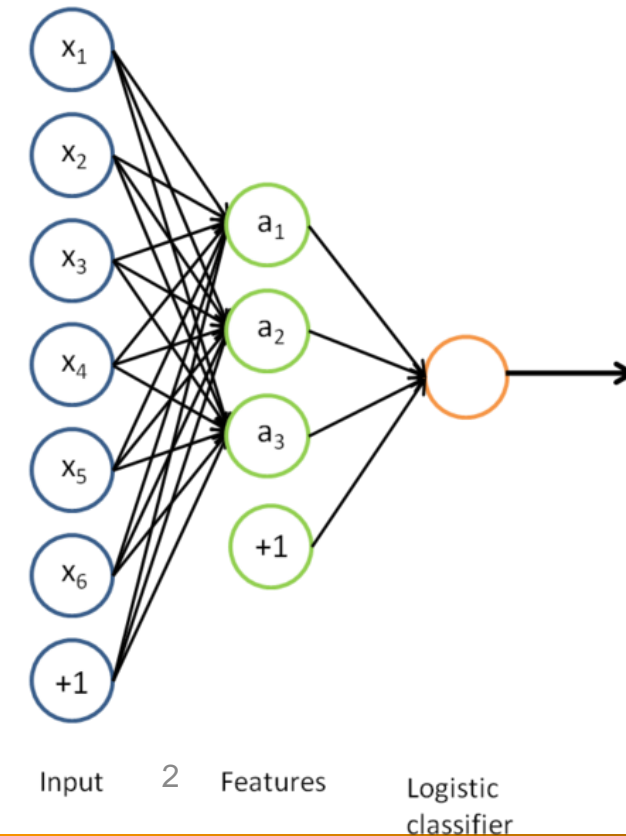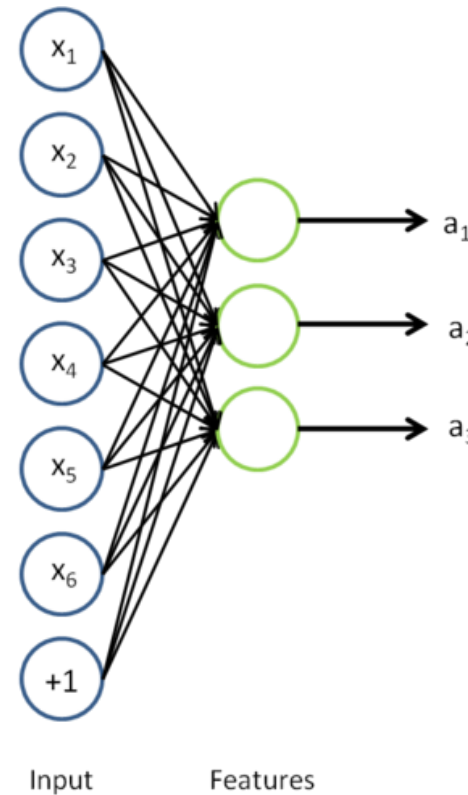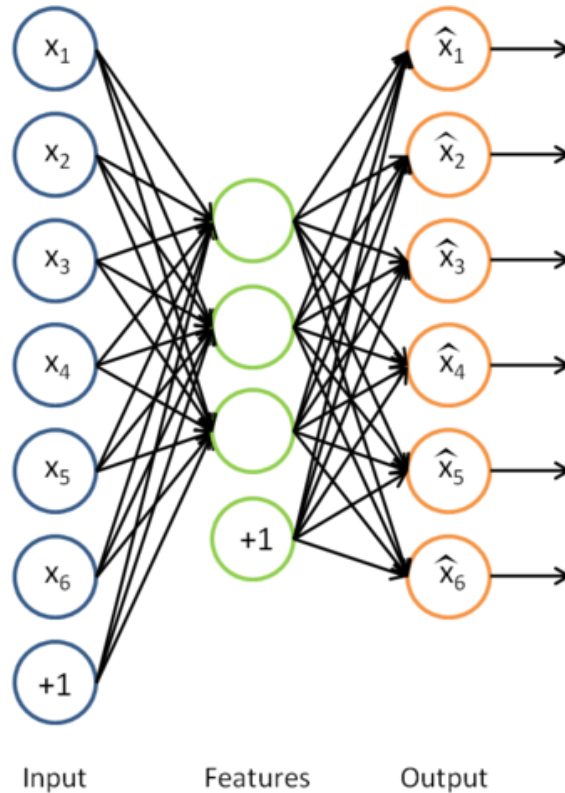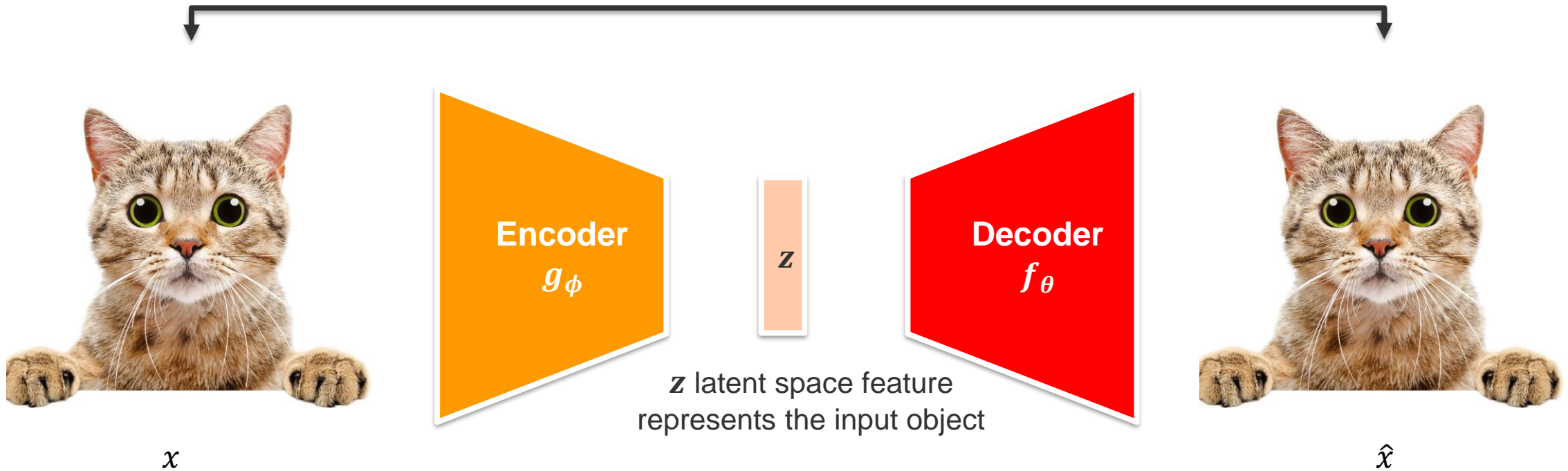# AutoEncoder

# Auto-Encoders

- A type of unsupervised learning which tries to discover generic features of the data
  - Learn identity function by learning important sub-features (not by just passing through data)
  - Compression, etc.
  - Can use just new features in the new training set or concatenate both

# Autoencoder

Reconstruction Loss $\qquad \mathcal{L}(\phi, \theta) = \left( x - f_\theta \left( g_\phi(x) \right) \right)^2$



**Encoder** $g_\phi$

**z**

**Decoder** $f_\theta$

**z** latent space feature represents the input object

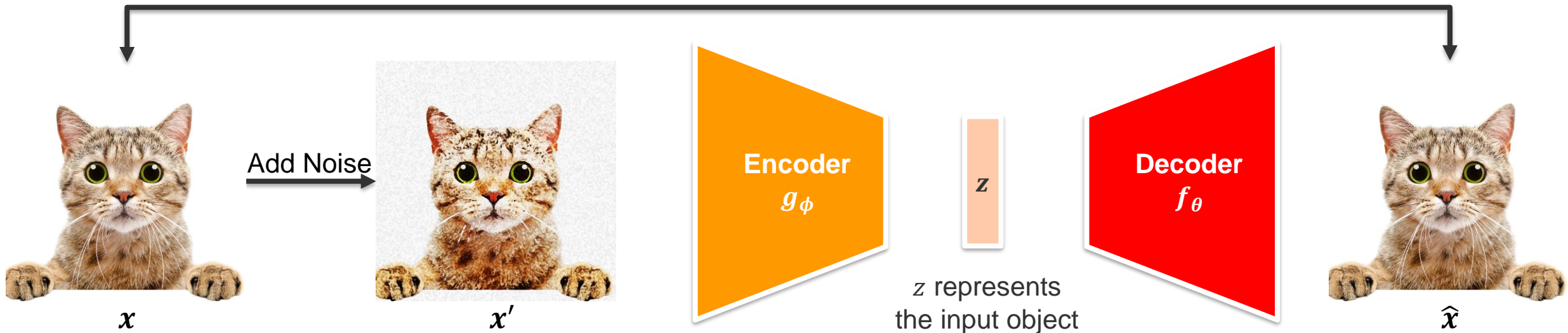$x$ $\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad \hat{x}$

- The decoder tries to reconstruct the input image based on the encoding of the input
- To reconstruct the image, the encoder <u>should</u> learn good representations of the input
- The self-supervised pre-trained model can then be used to train on the target task

# Denoising Autoencoder [2]

Reconstruction Loss $\quad \mathcal{L}(\phi, \theta) = \left(x - f_\theta\left(g_\phi(x')\right)\right)^2$ 即便加了雜訊所取出之特徵，仍要能還原



Add Noise

$x$       $x'$

**Encoder** $g_\phi$

$z$

$z$ represents the input object
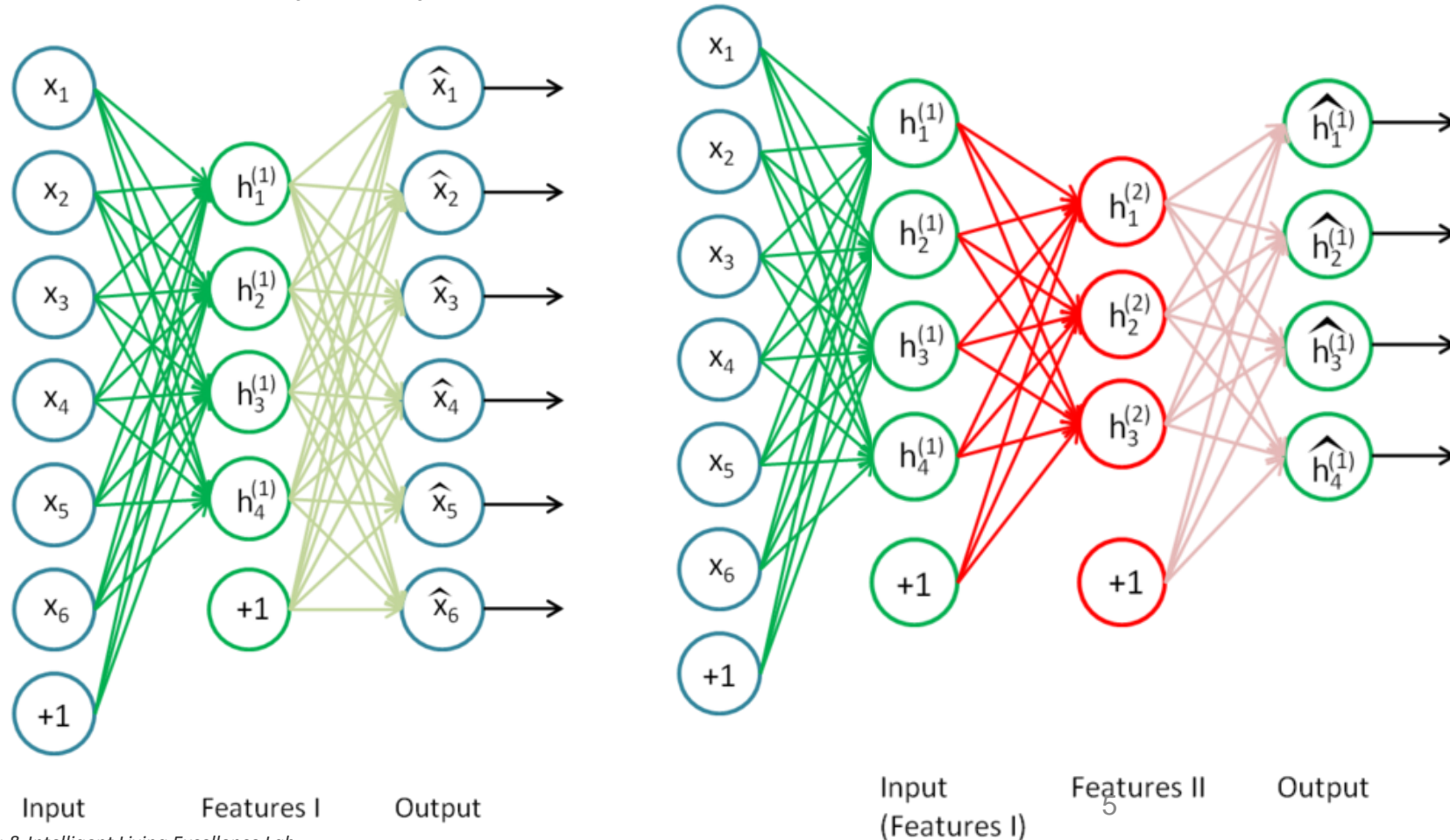
**Decoder** $f_\theta$

$\hat{x}$

- Sometimes the autoencoder with clean input only learns low level patterns instead of high level object representations
- Use noisy input instead and require the decoder to remove the noise
- To remove the noise, the decoder needs to know the object, and the encoder should learn better representations of the object
- The self-supervised pre-trained model can then be used to train on the target task
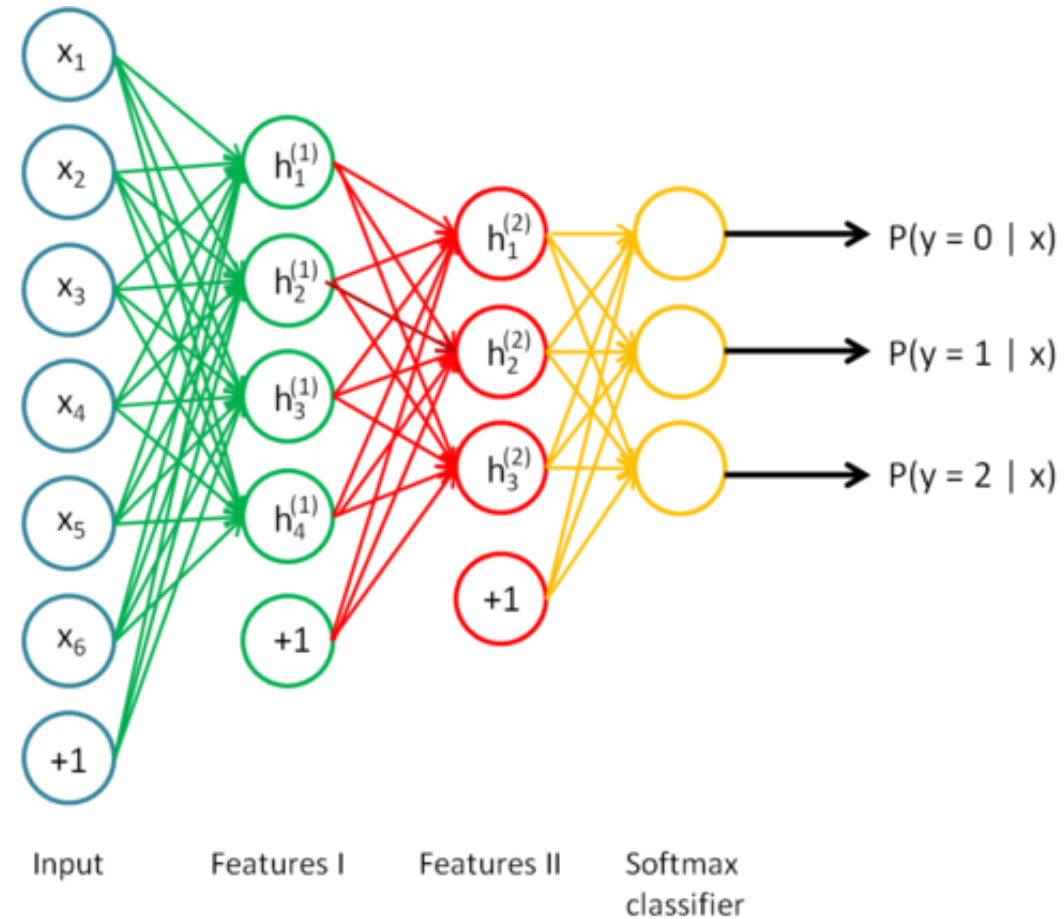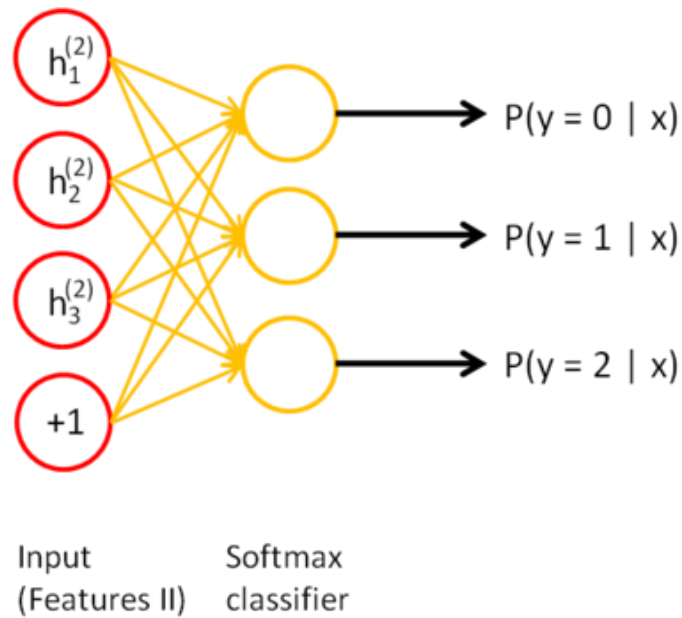
# Stacked Auto-Encoders

- Bengio (2007)
- Stack many (sparse) auto-encoders in succession and train them using greedy layer-wise training
- Drop the decode output layer each time

# Stacked Auto-Encoders

- Do supervised training on the last layer using final features
- Then do supervised training on the entire network to fine- tune all weights
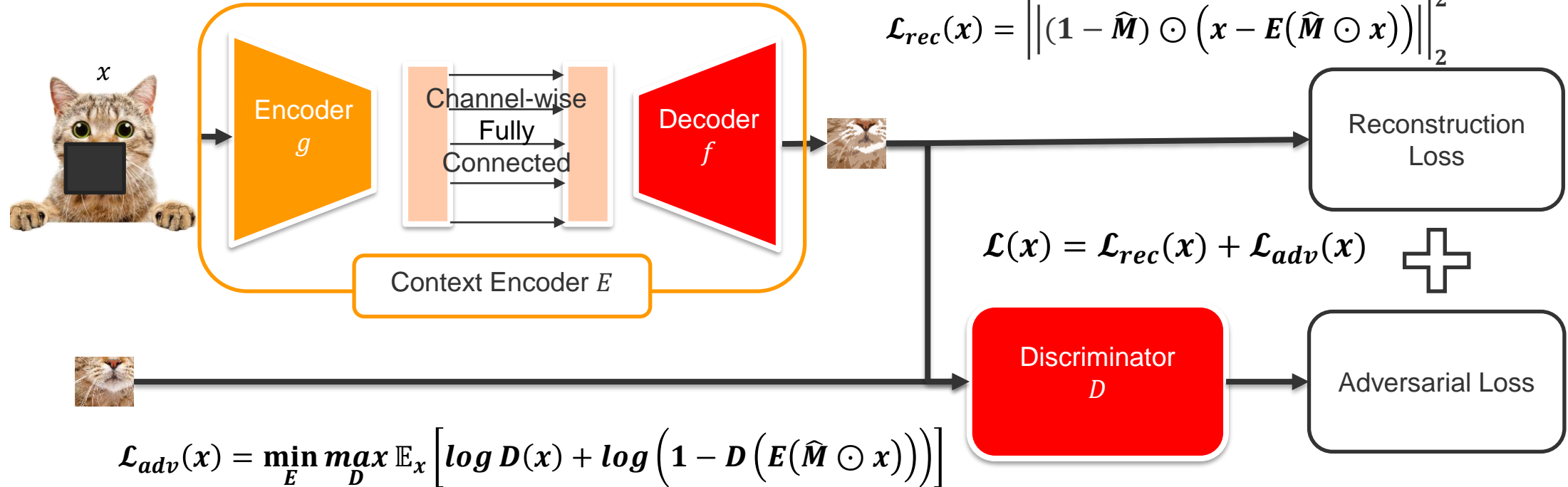
# Image Inpainting

**M: a bit map with masked 0, other regions 1**

$$\mathcal{L}_{rec}(x) = \left\| (1 - \widehat{M}) \odot \left( x - E(\widehat{M} \odot x) \right) \right\|_2^2$$

$x$

Encoder $g$

Channel-wise Fully Connected

Decoder $f$

Context Encoder $E$

Reconstruction Loss

$$\mathcal{L}(x) = \mathcal{L}_{rec}(x) + \mathcal{L}_{adv}(x)$$

Discriminator $D$

Adversarial Loss

$$\mathcal{L}_{adv}(x) = \min_E \max_D \mathbb{E}_x \left[ \log D(x) + \log \left( 1 - D\left( E(\widehat{M} \odot x) \right) \right) \right]$$

**Pixel wise reconstruction does not give semantic reconstruction,**
➔ **image 較不清晰**

- Part of the input image is cropped, and the decoder tries to reconstruct the missing part
- To make the reconstructed image looks more like real, a discriminator is used to identify real images and generated images, and the decoder should generate more realistic images to deceive the discriminator
- The self-supervised pre-trained model can then be used to train on the target task
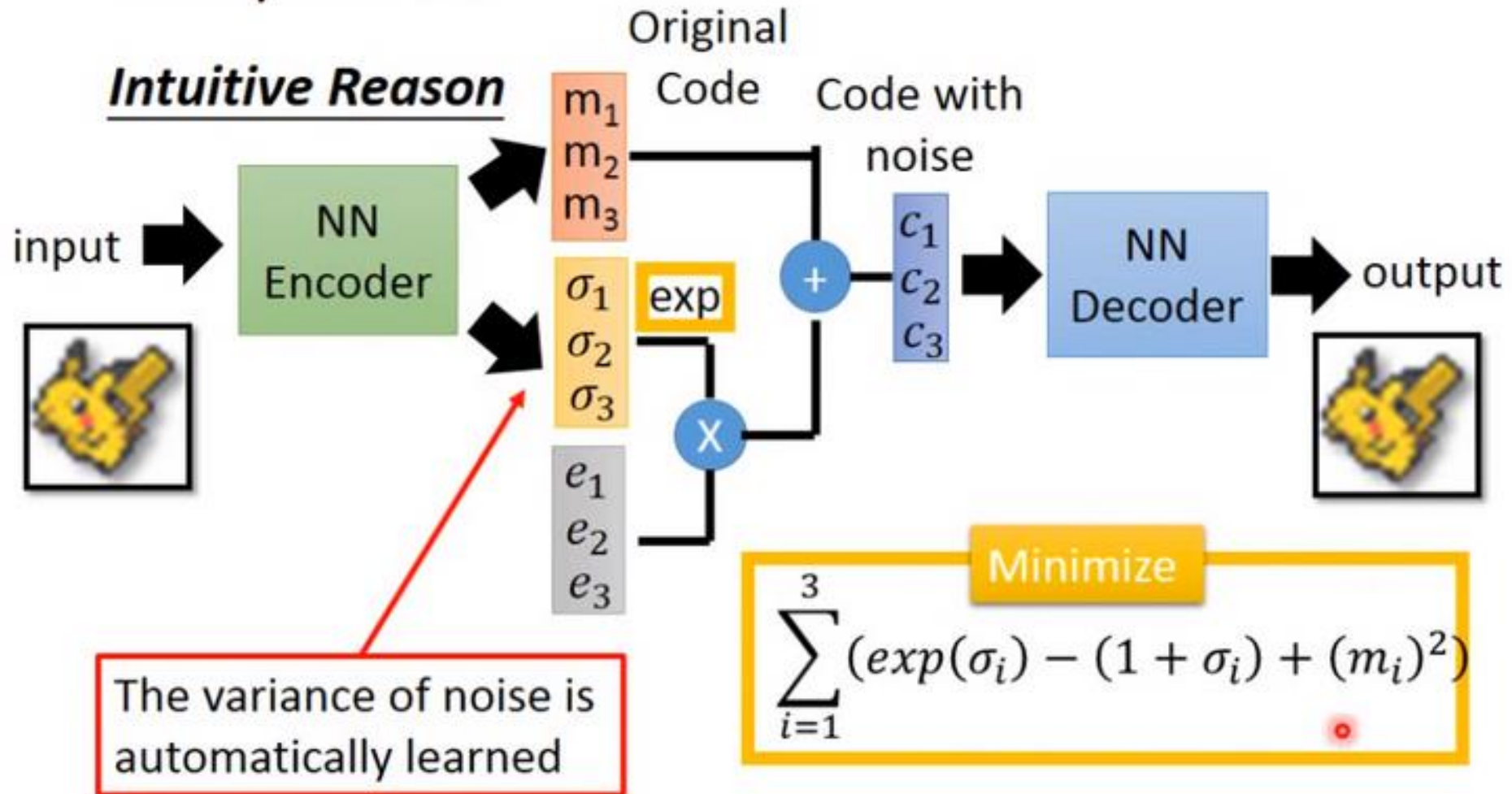
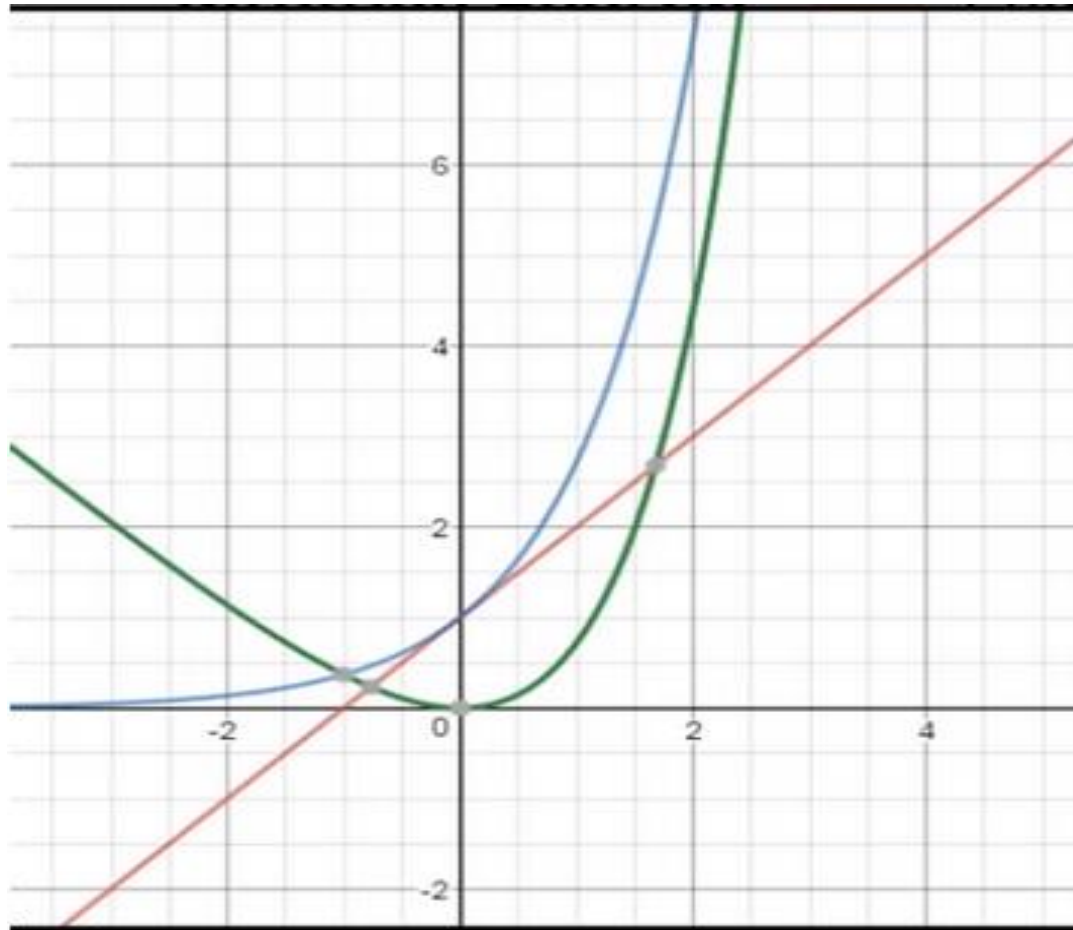# VAE(Variational AutoEncoder)

# Why VAE?

**What will happen if we only minimize reconstruction error?**

**Intuitive Reason**



Minimize

$$\sum_{i=1}^{3} (exp(\sigma_i) - (1 + \sigma_i) + (m_i)^2)$$
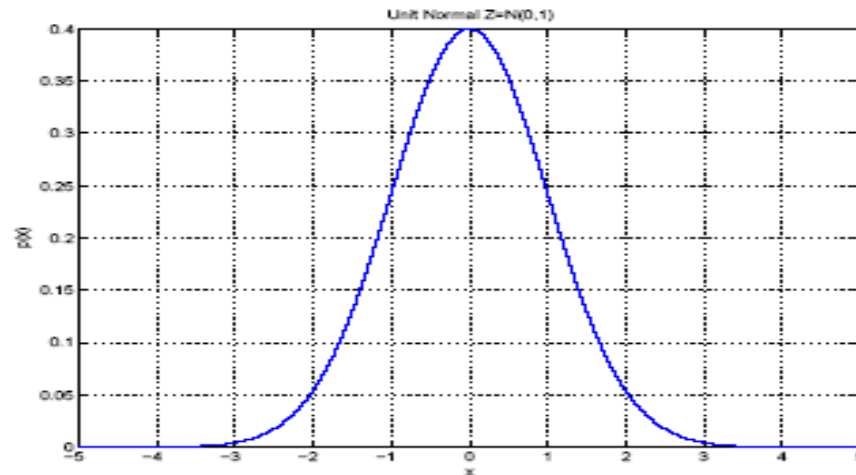
The variance of noise is automatically learned

$\delta_i = 0$, $\exp(\delta_i)=1$ has the minimum value

$C_i = \exp(\delta_i)$ x $e_i + m_i$ 依然加了 noises
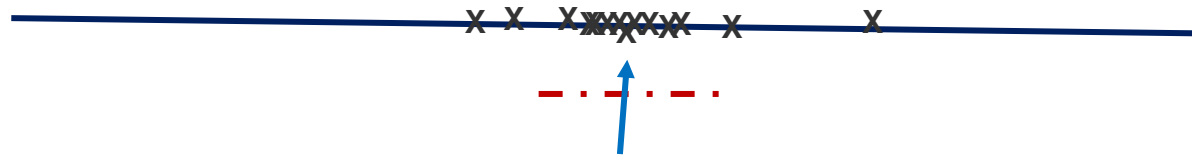
# Maximum likelihood estimation

When there is no reason to favor certain $\theta$



- $p(x) = \mathcal{N} ( \mu_1, \sigma_1{}^2)$

Most data cannot get the high prob.
Only small amount of data have high prob.

x  x  xxxxxx  x        x

More data points on this high prob. area

11

For a correct estimation, it is expected to have the more data achieving high probability

➔ Likelihood of $\theta$ given the sample $\mathcal{X}$

Maximize   $p\ (\mathcal{X}\,|\vartheta) = \prod_t p\ (x^t|\vartheta)$

should be maximized

Bayes rule:

$$P(\boldsymbol{\theta}\,|\,A) = \frac{P(\boldsymbol{\theta} \bigcap A)}{P(A)}$$

$$= \frac{P(A\,|\,\boldsymbol{\theta})P(\boldsymbol{\theta})}{P(A)}$$

- To model X with parametrized distribution $P_\theta$

- Let Z represent a latent encoding of X

- $P_\theta(x,z)$ represents the [joint distribution](#) under $P_\theta$ of the observable data and its latent space z, where
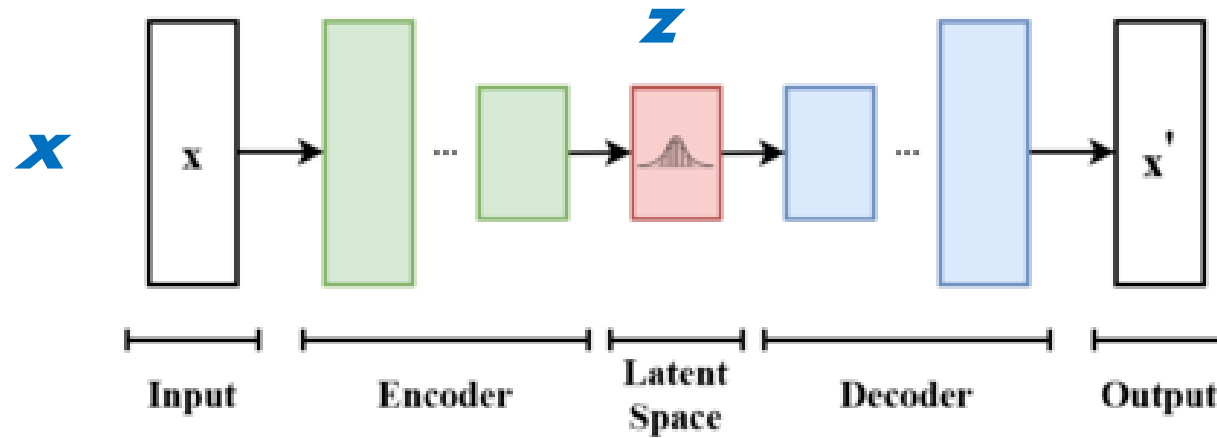
$$p_\theta(\mathbf{x}) = \int_{\mathbf{z}} p_\theta(\mathbf{x} \mid \mathbf{z}) p_\theta(\mathbf{z}) \, d\mathbf{z}$$

Define the set of relationships between the input data and its latent representation:

- Prior $p_\theta(\mathbf{z})$
- Likelihood $p_\theta(\mathbf{x} \mid \mathbf{z})$
- Posterior $p_\theta(\mathbf{z} \mid \mathbf{x})$

找個模型來實作 coder
$$q_\emptyset(z|x) \approx p_\theta(z|x)$$

但我們擁有的只有所有的 X
Or, P(x)

Z is the code,
in the latent space

$$D_{KL}(q_\emptyset(z|x)\|p_\theta(z|x)) = \int (q_\emptyset(z|x) \log \frac{q_\emptyset(z|x)}{p_\theta(z|x)} dz \quad = \int (q_\emptyset(z|x) \log \frac{q_\emptyset(z|x)p_\theta(x)}{p_\theta(z,x)} dz$$

$$= \int (q_\emptyset(z|x) \left( \log(p_\theta(x)) + \log \frac{q_\emptyset(z|x)}{p_\theta(z,x)} \right) dz$$

$$= \log(p_\theta(x)) + \int q_\emptyset(z|x) \log \frac{q_\emptyset(z|x)}{p_\theta(z,x)} dz$$

$$\log(p_\theta(x)) = - \int q_\emptyset(z|x) \log \frac{q_\emptyset(z|x)}{p_\theta(z,x)} dz + D_{KL}(q_\emptyset(z|x)\|p_\theta(z|x))$$

Maximize : $\underbrace{\log(p_\theta(x))}_{\text{Max}} - \underbrace{D_{KL}(q_\emptyset(z|x)\|p_\theta(z|x))}_{\text{Min}}$

$$= \underbrace{E_{Z\sim q_\emptyset(z|x)}\big(\log(p_\theta(x|z))\big)}_{\text{Max}} - \underbrace{D_{KL}(q_\emptyset(z|x)\|p_\theta(z))}_{\text{Min}}$$

lower bound

Minimize :

$$\underline{L_{\theta,\emptyset}} = -\log(p_\theta(x)) + D_{KL}(q_\emptyset(z|x)\|p_\theta(z|x))$$
$$= -E_{Z\sim q_\emptyset(z|x)}\big(\log(p_\theta(z|x))\big) + D_{KL}(q_\emptyset(z|x)\|p_\theta(z))$$

Evidence lower bound (ELBO) loss function

$$\theta^*, \emptyset^* = \underset{\theta,\emptyset}{\text{argmin}}\, L_{\theta,\emptyset}$$

$$-L_{\theta,\emptyset} = \log\big((p_\theta(x))\big) - D_{KL}(q_\emptyset(z|x)\|p_\theta(z|x)) \leq \log(p_\theta(x))$$

$KL$必定 $\geq 0$  =0 when $q_\emptyset(z|x) = p_\theta(z|x)$

# Maximizing Likelihood

**Connection with Network**

Minimizing $KL\big(q(z|x)||P(z)\big)$

$$\boxed{\sum_{i=1}^{3}\big(exp(\sigma_i) - (1+\sigma_i) + (m_i)^2\big)}$$

Minimize



$x \rightarrow$ NN' $\rightarrow \mu'(x),\ \sigma'(x)$

Maximizing

$$\int_z q(z|x)logP(x|z)dz = E_{q(z|x)}[logP(x|z)]$$

$x \rightarrow$ NN' $\rightarrow \mu'(x),\ \sigma'(x) \rightarrow z \rightarrow$ NN $\rightarrow \mu(x),\ \sigma(x)$

$\mu(x) \overset{close}{\longleftrightarrow} x$

**This is the auto-encoder**

Why VAE?

What will happen if we only minimize reconstruction error?

**Intuitive Reason**

input → NN Encoder → Original Code

$m_1, m_2, m_3$

Code with noise

$\sigma_1, \sigma_2, \sigma_3$  exp

$e_1, e_2, e_3$

$c_1, c_2, c_3$ → NN Decoder → output

The variance of noise is automatically learned

Minimize

$$\sum_{i=1}^{3} \left( exp(\sigma_i) - (1 + \sigma_i) + (m_i)^2 \right)$$

# Gaussian mixture Model

- GMM(Gaussian mixture Model)
- Considers each cluster as a different Gaussian distribution
- Mixture different distribution (blue line)
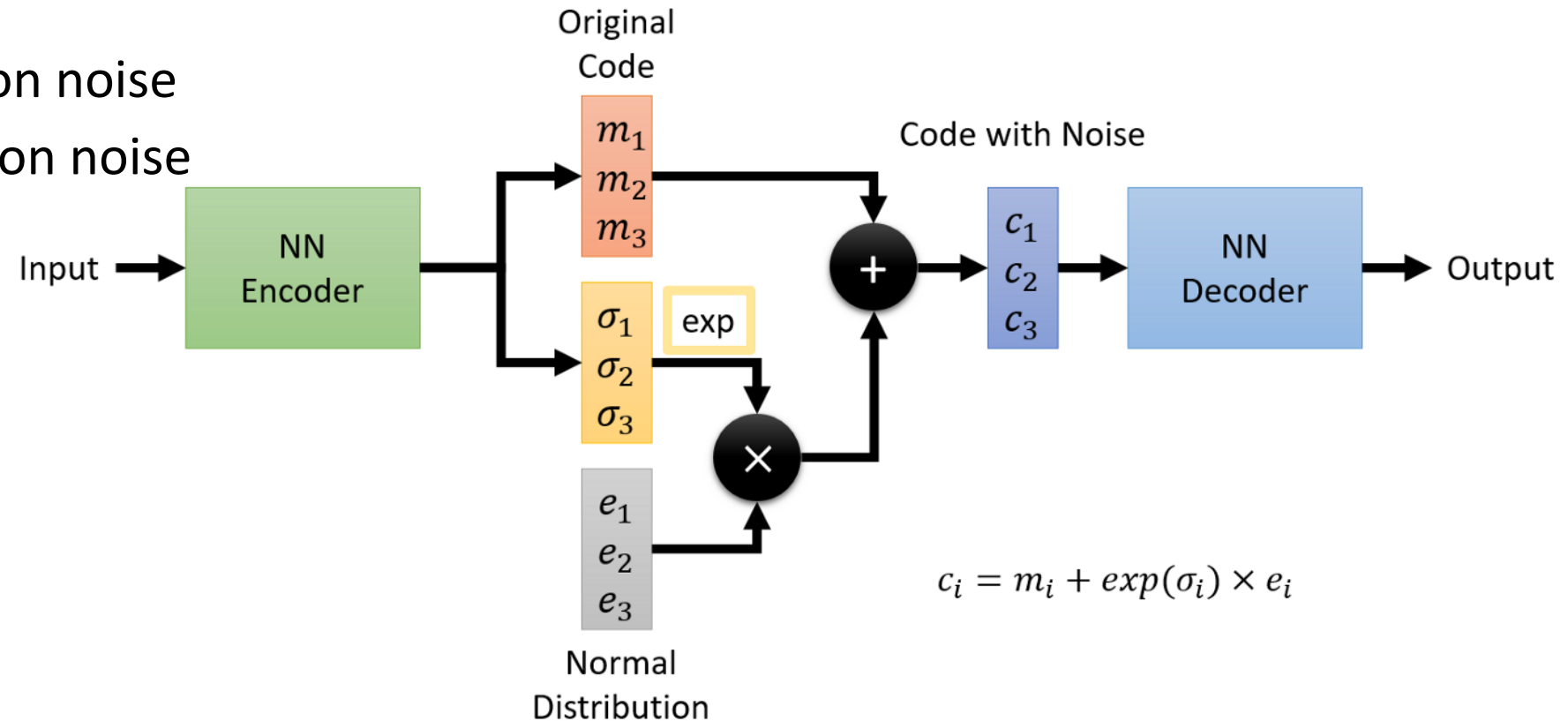  into new distribution(red line)

# VAE structure

利用 Normal distribution抽樣

σ:control weight of
   normal distribution noise
e: normal distribution noise



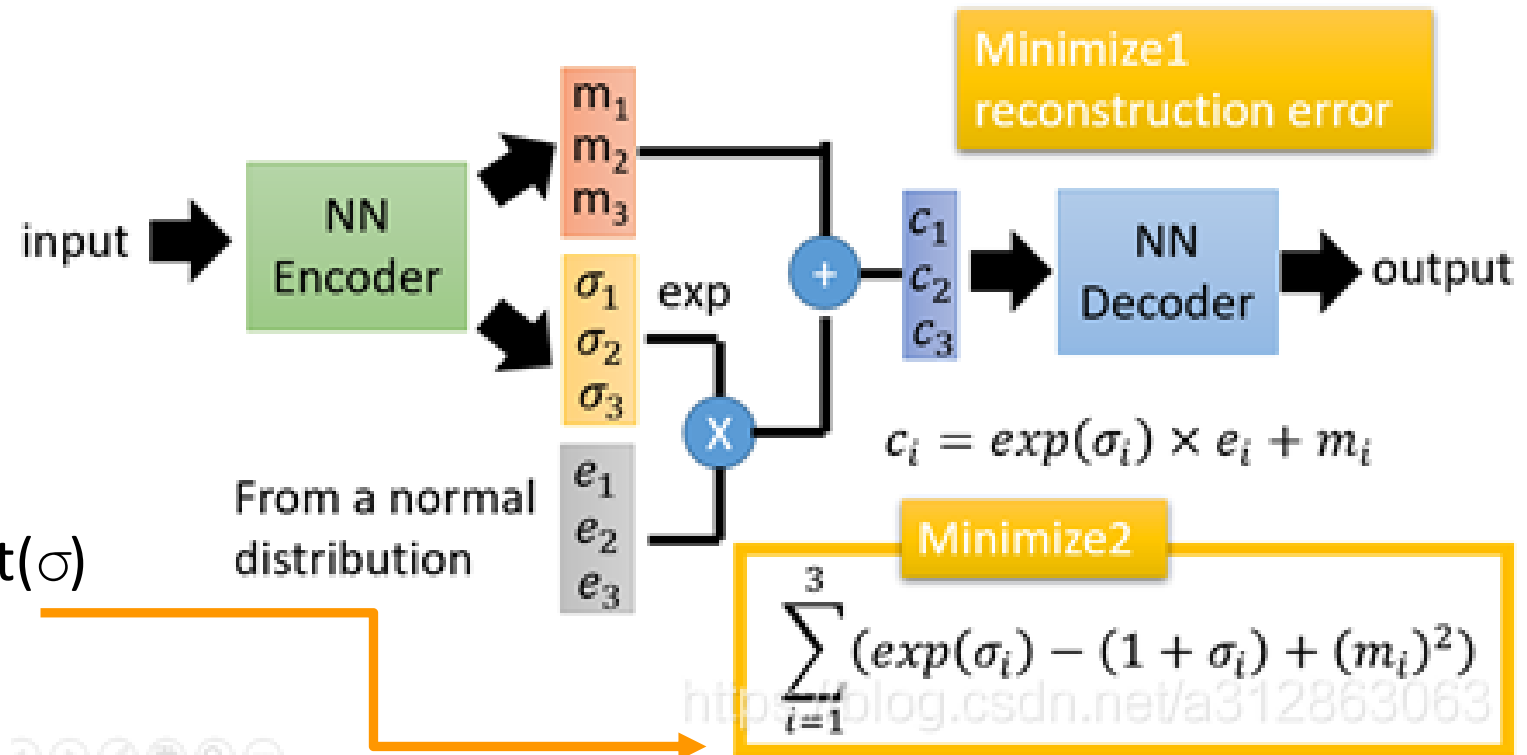$$c_i = m_i + exp(\sigma_i) \times e_i$$

# VAE structure

σ: control weight of
   normal distribution noise

e: normal distribution noise

Two loss for training
1. Reconstruction loss
2. Bounding of "noise" weight(σ)



$$c_i = exp(\sigma_i) \times e_i + m_i$$

Minimize1
reconstruction error

Minimize2

$$\sum_{i=1}^{3} (exp(\sigma_i) - (1 + \sigma_i) + (m_i)^2)$$

## Appendix:

$$\mathrm{D_{KL}}(q_\emptyset(z|x)\|p_\theta(z|x)) = \int (q_\emptyset(z|x)\log\frac{q_\emptyset(z|x)}{p_\theta(z|x)}dz$$

$$= \log(p_\theta(x)) + \int q_\emptyset(z|x)\log\frac{q_\emptyset(z|x)}{p_\theta(z,x)}dz$$

$$= \log(p_\theta(x)) + \int q_\emptyset(z|x)\log\frac{q_\emptyset(z|x)}{p_\theta(x|z)p_\theta(z)}dz$$

$$= \log(p_\theta(x)) + E_{Z \sim q_\emptyset(z|x)}\left(\log\frac{q_\emptyset(z|x)}{p_\theta(z)} - \log(p_\theta(x|z))\right)$$

$$= \log(p_\theta(x)) + E_{Z \sim q_\emptyset(z|x)}\left(\log\frac{q_\emptyset(z|x)}{p_\theta(z)} - \log(p_\theta)(x|z)\right)$$

$$= \log(p_\theta(x)) + D_{KL}(q_\emptyset(z|x)\|p_\theta(z)) - E_{Z \sim q_\emptyset(z|x)}(\log(p_\theta(x|z)))$$