



EECS 230 Deep Learning

Lecture 16: Variational Autoencoder

Some slides from Frank Noe and Pascal Poupart

Outline

- ❑ Generative model
- ❑ Variational autoencoder
 - ❑ Autoencoder
 - ❑ Variational autoencoder



Generative model

So far...

- ❑ Discriminative model $P(y|x)$

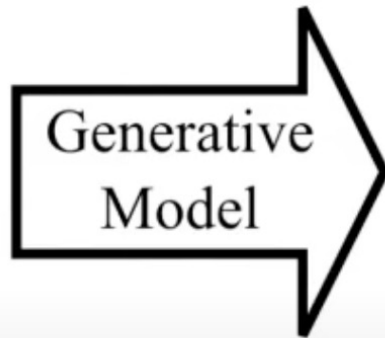
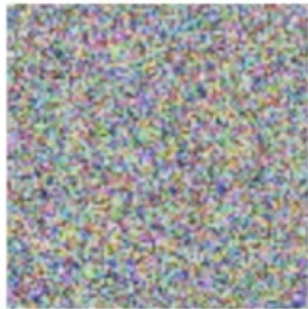
- ❑ Given input data x , predict y
 - ❑ E.g., classification, regression

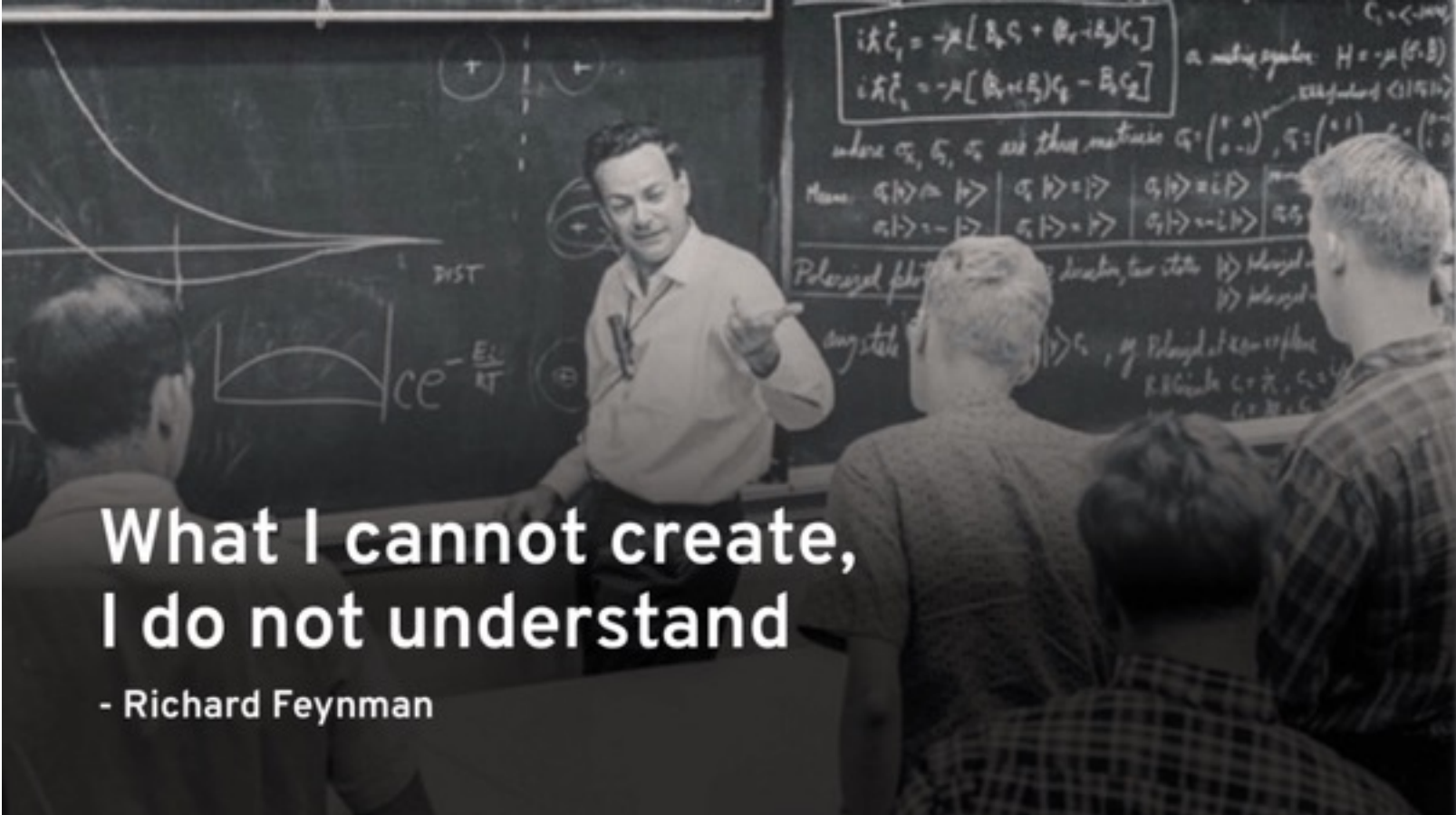
- ❑ Generative model

- ❑ Model data distribution $P(x)$
 - ❑ Sample from $P(x)$ to generate new data

Generative model

Noise $\sim N(0,1)$



A black and white photograph of Richard Feynman in a lecture hall. He is standing in front of a chalkboard, gesturing with his right hand. The chalkboard is filled with mathematical equations and diagrams. To the left, there is a graph with a curve and the word "DIST" written below it. To the right, there are several equations, including $i\hbar \dot{c}_1 = -\mu [B_0 c_1 + (B_0 - iB_1) c_2]$ and $i\hbar \dot{c}_2 = -\mu [B_0 c_2 + (B_0 + iB_1) c_1]$. Feynman is wearing a light-colored shirt and dark trousers. In the foreground, the backs of several students' heads are visible, showing they are listening to the lecture.

What I cannot create, I do not understand

- Richard Feynman

Types of generative neural networks

- ☐ Boltzmann machines
- ☐ Sigmoid belief networks
- ☐ **Variational autoencoders (inference net + generator net)**
- ☐ **Generative adversarial networks (generator net + discriminator net)**
- ☐ Normalizing flows
- ☐ **Diffusion Models**
- ☐ ...



Variational Autoencoder

Autoencoder

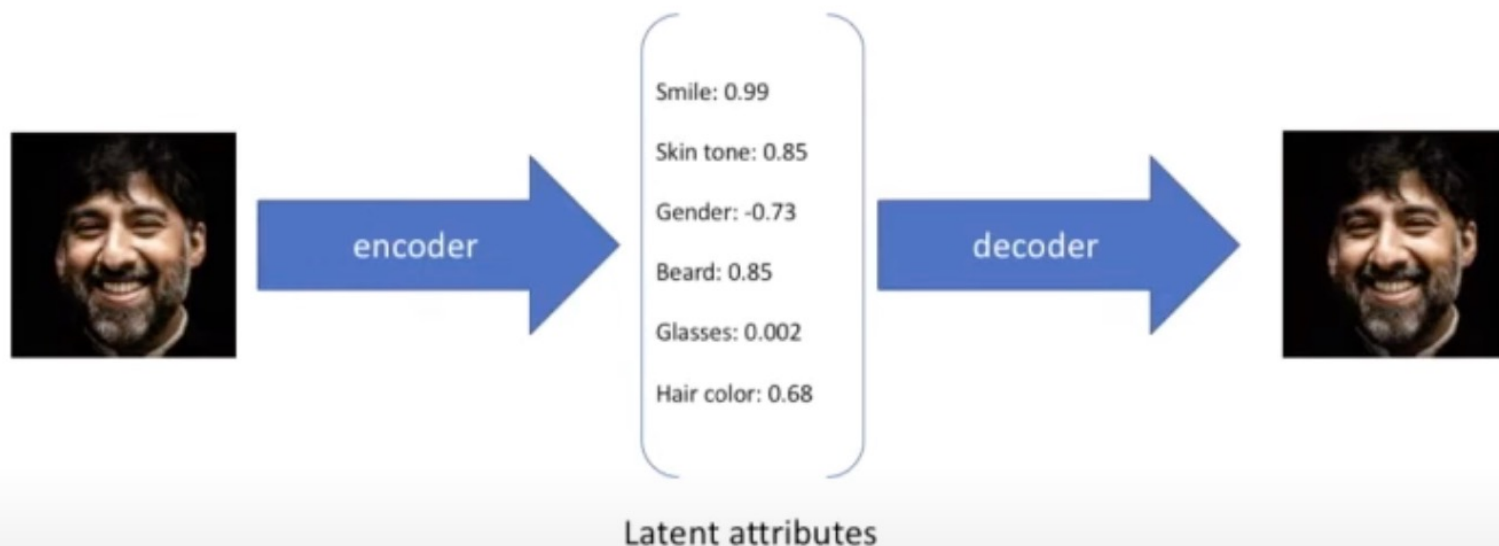
- ❑ Special type of feed forward network for
 - ❑ Compression
 - ❑ Denoising
 - ❑ Sparse representation
 - ❑ Data generation

Autoencoder

❑ Encoder: $f(\cdot)$

❑ Decoder: $g(\cdot)$

❑ Autoencoder: $g(f(x)) = x$



Latent variable encodes “essential” information about input data

Linear autoencoder

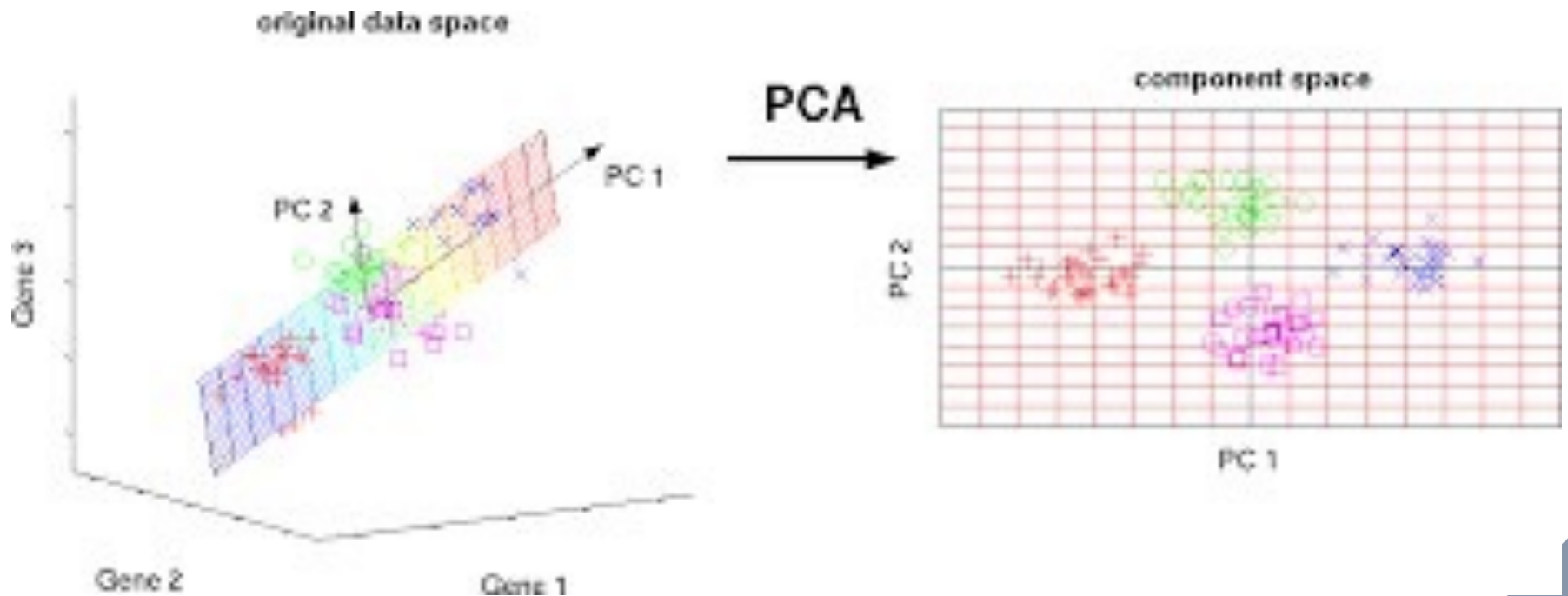
- ❑ Objective: find weights W_f and W_g that minimize reconstruction error

$$\min_{\mathbf{W}} \frac{1}{2} \sum_n \left\| \mathbf{W}_g \mathbf{W}_f \mathbf{x}_n - \mathbf{x}_n \right\|_2^2$$

- ❑ When using Euclidean norm (i.e., squared loss), solution is the same as principal component analysis (PCA)

Recap: principle component analysis

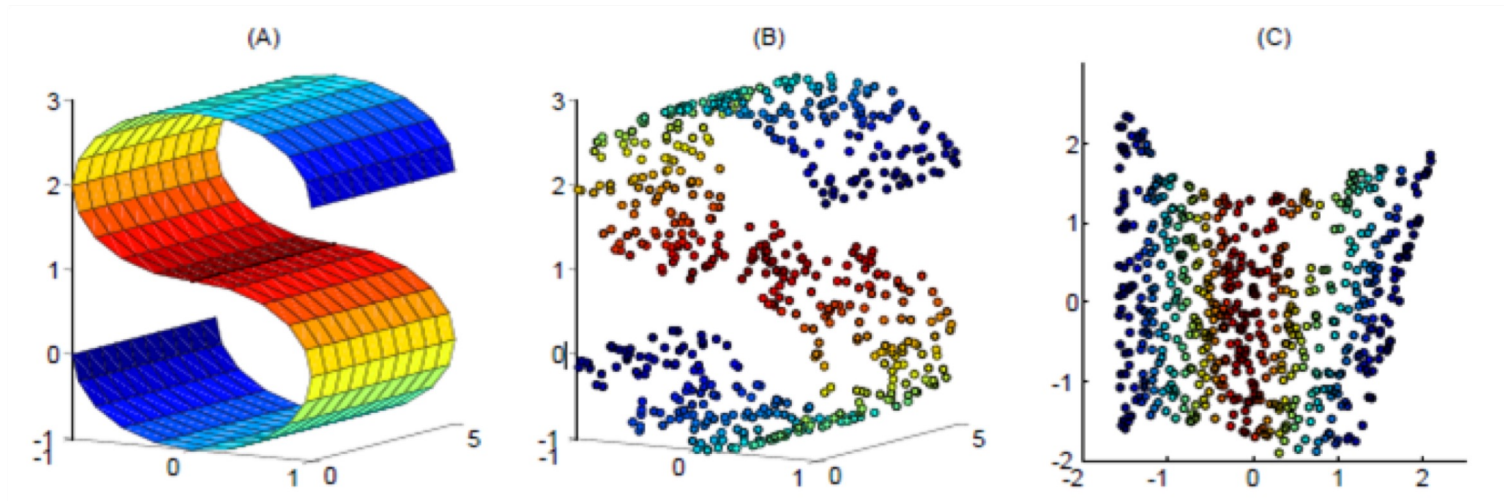
□ Components with maximum variance



Non-linear autoencoder

□ $f(\cdot)$ and $g(\cdot)$ are both non-linear functions

$$\min_W \frac{1}{2} \sum_n \left\| g(f(\mathbf{x}_n; \mathbf{W}_f); \mathbf{W}_g) - \mathbf{x}_n \right\|_2^2$$



Sparse representation

□ When more hidden nodes than inputs, use regularization to constrain autoencoder

□ Example: force hidden nodes to be sparse

$$\min_{\mathbf{W}} \frac{1}{2} \sum_n \left\| g(f(\mathbf{x}_n; \mathbf{W}_f); \mathbf{W}_g) - \mathbf{x}_n \right\|_2^2 + c \underbrace{nnz(f(\mathbf{x}_n; \mathbf{W}_f))}_{\text{Sparse hidden nodes}}$$

where $nnz(f(\mathbf{x}_n; \mathbf{W}_f))$ is the number of non-zero entries in the vector produced by f .

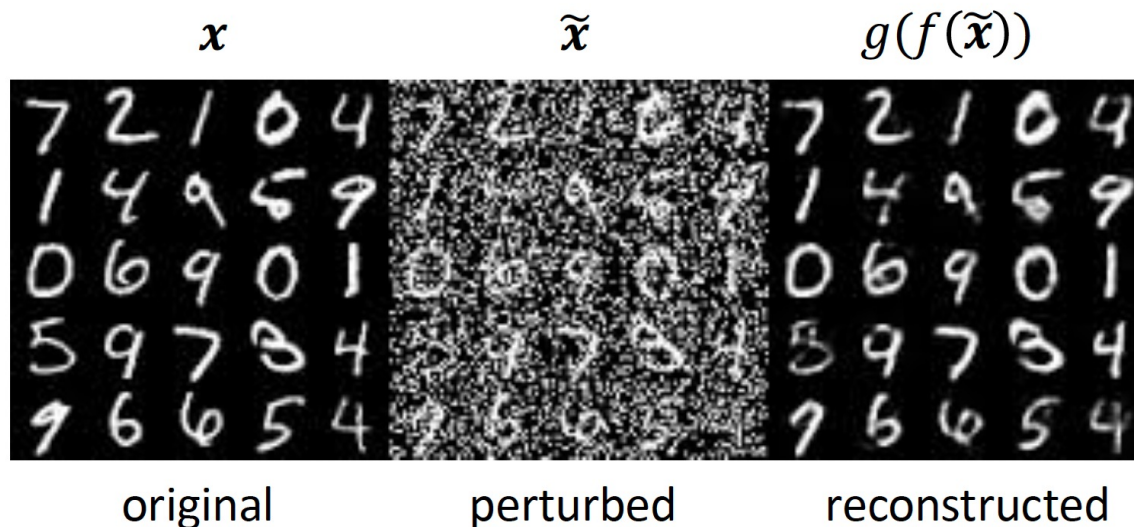
□ Approximate objective: L_1 regularization

$$\min_{\mathbf{W}} \frac{1}{2} \sum_n \left\| g(f(\mathbf{x}_n; \mathbf{W}_f); \mathbf{W}_g) - \mathbf{x}_n \right\|_2^2 + c \left\| f(\mathbf{x}_n; \mathbf{W}_f) \right\|_1$$

Denoising autoencoder

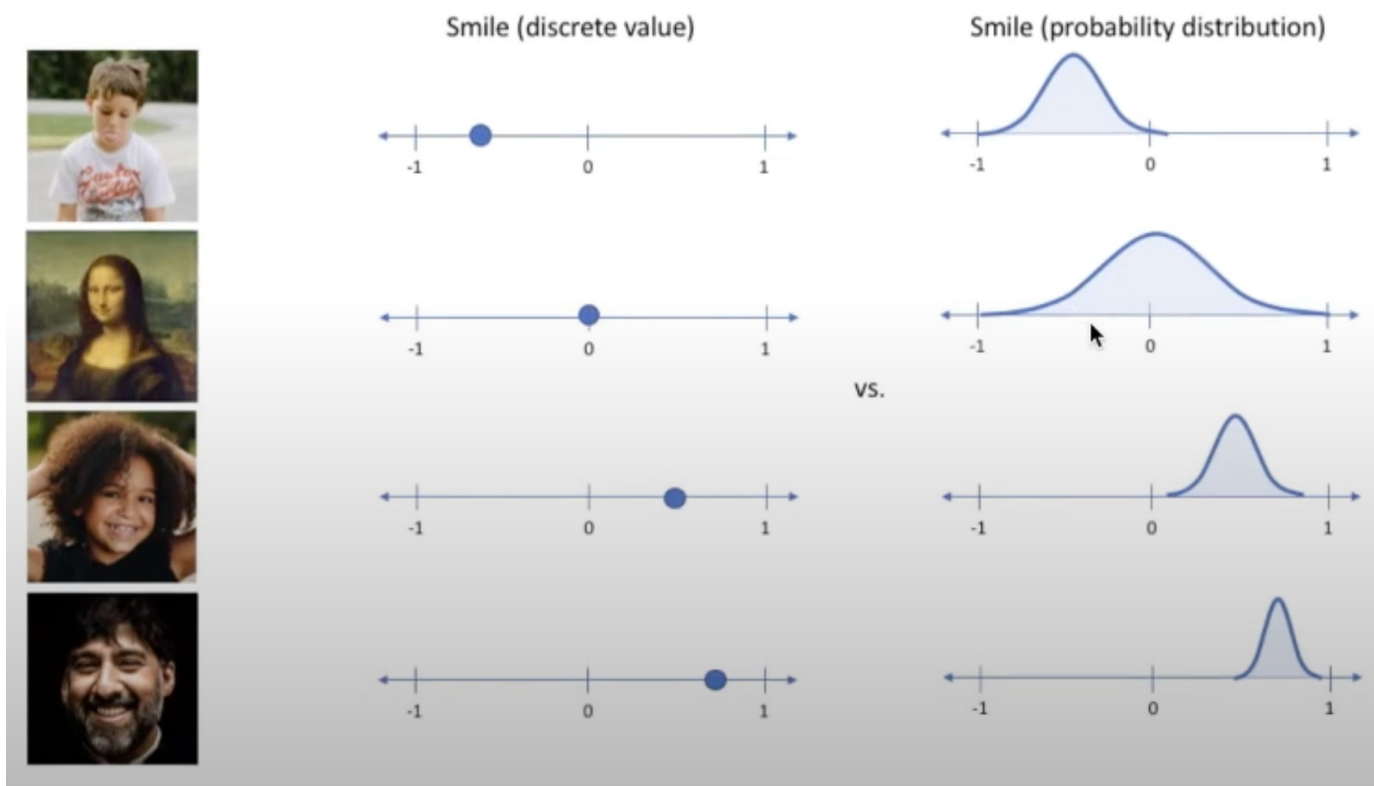
□ Consider noisy version \tilde{x} of the input x

$$\min_{\mathbf{W}} \frac{1}{2} \sum_n \left\| g(f(\tilde{x}_n; \mathbf{W}_f); \mathbf{W}_g) - \mathbf{x}_n \right\|_2^2 + c \left\| f(\tilde{x}_n; \mathbf{W}_f) \right\|_1$$

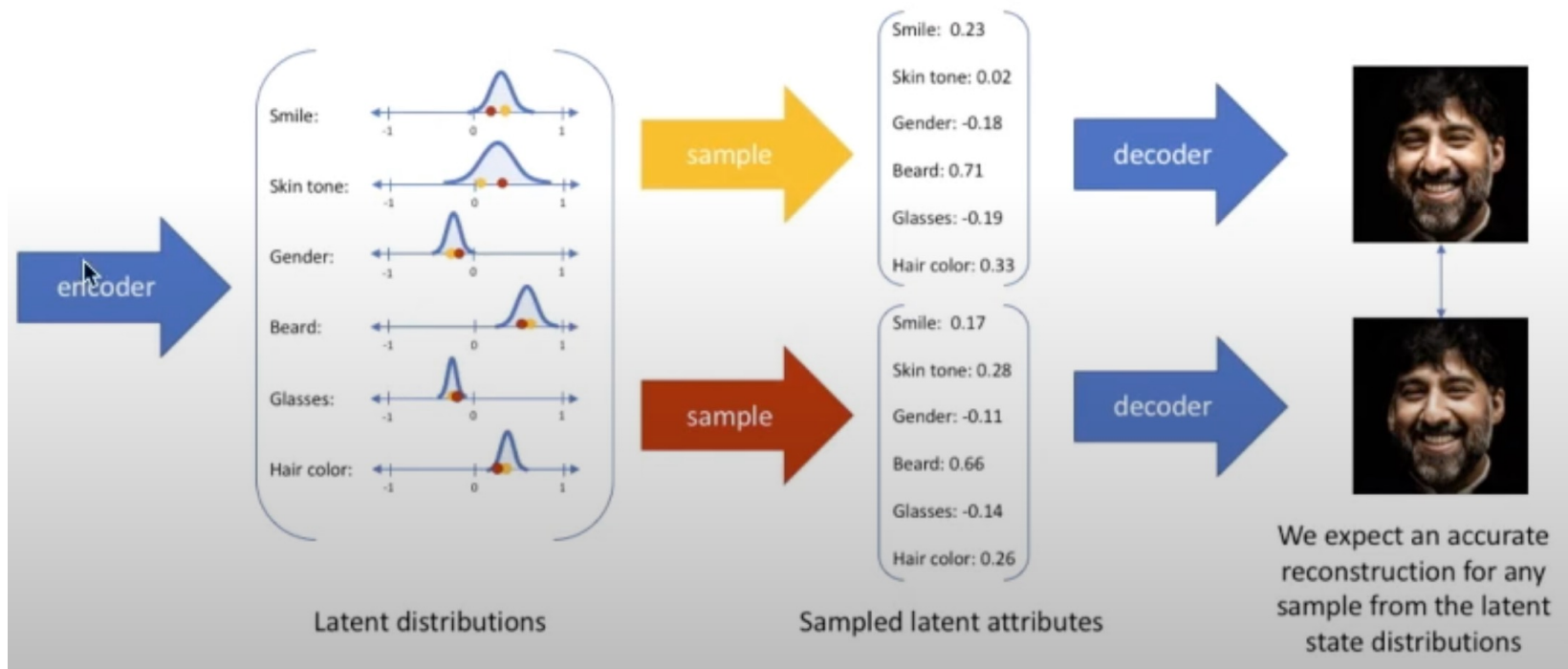


Variational/Probabilistic autoencoder

- ❑ Instead of a single value for each attribute, represent each attribute as a range of values
- ❑ VAE: Describe latent attribute in probabilistic terms



Variational/Probabilistic autoencoder



Variational/Probabilistic autoencoder

- ❑ Let $f()$ and $g()$ represent conditional distributions
 - ❑ $f: \Pr(h|x; w_f)$
 - ❑ $g: \Pr(x|h; w_g)$
- ❑ The decoder $g()$ can be treated as a generative model
 - ❑ First sample h from $\Pr(h)$
 - ❑ Then sample x from $\Pr(x|h; w_g)$

Variational autoencoder

- Idea: train encoder $\Pr(h|x; w_f)$ to approach a simple and fixed distribution, e.g., $N(h; 0, I)$

$$\max_{\mathbf{W}} \sum_n \log \Pr(\mathbf{x}_n; \mathbf{W}_f, \mathbf{W}_g) - \underbrace{c \, KL(\Pr(\mathbf{h}|\mathbf{x}_n; \mathbf{W}_f) || N(\mathbf{h}; \mathbf{0}, I))}_{\text{Kullback-Leibler divergence}}$$

Kullback-Leibler divergence
Distance measure for distributions

Variational Autoencoder Likelihood

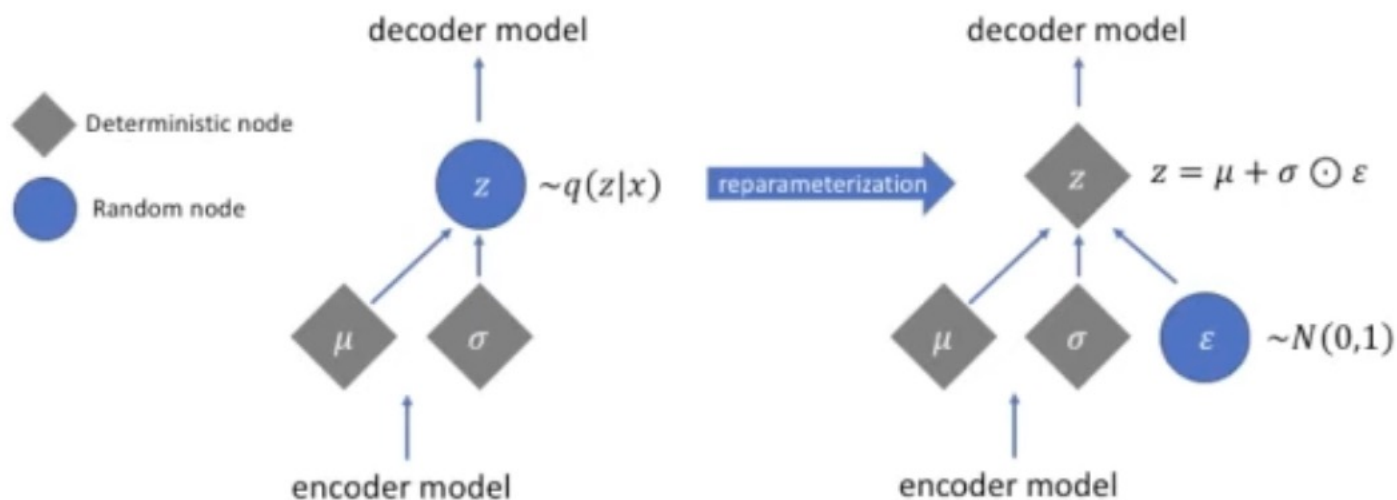
□ How to compute $\Pr(\mathbf{x}_n; \mathbf{W}_f, \mathbf{W}_g)$?

$$\Pr(\mathbf{x}_n; \mathbf{W}_f, \mathbf{W}_g) = \int_{\mathbf{h}} \Pr(\mathbf{x}_n | \mathbf{h}; \mathbf{W}_g) \Pr(\mathbf{h} | \mathbf{x}_n; \mathbf{W}_f) d\mathbf{h}$$

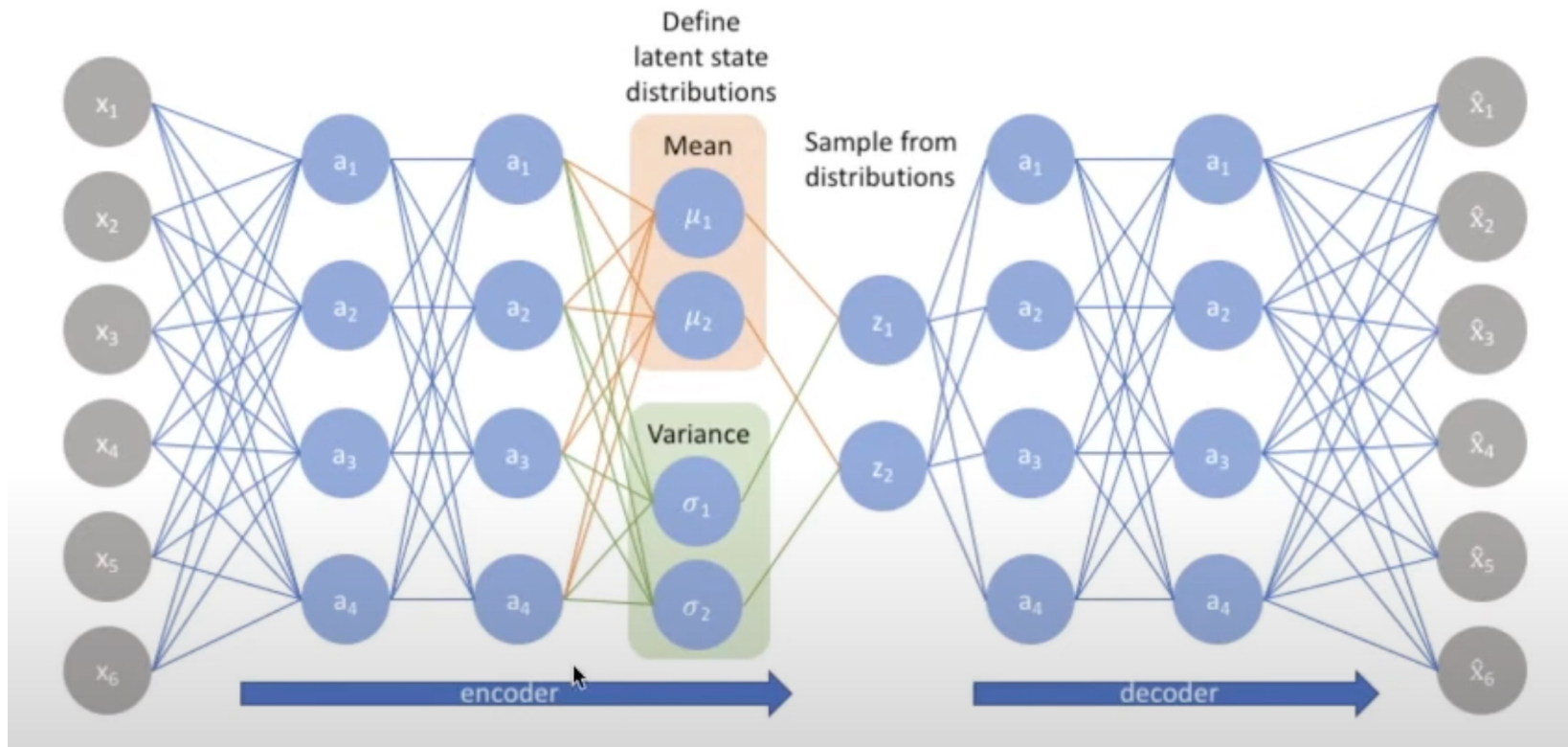
□ Obtain mean and variance of $\Pr(\mathbf{h})$ by neural network

$$\Pr(\mathbf{h} | \mathbf{x}_n; \mathbf{W}_f) = N(\mathbf{h}; \mu_n(\mathbf{x}_n; \mathbf{W}_f), \sigma_n(\mathbf{x}_n; \mathbf{W}_f) \mathbf{I})$$

Reparameterization trick



VAE implementation

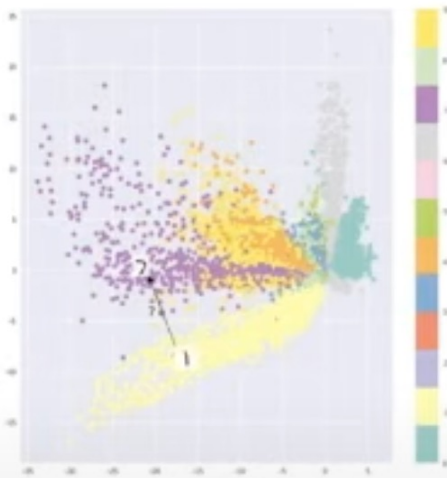


MNIST VAE

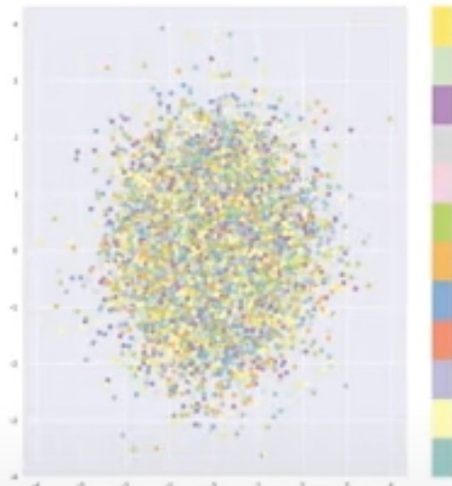
$$\max_{\mathbf{W}} \sum_n \log \Pr(\mathbf{x}_n; \mathbf{W}_f, \mathbf{W}_g) - \underbrace{c \, KL(\Pr(\mathbf{h}|\mathbf{x}_n; \mathbf{W}_f) || N(\mathbf{h}; \mathbf{0}, \mathbf{I}))}_{\text{Kullback-Leibler divergence}}$$

Kullback-Leibler divergence

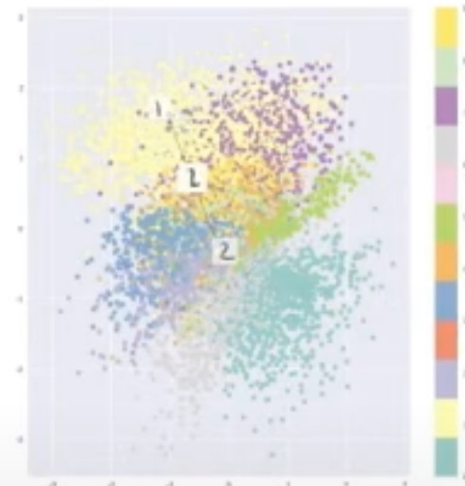
Only reconstruction loss



Only KL divergence



Combination



Examples from VAE

