



# Variational Autoencoders

Gunhee Kim

Computer Science and Engineering



서울대학교

SEOUL NATIONAL UNIVERSITY

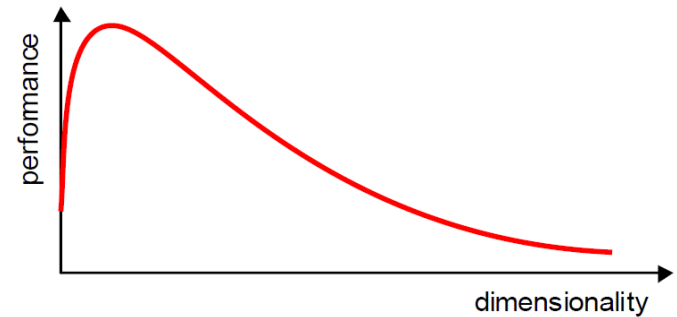
# Outline

- Autoencoders
- Variational Autoencoders

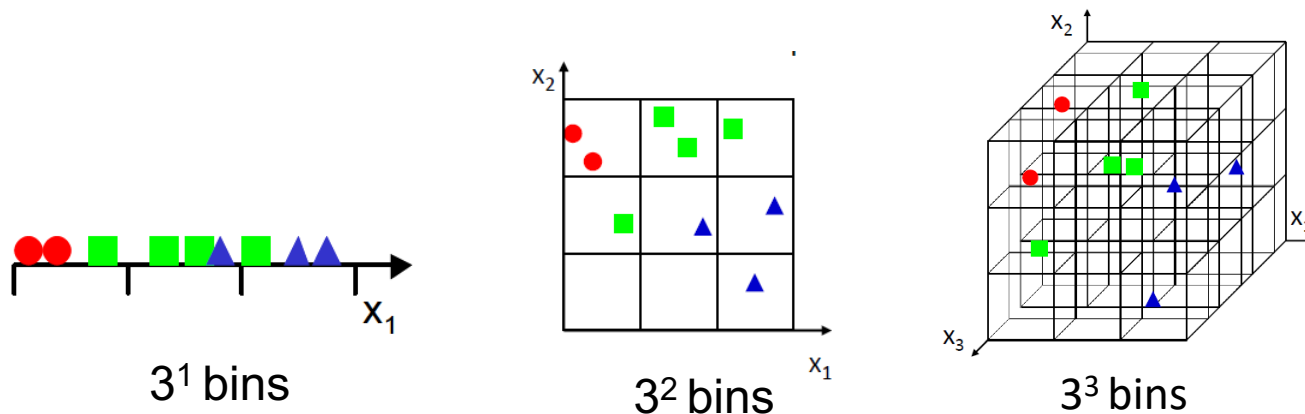
# Curse of Dimensionality

Increasing the number of features will not always improve classification accuracy

- In practice, the inclusion of more features might lead to worse performance



The number of training examples required increases exponentially with dimensionality  $d$  (i.e.,  $k^d$ )

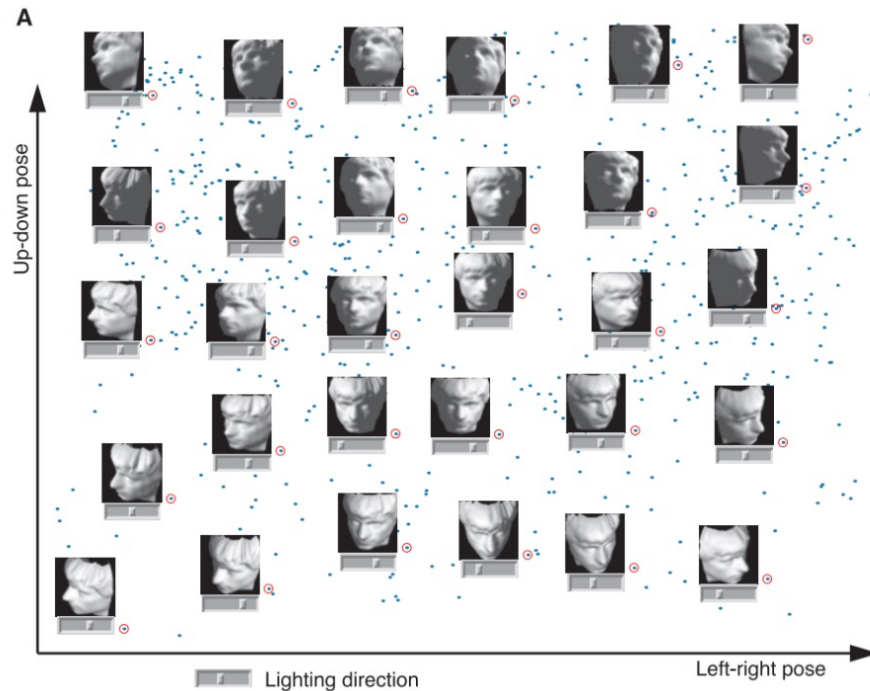


# Why Dimensionality Reduction?

Feature selection: some features may be irrelevant

Visualization: especially for high dimensional data

Intrinsic dimensionality: smaller than # of features



# Feature Selection vs Extraction

## Feature selection (supervised)

- Chooses a subset of the original features

$$\mathbf{x} = \begin{bmatrix} x_1 \\ x_2 \\ \cdot \\ \cdot \\ \cdot \\ x_p \end{bmatrix} \longrightarrow \mathbf{y} = \begin{bmatrix} x_{i1} \\ x_{i2} \\ \cdot \\ x_{ik} \end{bmatrix} \\ k \ll p$$

## Feature extraction (unsupervised)

- Finds a set of new features (i.e., through some mapping  $f(x)$ ) from the existing features
- The mapping  $f(x)$  could be linear or non-linear

$$\mathbf{x} = \begin{bmatrix} x_1 \\ x_2 \\ \cdot \\ \cdot \\ \cdot \\ x_p \end{bmatrix} \xrightarrow{f(x)} \mathbf{y} = \begin{bmatrix} y_1 \\ y_2 \\ \cdot \\ y_k \end{bmatrix} \\ k \ll p$$

# Feature Extraction

Often called dimensionality reduction or manifold learning

How to find an optimum mapping  $y = f(x)$  is equivalent to optimizing an objective function?

Minimize information loss

- The goal is to represent the data as accurately as possible (i.e., no loss of information) in the lower-dimensional space

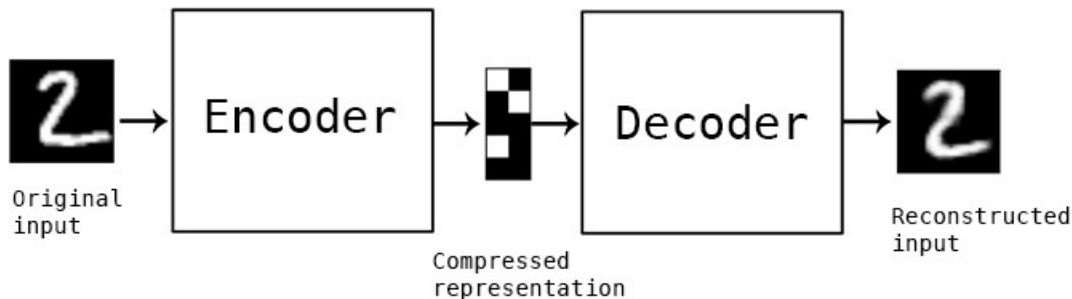
Maximize discriminatory information

- The goal is to enhance the class-discriminatory information in the lower-dimensional space

# Autoencoders

## An unsupervised neural network model

- Used for dimensionality reduction (e.g., feature selection and extraction)
- Lossy dimensionality reduction with few hidden units
- e.g.,  $10 \times 10$  images as input, and 50 hidden units  
→ compressed representation of images



- Encoder: represent (or compress) input data into a low-dim code
- Decoder: decompress a code into a data
- Encoder and decoder are implemented by neural networks

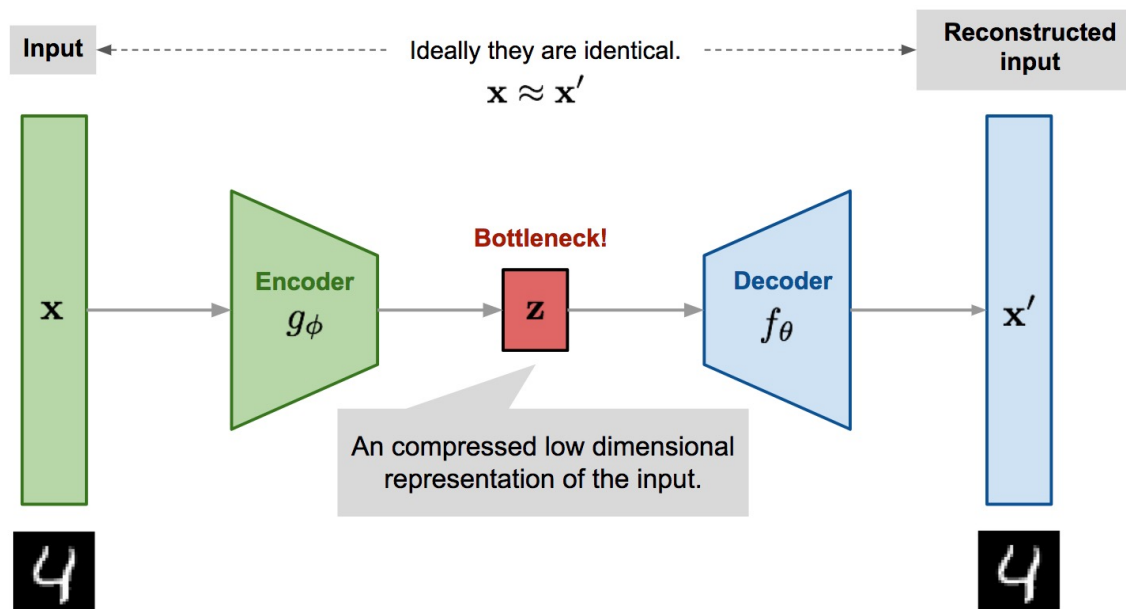
# Formulation

## Objective for dimensionality reduction

- Encoder  $g_\phi: \mathcal{X} \rightarrow \mathcal{Z}$  and decoder  $f_\theta: \mathcal{Z} \rightarrow \mathcal{X}$

$$\phi, \theta = \min_{\phi, \theta} \|X - (f_\theta \circ g_\phi)X\|^2$$

- Feature space  $\mathcal{Z}$  has often lower dimensionality than input space  $\mathcal{X}$





# Formulation

## A simple NN model

- An input  $\mathbf{x} \in \mathbb{R}^d = \mathcal{X}$  maps to a code (i.e., latent variable)  $\mathbf{z} \in \mathbb{R}^p = \mathcal{Z}$ , which is reconstructed to  $\mathbf{x}'$

$$\mathbf{z} = \sigma(\mathbf{W}\mathbf{x} + \mathbf{b})$$

$$\mathbf{x}' = \sigma(\mathbf{W}'\mathbf{z} + \mathbf{b}')$$

- Encoder's weight matrix  $\mathbf{W}$  and bias  $\mathbf{b}$  could be the same with those of decoder  $\mathbf{W}'$  and  $\mathbf{b}'$

## Training

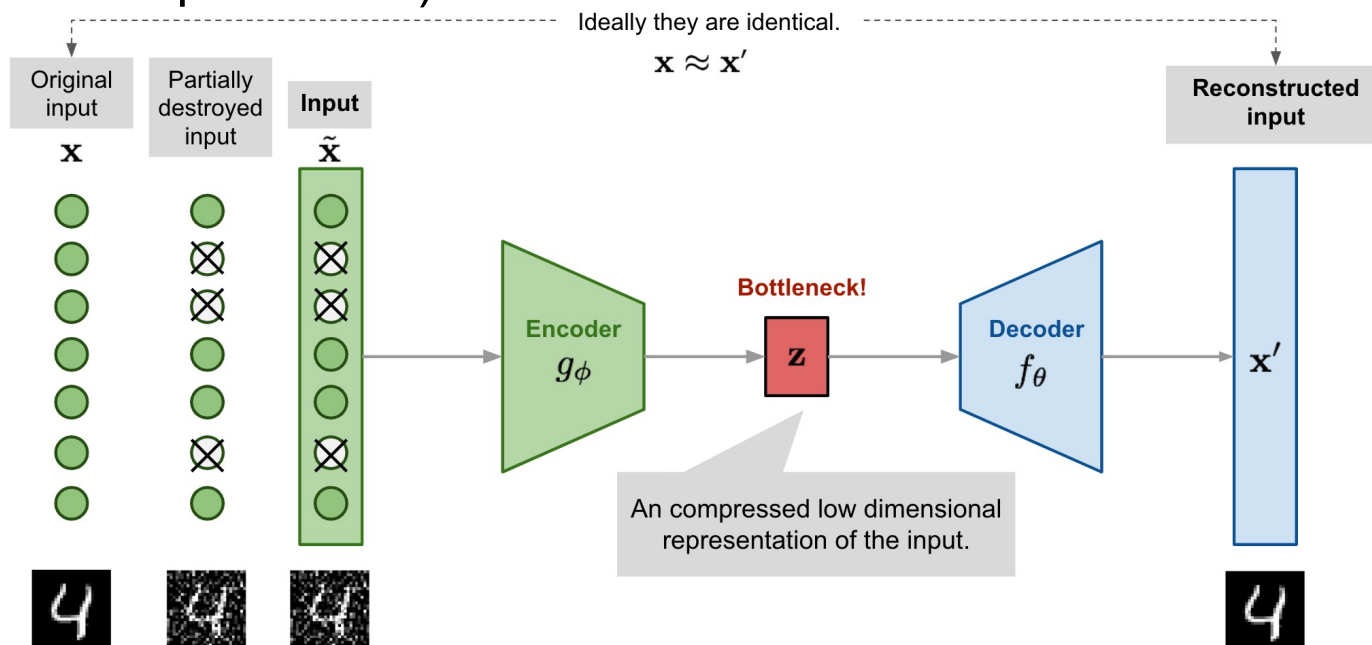
- AEs are also trained to minimize reconstruction errors (averaged over some input training set)

$$\mathcal{L}(\mathbf{x}, \mathbf{x}') = \|\mathbf{x} - \mathbf{x}'\|^2 = \|\mathbf{x} - \sigma(\mathbf{W}'\sigma(\mathbf{W}\mathbf{x} + \mathbf{b}) + \mathbf{b}')\|^2$$

# Variants – Denoising Autoencoder

A stochastic and more robust extension

- Can avoid the risk of overfitting
- Randomly corrupt input (by adding noises to or masking some values) and let AE reconstruct its denoising one
- With fewer data, add more noise (e.g., adding noise 30% corruption level)



# Variants – Sparse Autoencoder

A sparse constraint on the hidden unit activation

- Can avoid the risk of overfitting
- Only a small number of hidden units are activated simultaneously (i.e., one hidden neuron should be inactivated most of time)
- Average activation of hidden unit  $i$  in layer  $l$  over training data

$$\hat{\rho}_i^{(l)} = \frac{1}{m} \sum_{j=1}^m a_i^{(l)}(x^{(j)})$$

- Set sparsity parameter  $\hat{\rho}_i^{(l)} = \rho = 0.05$

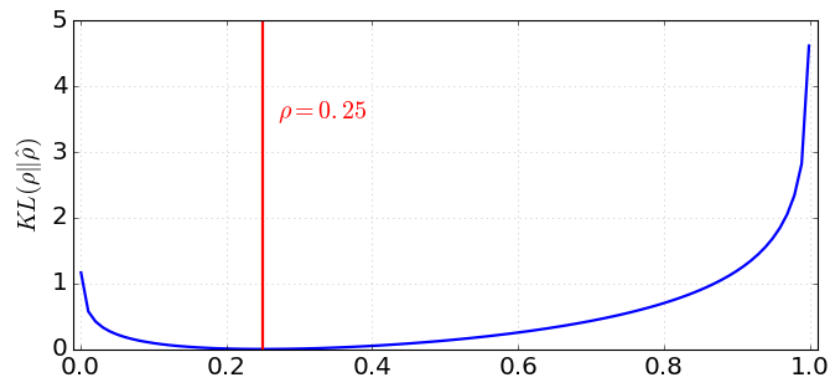
# Variants – Sparse Autoencoder

## The objective

- The constraint is achieved by adding a penalty term
- The KL-divergence between two Bernoulli distributions, one with mean  $\rho$  and the other with mean  $\hat{\rho}_i^{(l)}$

$$\begin{aligned}\mathcal{L} &= \frac{1}{m} \sum_{j=1}^m (\mathbf{x}_i - f_{\theta}(g_{\phi}(\tilde{\mathbf{x}}_i)))^2 + \beta \sum_{l=1}^L \sum_{j=1}^{s_l} KL(\rho \parallel \hat{\rho}_i^{(l)}) \\ &= \frac{1}{m} \sum_{j=1}^m (\mathbf{x}_i - f_{\theta}(g_{\phi}(\tilde{\mathbf{x}}_i)))^2 + \beta \sum_{l=1}^L \sum_{j=1}^{s_l} \rho \log \frac{\rho}{\hat{\rho}_i^{(l)}} + (1 - \rho) \log \frac{1 - \rho}{1 - \hat{\rho}_i^{(l)}}\end{aligned}$$

- The KL-divergence when  $\rho = 0.25$  and  $0 \leq \hat{\rho}_i^{(l)} \leq 1$



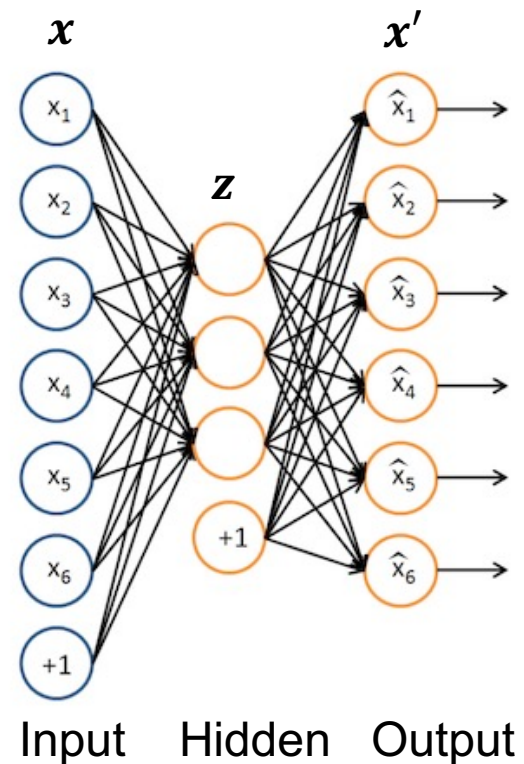
# Visualization of Autoencoder

Train a sparse autoencoder on

- Use 10x10 images and 100 hidden units
- Each hidden unit  $i$  computes a function of the input

$$z_i = \sigma\left(\sum_{j=1}^{100} W_{ij}x_j + b_i\right)$$

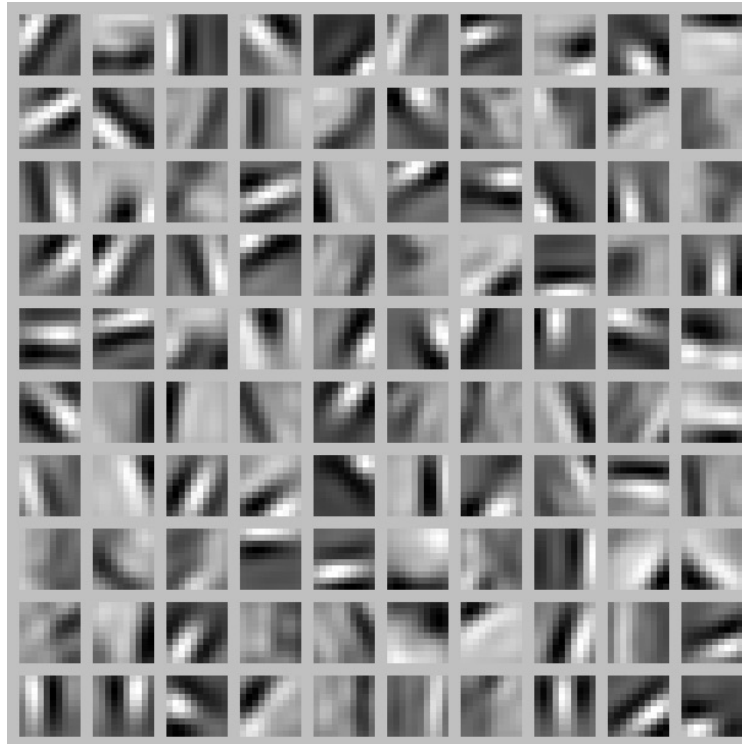
- $z_i$  is a non-linear feature of input  $x$
- What input image  $x$  would cause  $z_i$  to be maximally activated?  
(i.e., what is the feature that hidden unit  $i$  is looking for?)



# Visualization of Autoencoder

The input that maximally activates hidden unit  $i$  is given by setting pixel  $x_j$  (for all  $j = 1, \dots, 100$  pixels)

$$x_j = \frac{W_{ij}}{\sqrt{\sum_{j=1}^{100} (W_{ij})^2}} \quad \text{with} \quad \|x\|^2 \leq 1$$



100 such images,  
one per hidden unit

# Outline

- Autoencoders
- Variational Autoencoders

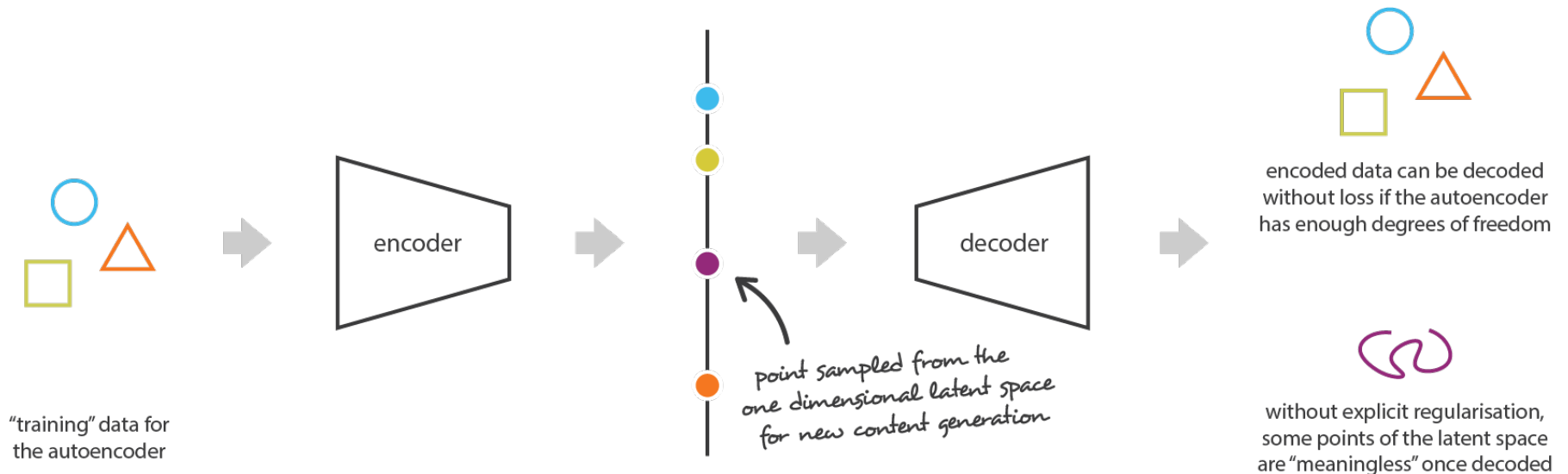
# Limitation of AE for Generation

AE is trained to encode and decode with the minimum loss

- Do NOT care how the latent space is organized

## Severe overfitting

- Due to too large degree-of-freedom, some points of the latent space will give meaningless content once decoded





# Idea of VAE

VAE follows AE

- Composed of both an encoder and a decoder
- Trained to minimize the reconstruction error between the input and the encoded-decoded data

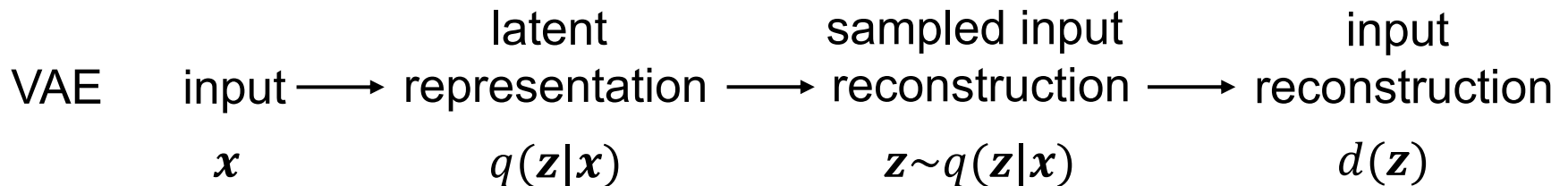
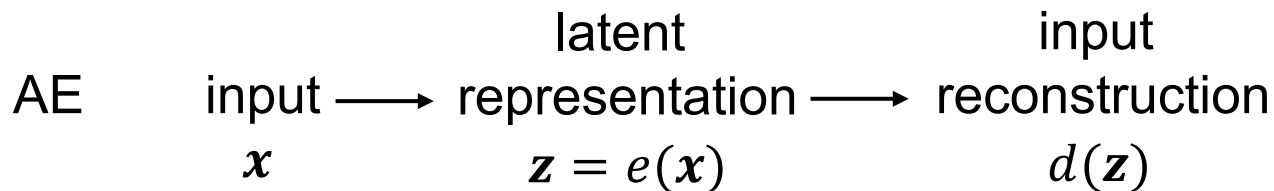
Instead, its training is regularized

- The latent space has good properties that enable generative process
- One key difference: instead of encoding an input as a single point, we encode it as a **distribution** over the latent space

# Idea of VAE

## Basic steps

- (1) Encode the input as distribution over the latent space
- (2) Sample a point from the latent distribution
- (3) Decode the sampled point
- (4) Backpropagate the reconstruction error through the network



# VAE from AE Perspective

Consists of an encoder, a decoder, and a loss function

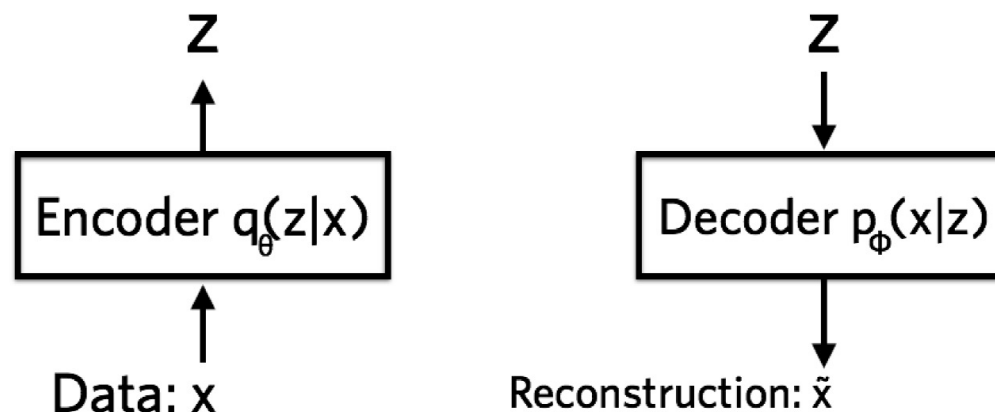
- Encoder and decoder are neural networks with parameter  $\theta$  and  $\phi$

Encoder  $q_{\theta}(z|x)$

- Input: a data point  $x$ , output: its low-dim representation  $z$

Decoder  $p_{\phi}(x|z)$

- Input: a representation  $z$ , output: a data point  $x$



# VAE

Loss function = reconstruction loss + regularization

- The reconstruction term is the same with AE
- The regularization term organizes the latent space by making distributions returned by the encoder close to normal distribution
- The loss for a data point is

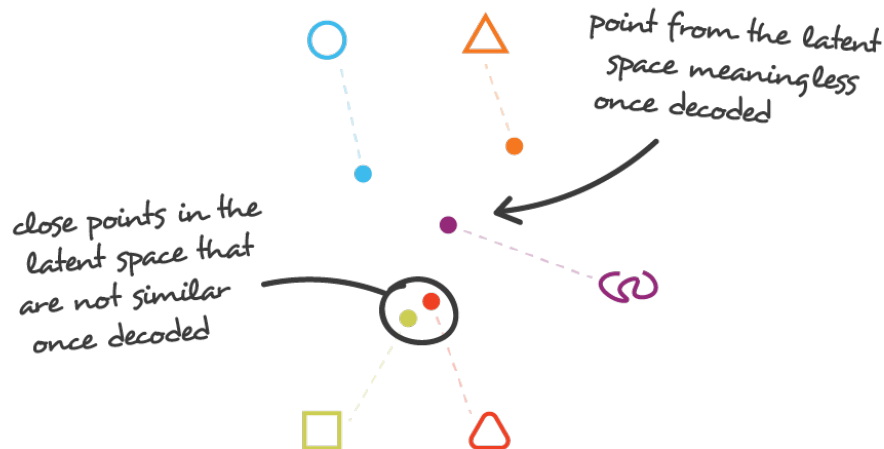
$$l_i(\theta, \phi) = -\mathbb{E}_{\mathbf{z} \sim q_\theta(\mathbf{z}|\mathbf{x}_i)} [\log p_\phi(\mathbf{x}_i|\mathbf{z})] + KL(q_\theta(\mathbf{z}|\mathbf{x}_i) || p(\mathbf{z}))$$

- Total loss is a summation of individual loss  $\sum_{i=1}^m l_i$

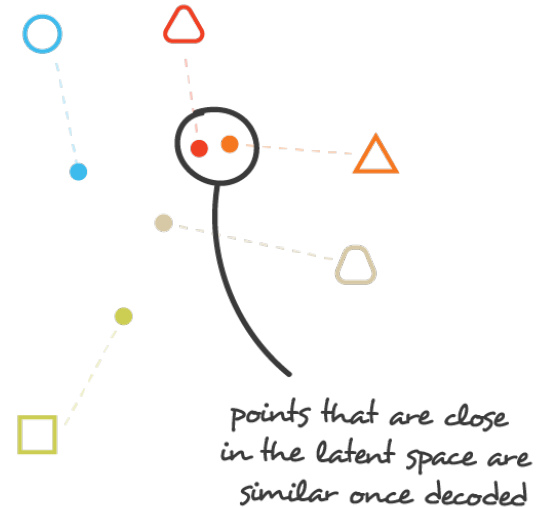
# Intuition about Regularization

## Two main properties for regularity of latent space

- Continuity: two close points in the latent space should not give two completely different contents once decoded
- Completeness: a point sampled from the latent distribution should give “meaningful” content once decoded



irregular latent space



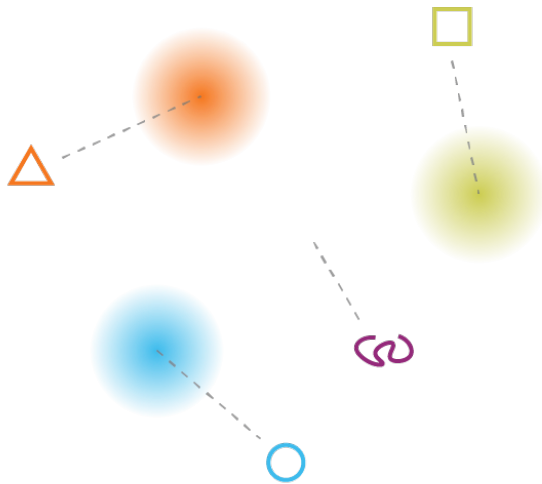
regular latent space



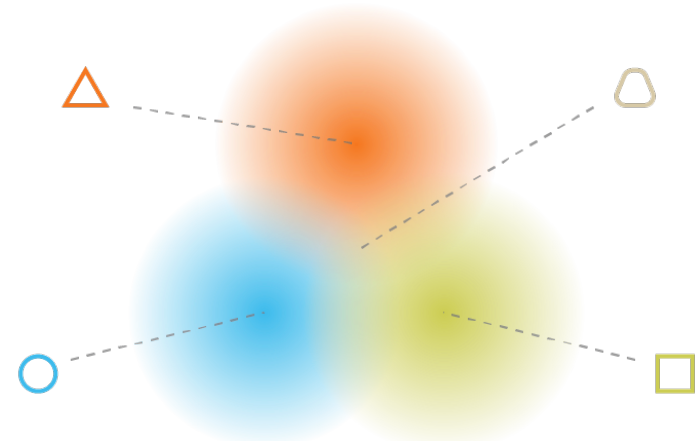
# Intuition about Regularization

Regularization enforces distributions to be close to a standard normal (centered and reduced)

- The covariance matrices to be close to the identity (not to be skinny or focused on a point)
- The mean to be close to 0 (prevent too far apart from each others)



what can happen without regularisation

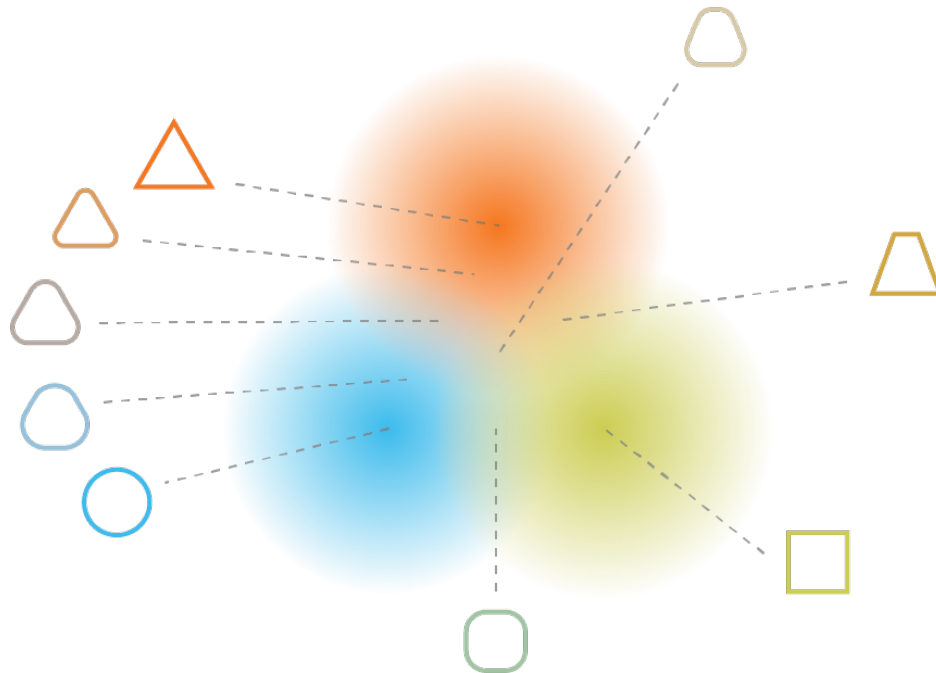


what we want to obtain with regularisation

# Intuition about Regularization

Overlapped distribution is encouraged

- To satisfy the expected continuity and completeness conditions
- However, any regularization including this comes at the price of a higher reconstruction error on the training data



# VAE

Loss function = reconstruction loss + regularization

- The loss for a data point is

$$l_i(\theta, \phi) = -\mathbb{E}_{\mathbf{z} \sim q_\theta(\mathbf{z}|\mathbf{x}_i)} [\log p_\phi(\mathbf{x}_i|\mathbf{z})] + KL(q_\theta(\mathbf{z}|\mathbf{x}_i) || p(\mathbf{z}))$$

The 1st term: negative log-likelihood

- The expectation is taken w.r.t the encoder's distribution over the representations
- If the decoder's output does not reconstruct the data well, it will incur a large cost in this loss function



# VAE

Loss function = reconstruction loss + regularization

- The loss for a data point is

$$l_i(\theta, \phi) = -\mathbb{E}_{\mathbf{z} \sim q_\theta(\mathbf{z}|\mathbf{x}_i)} [\log p_\phi(\mathbf{x}_i|\mathbf{z})] + KL(q_\theta(\mathbf{z}|\mathbf{x}_i) || p(\mathbf{z}))$$

The 2nd term: regularization

- KL divergence between the decoder's distribution  $q_\theta(\mathbf{z}|\mathbf{x}_i)$  and actual  $p(\mathbf{z})$
- Measure how much information is lost when using  $q$  to represent  $p$  (i.e., how close  $q$  is to  $p$ )
- In VAE,  $p(\mathbf{z}) = \text{Normal}(\mathbf{0}, \mathbf{1})$
- Make the representation space of  $\mathbf{z}$  meaningful (i.e., if the encoder output is different from standard normal, it is penalized)

# VAE from PGM Perspective

A probabilistic model of data  $\mathbf{x}$  and latent variable  $\mathbf{z}$

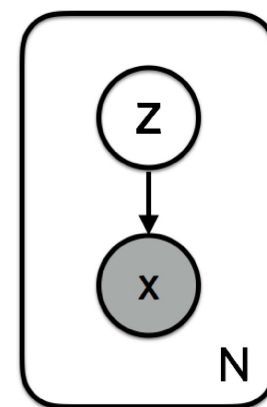
- Joint probability of the model

$$\begin{aligned} P(\mathbf{x}, \mathbf{z}) &= p(\mathbf{x}|\mathbf{z})p(\mathbf{z}) \\ &= (\text{likelihood}) \times (\text{prior}) \end{aligned}$$

- Data generating process

For each data point  $i$

- Draw latent variable  $\mathbf{z}_i \sim p(\mathbf{z})$
- Draw data point  $\mathbf{x}_i \sim p(\mathbf{x}|\mathbf{z}_i)$



# VAE from PGM Perspective

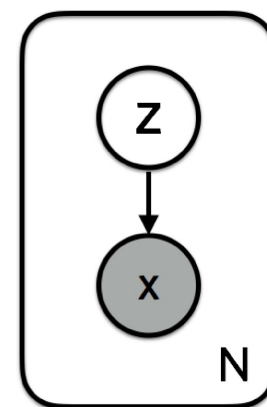
Goal: calculate posterior  $p(\mathbf{z}|\mathbf{x})$

- Infer good values of the latent variables given observed data

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

- Computing denominator  $p(\mathbf{x})$  is intractable!

$$P(\mathbf{x}) = \int p(\mathbf{x}|\mathbf{z})p(\mathbf{z})d\mathbf{z}$$



Solution: use variational inference

- Approximate the posterior  $P(\mathbf{z}|\mathbf{x})$  with  $q_{\theta}(\mathbf{z}|\mathbf{x})$  of a known distribution with parameter  $\theta$
- How to decide whether  $q_{\theta}(\mathbf{z}|\mathbf{x})$  approximates  $P(\mathbf{z}|\mathbf{x})$  well?

# VAE from PGM Perspective

Find  $\theta$  that minimizes  $KL(q_\theta(\mathbf{z}|\mathbf{x})||p(\mathbf{z}|\mathbf{x}))$

$$q_\theta^*(\mathbf{z}|\mathbf{x}) = \underset{\theta}{\operatorname{argmin}} KL(q_\theta(\mathbf{z}|\mathbf{x})||p(\mathbf{z}|\mathbf{x}))$$

- The KL divergence becomes

$$\begin{aligned} KL(q_\theta(\mathbf{z}|\mathbf{x})||p(\mathbf{z}|\mathbf{x})) &= \sum q_\theta(\mathbf{z}|\mathbf{x}) \log \frac{q_\theta(\mathbf{z}|\mathbf{x})}{p(\mathbf{z}|\mathbf{x})} = \mathbb{E}_q \left[ \log \frac{q_\theta(\mathbf{z}|\mathbf{x})p(\mathbf{x})}{p(\mathbf{x}, \mathbf{z})} \right] \\ &= \mathbb{E}_q [\log q_\theta(\mathbf{z}|\mathbf{x})] - \mathbb{E}_q [\log p(\mathbf{x}, \mathbf{z})] + \log p(\mathbf{x}) \end{aligned}$$

- Now re-organize it

$$\log p(\mathbf{x}) = \underbrace{\mathbb{E}_q [\log p(\mathbf{x}, \mathbf{z})] - \mathbb{E}_q [\log q_\theta(\mathbf{z}|\mathbf{x})]}_{\text{ELBO (Evidence Lower BOund)}} + KL(q_\theta(\mathbf{z}|\mathbf{x})||p(\mathbf{z}|\mathbf{x})) \geq 0$$

Minimizing KL divergence = Maximizing ELBO!
---

# VAE from PGM Perspective

Goal: maximize ELBO (tractable)

- For each datapoint  $i$

$$\begin{aligned}\text{ELBO}_i(\lambda) &= \mathbb{E}_q[\log p(\mathbf{x}_i, \mathbf{z})] - \mathbb{E}_q[\log q_\theta(\mathbf{z}|\mathbf{x}_i)] \\ &= \mathbb{E}_q[\log p(\mathbf{x}_i|\mathbf{z})] + \mathbb{E}_q[\log p(\mathbf{z})] - \mathbb{E}_q[\log q_\theta(\mathbf{z}|\mathbf{x}_i)] \\ &= \mathbb{E}_q[\log p(\mathbf{x}_i|\mathbf{z})] - KL(q_\theta(\mathbf{z}|\mathbf{x}_i)||p(\mathbf{z}))\end{aligned}$$

Learn  $\theta, \phi$  that maximize ELBO

$$\text{ELBO}_i(\theta, \phi) = \mathbb{E}_q[\log p_\phi(\mathbf{x}_i|\mathbf{z})] - KL(q_\theta(\mathbf{z}|\mathbf{x}_i)||p(\mathbf{z}))$$

- Remind that  $p_\phi(\mathbf{x}_i|\mathbf{z})$  is the decoder, and  $q_\theta(\mathbf{z}|\mathbf{x}_i)$  is the encoder (both are neural networks)

# PGM and AE Perspective

AE perspective: minimize loss

$$l_i(\theta, \phi) = -E_{\mathbf{z} \sim q_\theta(\mathbf{z}|\mathbf{x}_i)}[\log p_\phi(\mathbf{x}_i|\mathbf{z})] + KL(q_\theta(\mathbf{z}|\mathbf{x}_i)||p(\mathbf{z}))$$

PGM perspective: maximize ELBO

$$\text{ELBO}_i(\theta, \phi) = \mathbb{E}_q[\log p_\phi(\mathbf{x}_i|\mathbf{z})] - KL(q_\theta(\mathbf{z}|\mathbf{x}_i)||p(\mathbf{z}))$$

They are equivalent  $\text{ELBO}_i(\theta, \phi) = -l_i(\theta, \phi)$

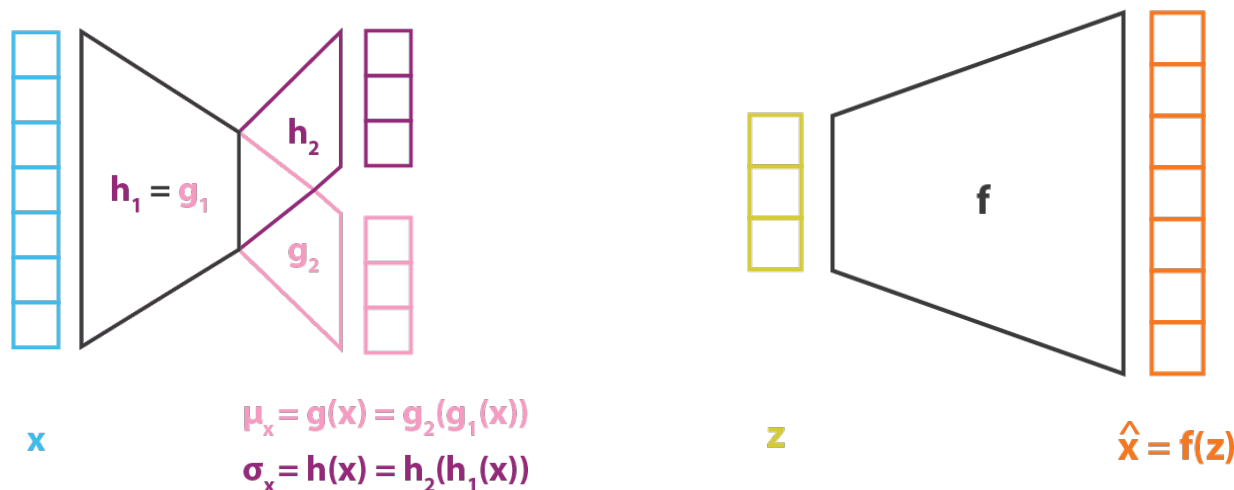
# Modeling Encoder and Decoder

## Encoder $q_{\theta}(\mathbf{z}|\mathbf{x}_i)$

- The encoder is a neural network that takes  $x_i$  and outputs the mean  $\mu_i$  and covariance  $\sigma_i$  of the multivariate Gaussian

## Decoder $p_\phi(\mathbf{x}_i|\mathbf{z})$

- The decoder is another neural network that takes a sampled  $\mathbf{z}$  from Gaussian and outputs a reconstructed  $\hat{\mathbf{x}}_i$



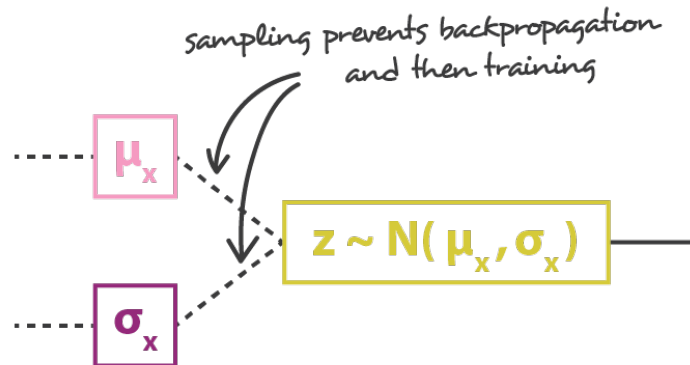
# Reparametrization Trick

VAE is a concatenation of encoder and decoder

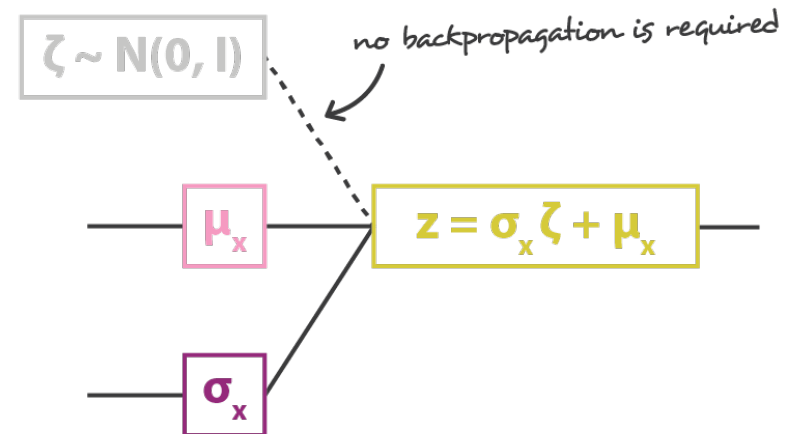
- Sampling is a stochastic process and therefore we cannot backpropagate the gradient
- A simple trick makes the gradient descent possible despite the random sampling for  $z$  that occurs halfway of the architecture

—— no problem for backpropagation

----- backpropagation is not possible due to sampling



sampling without reparametrisation trick



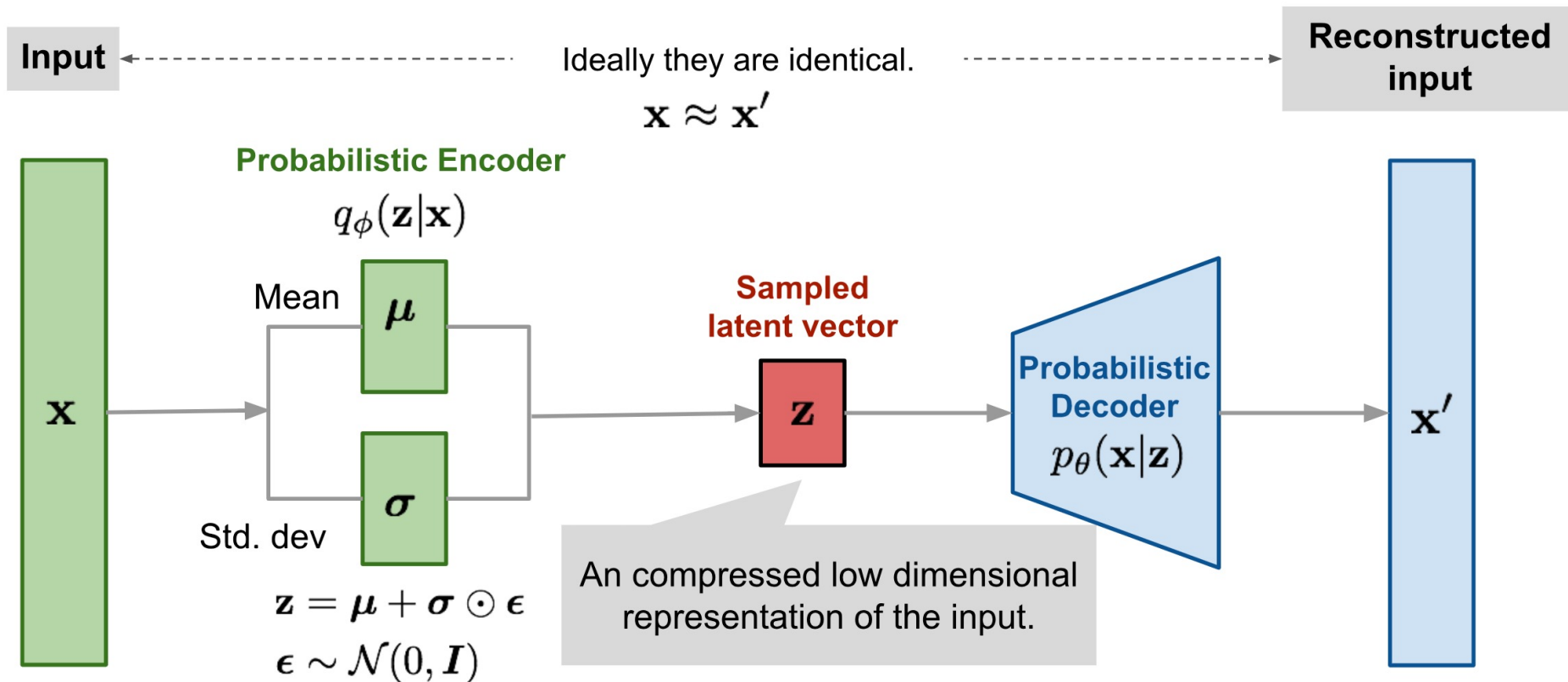
sampling with reparametrisation trick



# The Overall Architecture

End-to-end training is possible

- Use L2 distance between  $x_i$  and  $\hat{x}_i$  for reconstruction loss



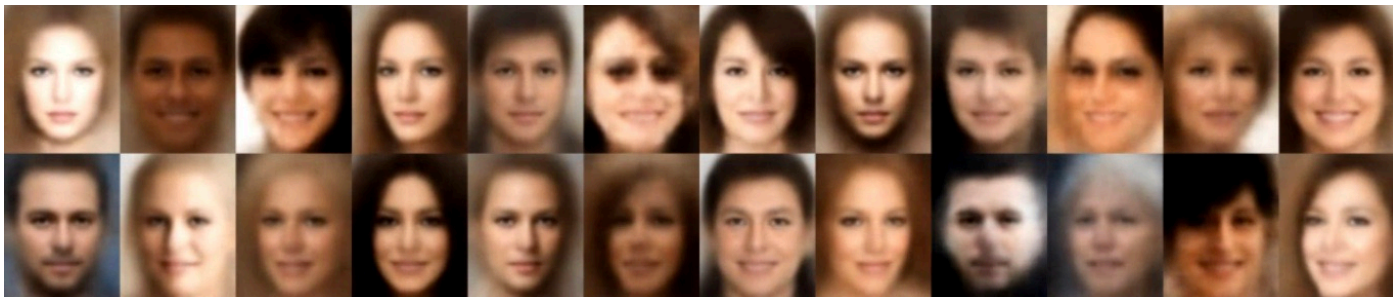
# VAE vs. GAN

## Downside of GAN

- Images are generated off some arbitrary noise
  - It is not straightforward to generate with specific features
- GAN only discriminates between real and fake images
  - No constraints that an image of a cat must look like a cat
  - No actual object in a generated image, but the style just looks like a cat picture

## Downside of VAE

- Output images are blurry; it uses direct mean squared errors



# Summary

AE is good for dimensionality reduction

- AE is a neural network composed of an encoder and a decoder
- Create a bottleneck to go through for data
- Trained to lose a minimal quantity of information during the encoding-decoding process (i.e., reduce the reconstruction error)
- Due to overfitting, the latent space can be extremely irregular
- Hard to use for a generative process

VAE tackles the AE's problem of latent space irregularity

- Encoder returns a distribution over the latent space instead of a single point
- Loss function additionally includes a regularization for better organization of the latent space
- Derived from the technique of variational inference