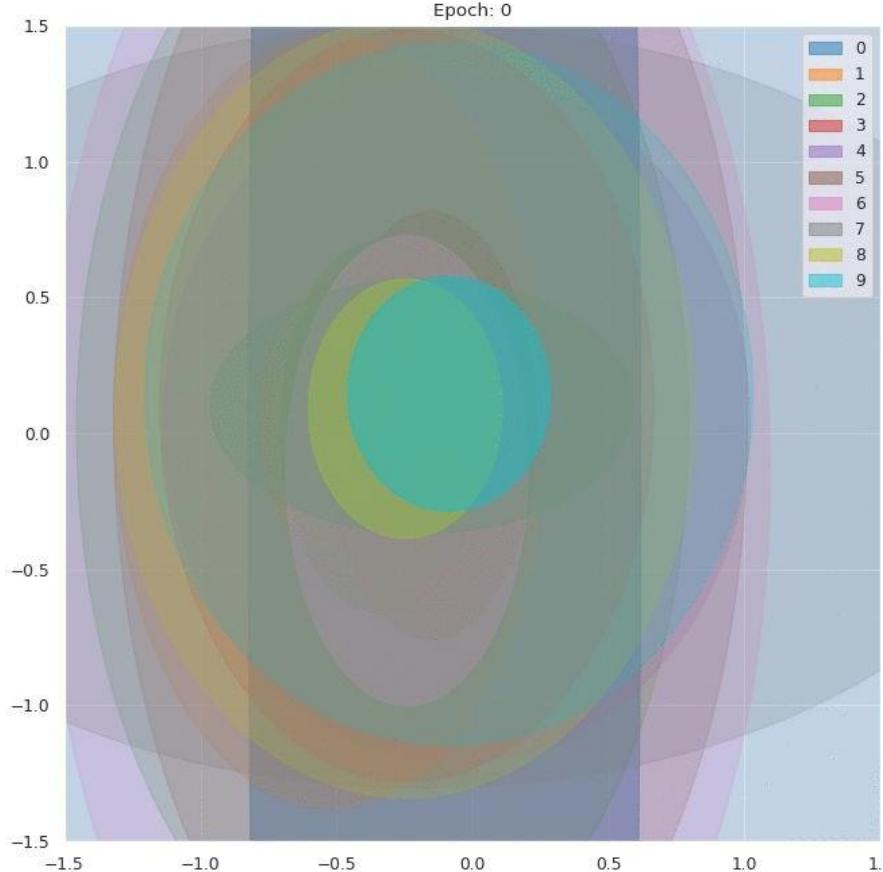
0	0	0	0	0	0	0	0	0	0
1	1	1	1	1	1	1	1	1	1
2	2	2	2	2	2	2	2	2	2
3	3	3	3	3	3	3	3	3	3
4	4	4	4	4	4	4	4	4	4
5	5	5	5	5	5	5	5	5	5
6	6	6	6	6	6	6	6	6	6
7	7	7	7	7	7	7	7	7	7
8	8	8	8	8	8	8	8	8	8
9	9	9	9	9	9	9	9	9	9

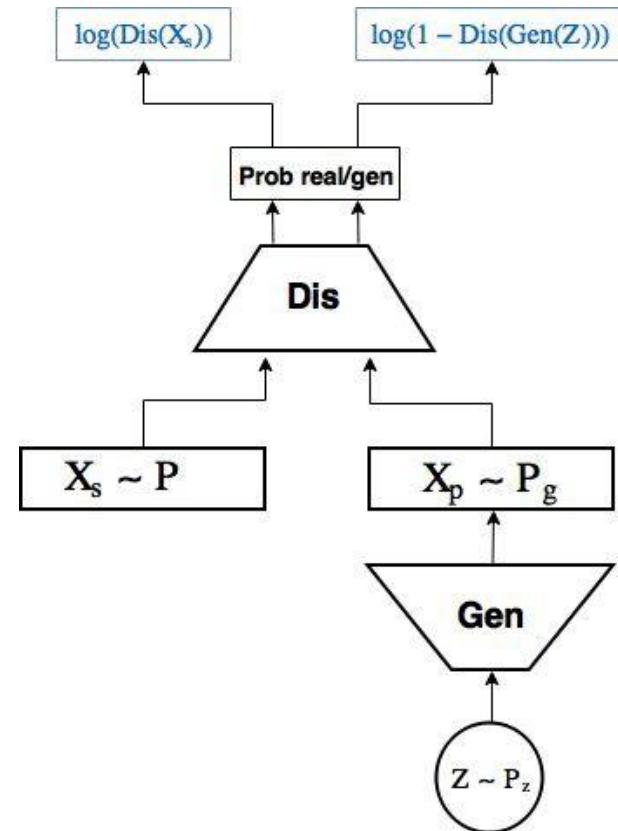
Seed label is 1

Once again

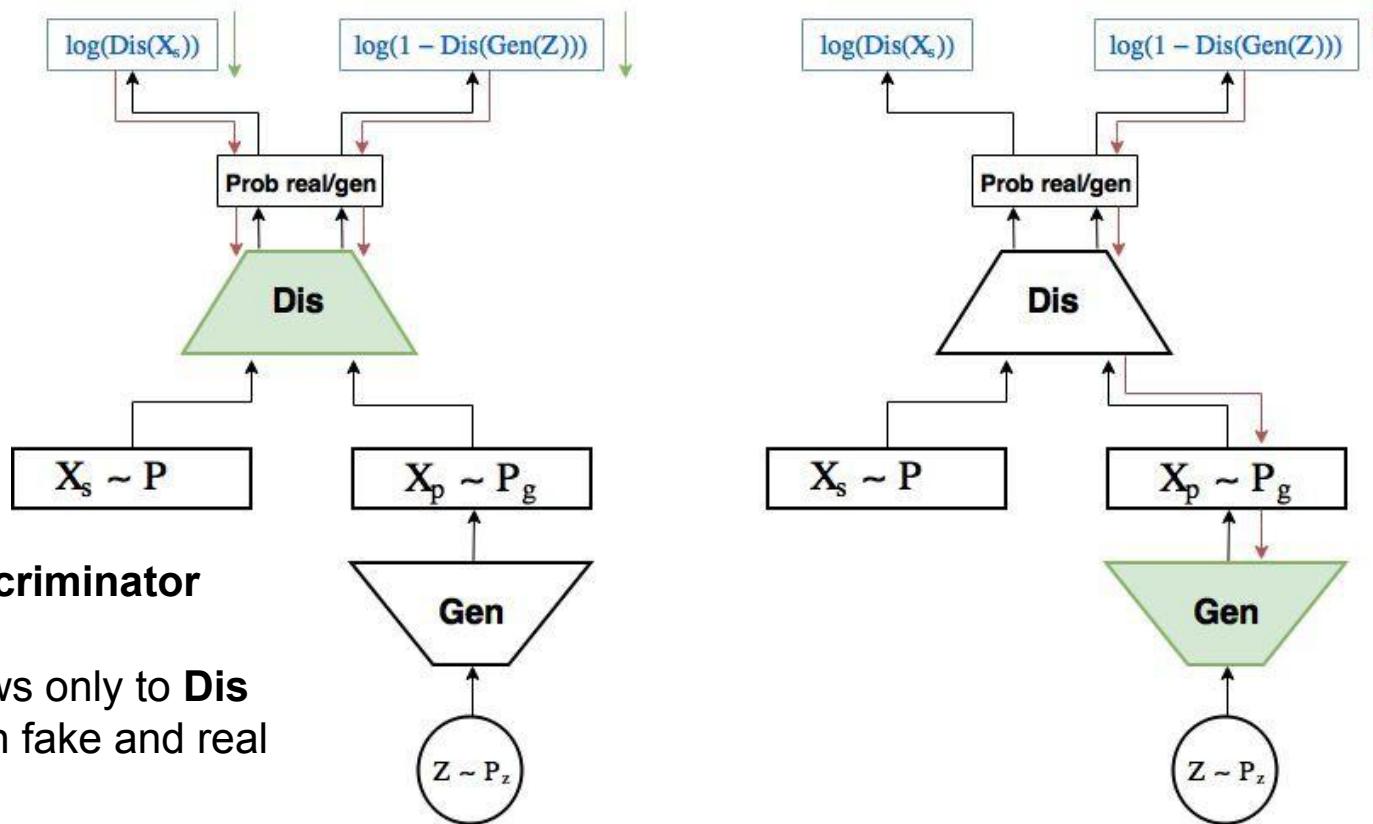


cVAE latent space distribution





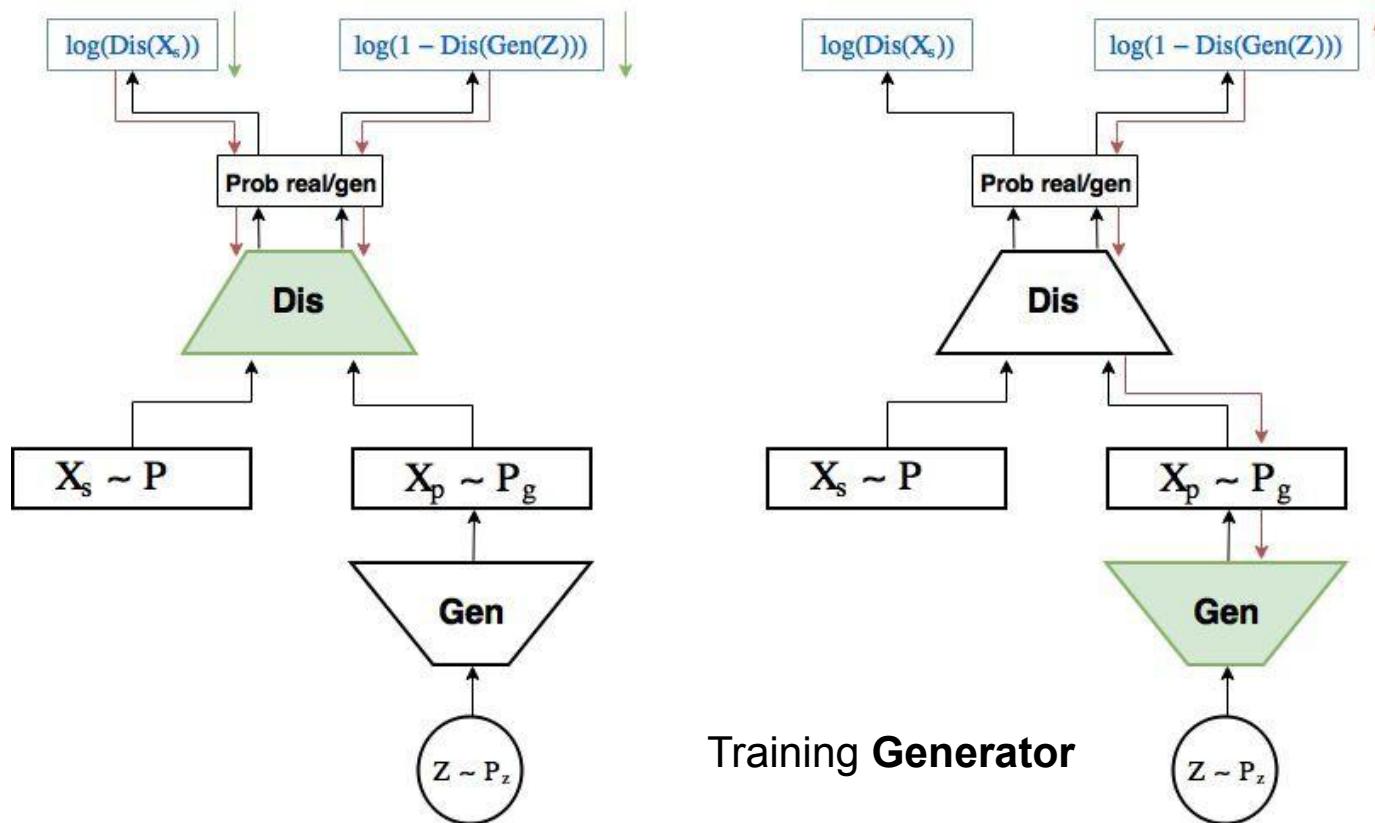
Training GAN



Training Discriminator

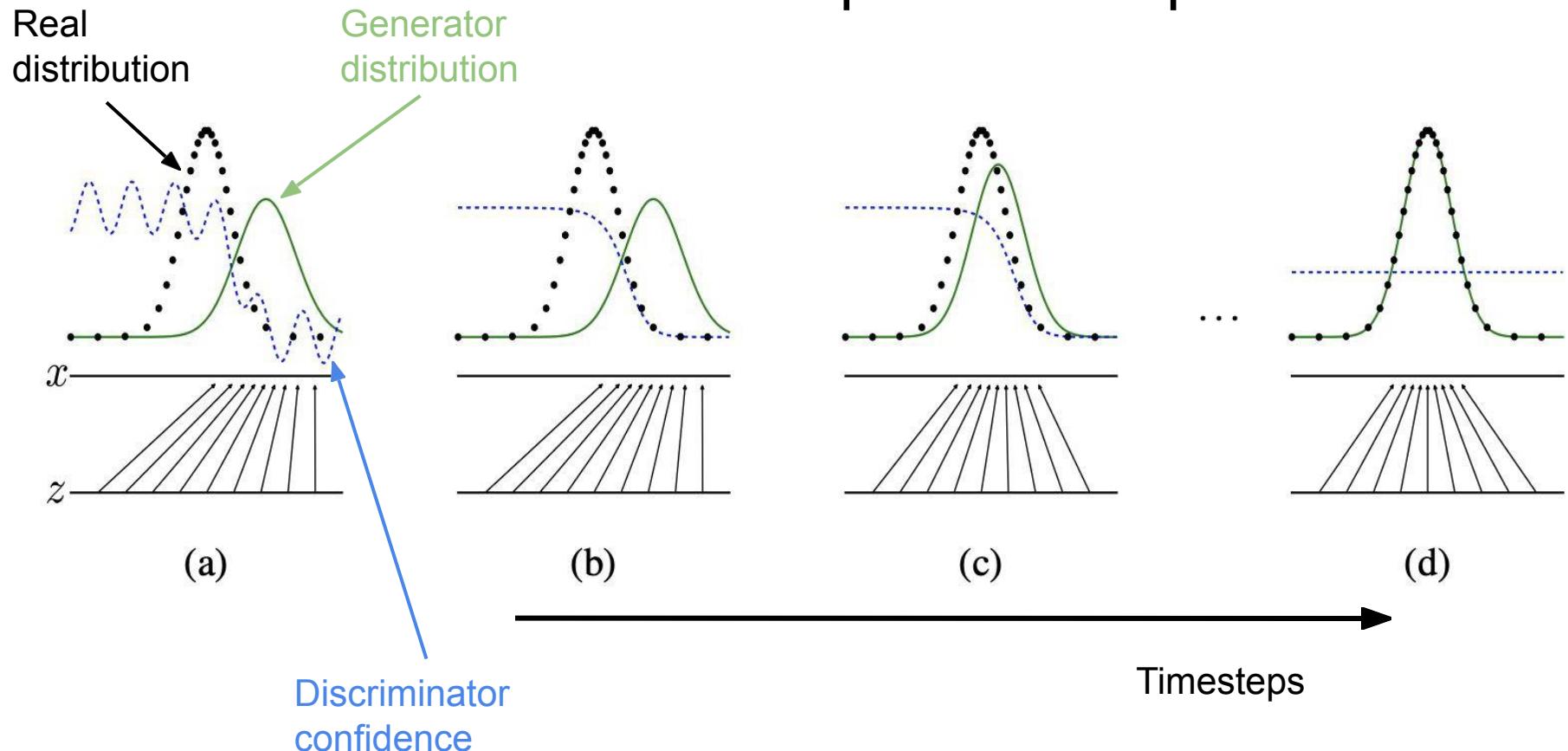
Gradient flows only to **Dis**
to distinguish fake and real
examples

Training GAN



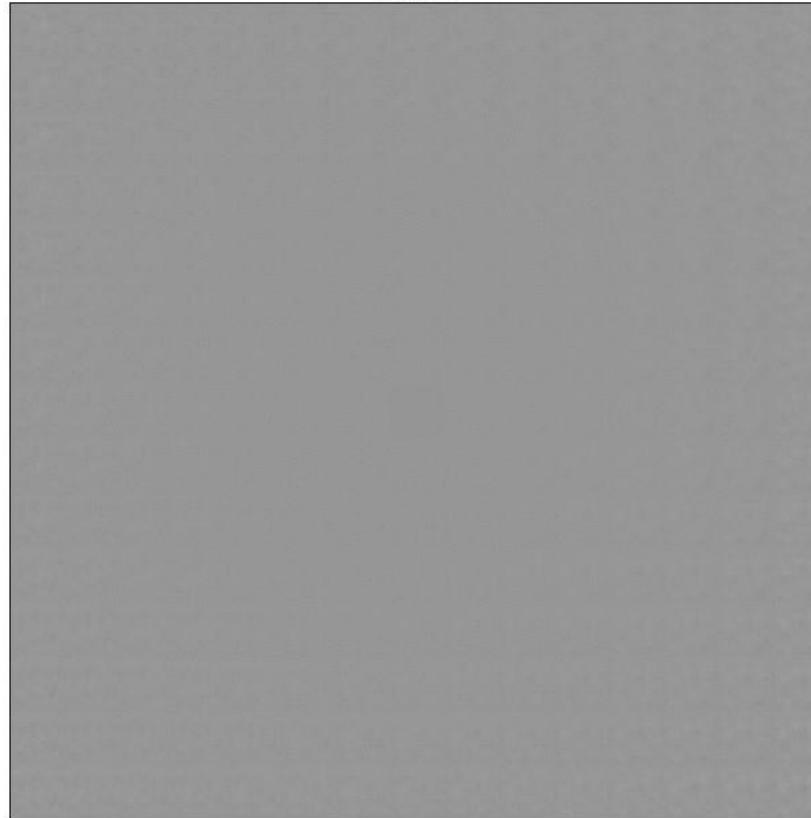
Gradient flows to **Gen** with **Dis** weights freezed
to fool the Discriminator

Optimization process in GAN

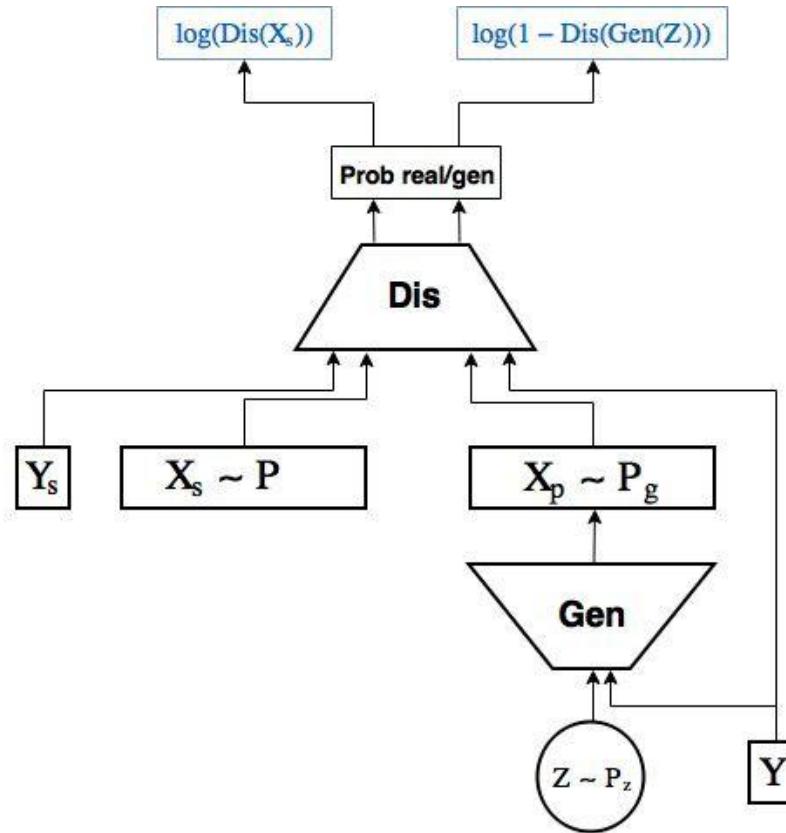


GAN manifold

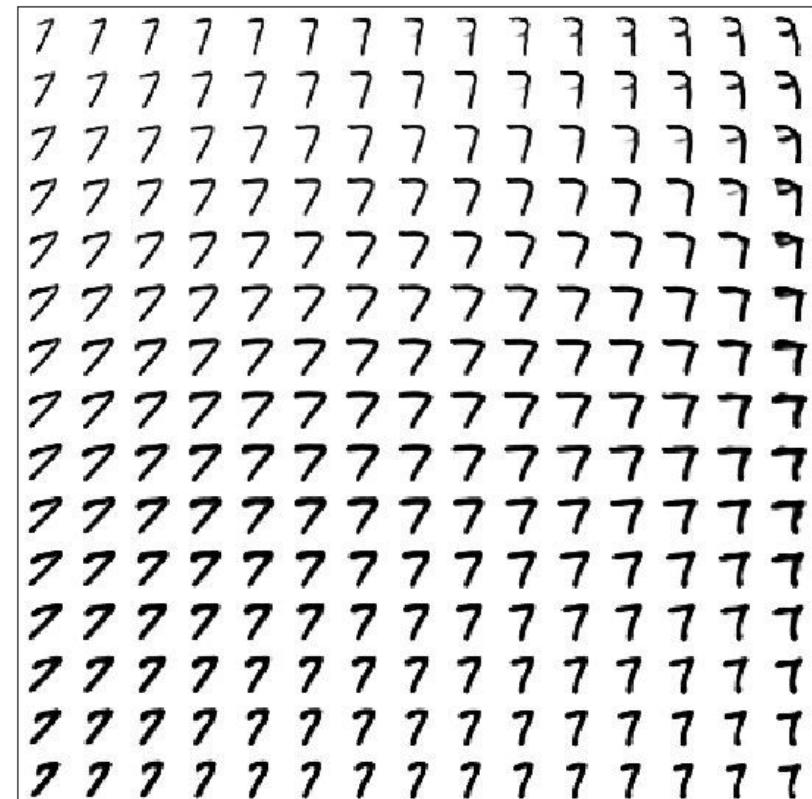
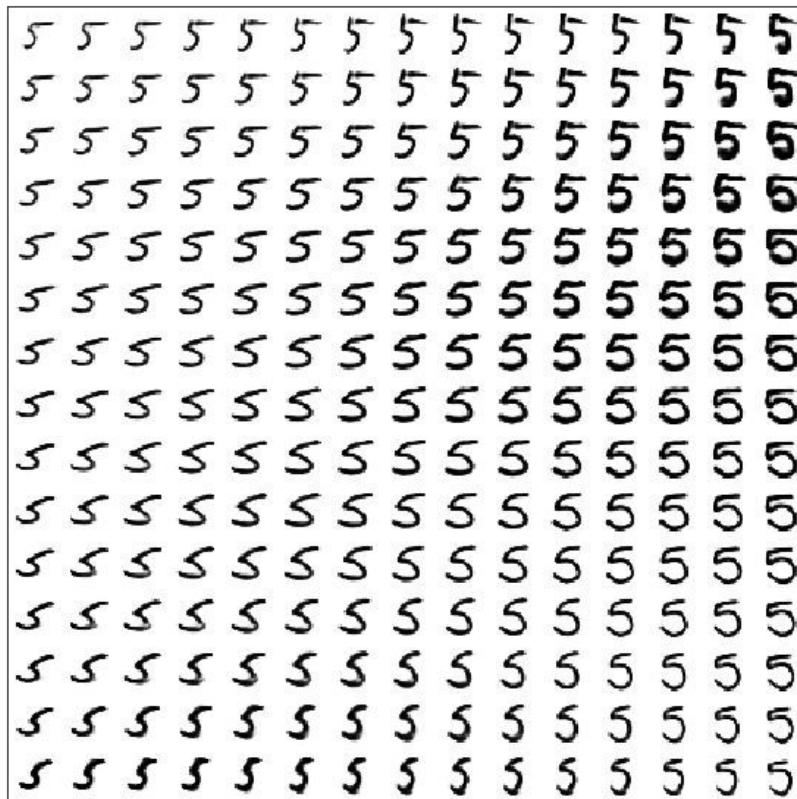
Label: all
Batch: 0



Conditional GAN



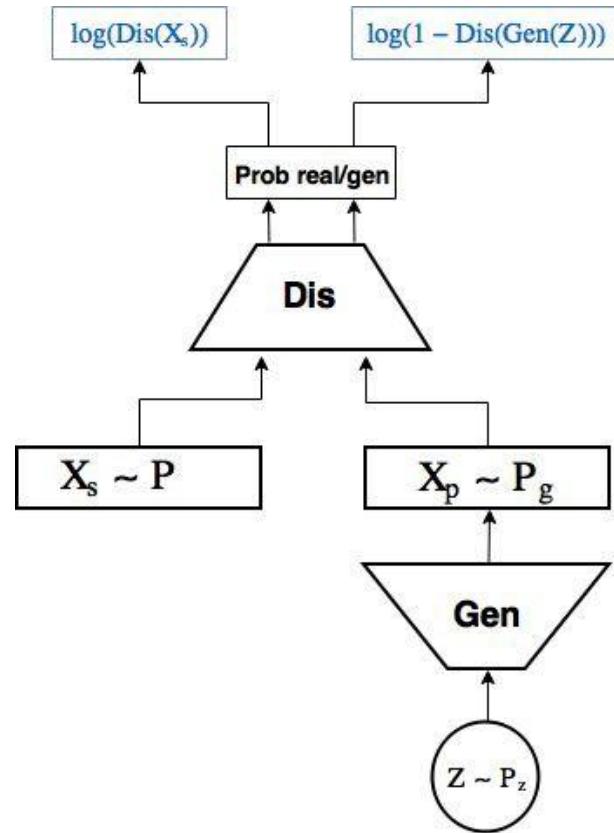
cGAN manifolds



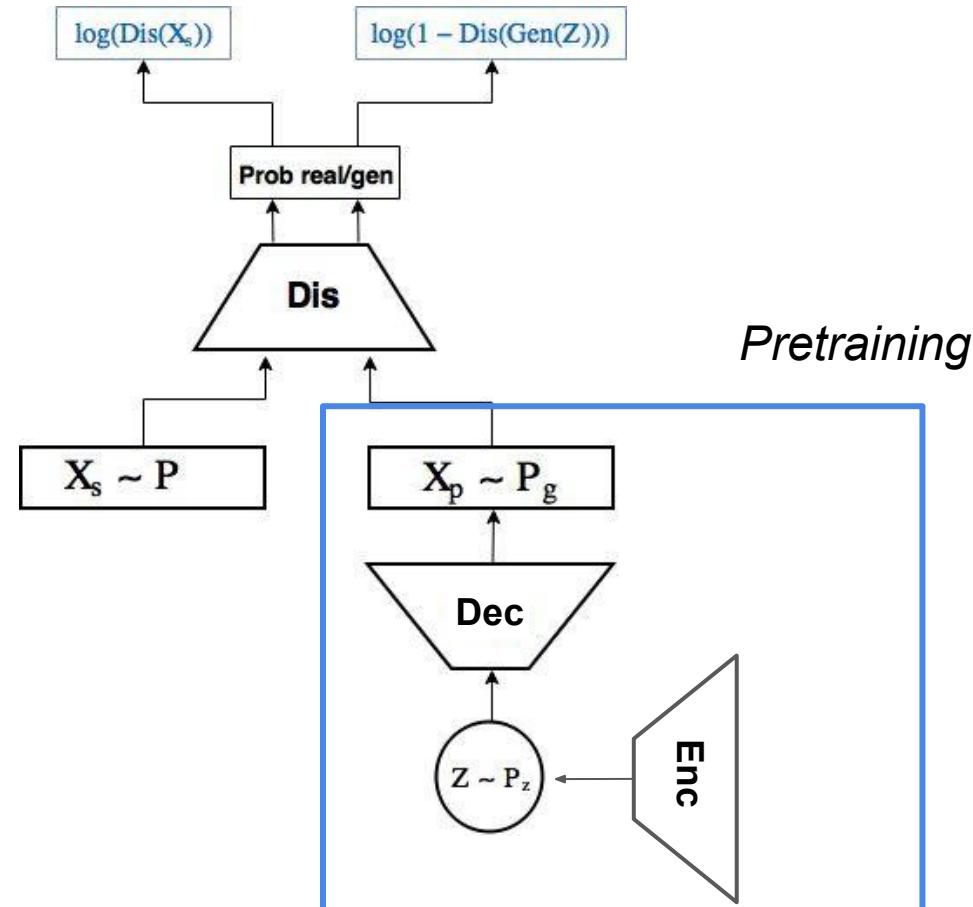
Some more combinations



Simple GAN



VAE/GAN



VAE/GAN original illustration

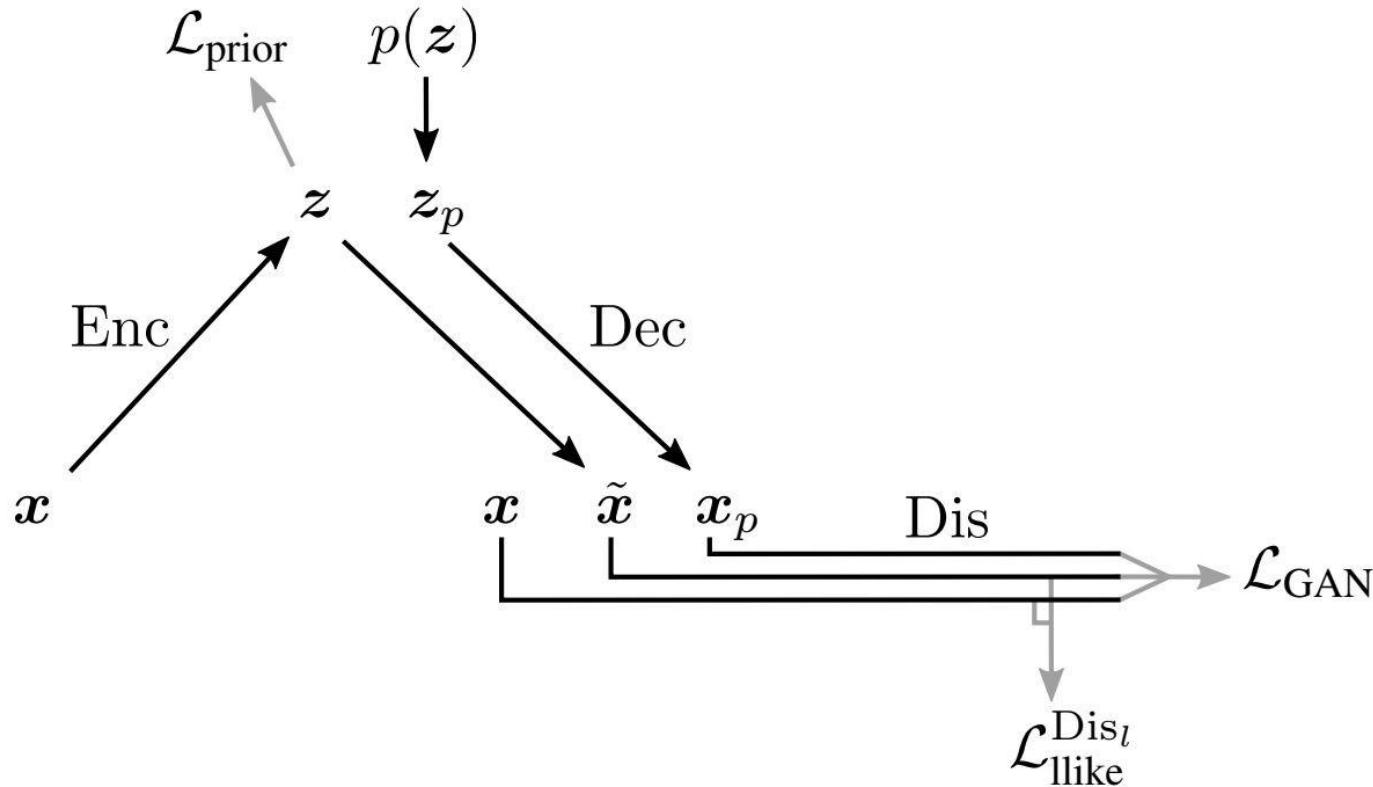


Image source: [Autoencoding beyond pixels using a learned similarity metric.](#)

Anders Boesen Lindbo Larsen et al, 2016

Slides source

https://web.eecs.umich.edu/~justincj/slides/eecs498/498_FA2019_lecture19.pdf

<https://web.eecs.umich.edu/~justincj/teaching/eecs498/WI2022/>