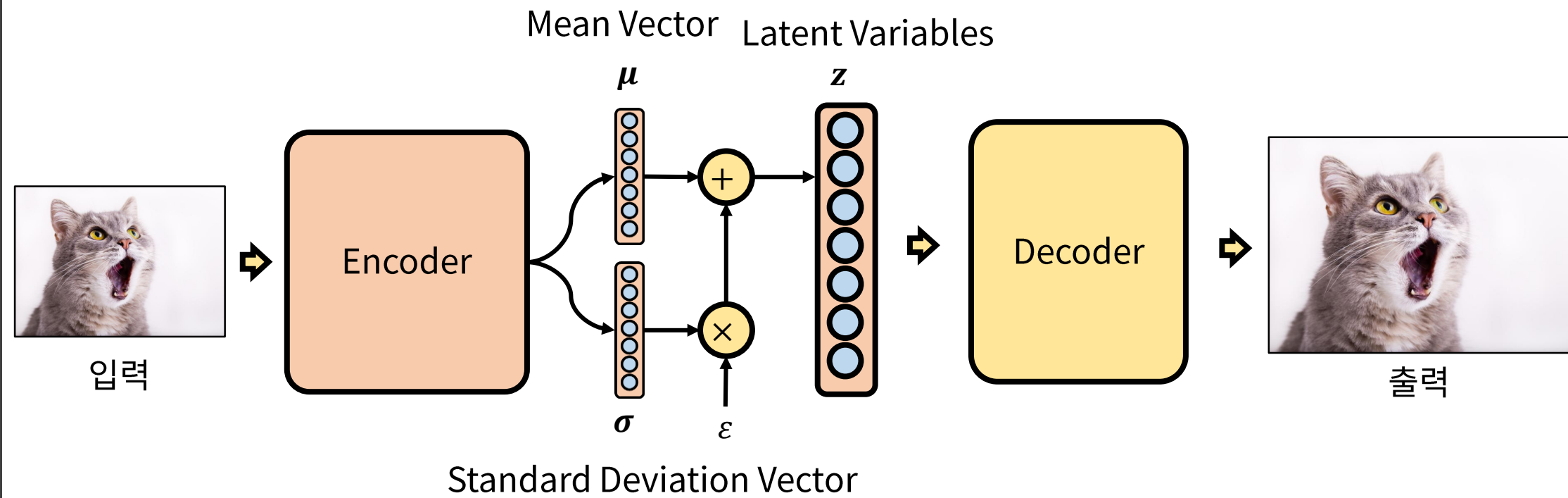


Chapter 06. 무엇이든 진짜처럼 생성하는 생성 모델(Generative Networks)

# Variational Autoencoder

# VAE Variational Autoencoder

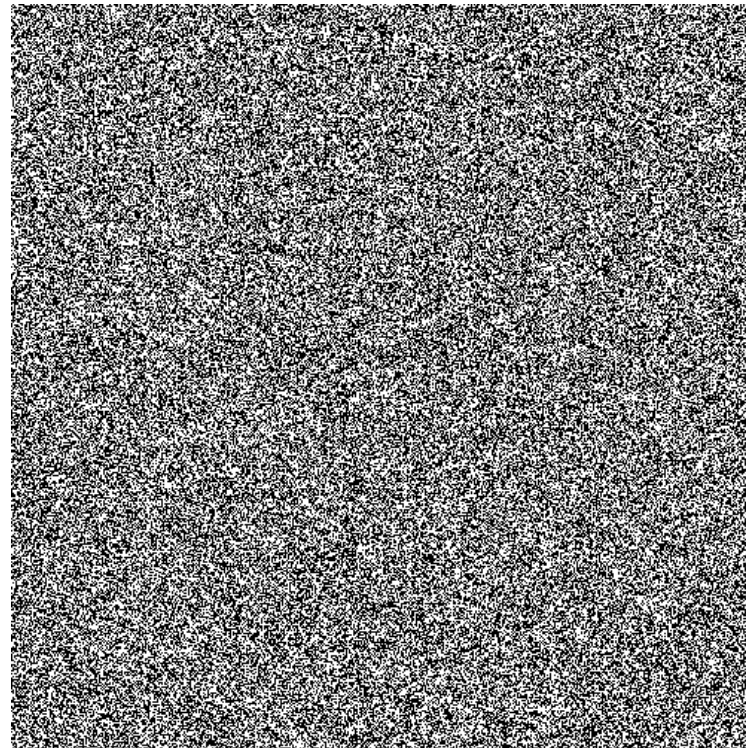


GAN에 대해서 먼저 알아보았기 때문에, VAE는 조금 더 쉽게 이해할 수 있을 것이다!

# Why VAE?



$$P(X = x_1)$$



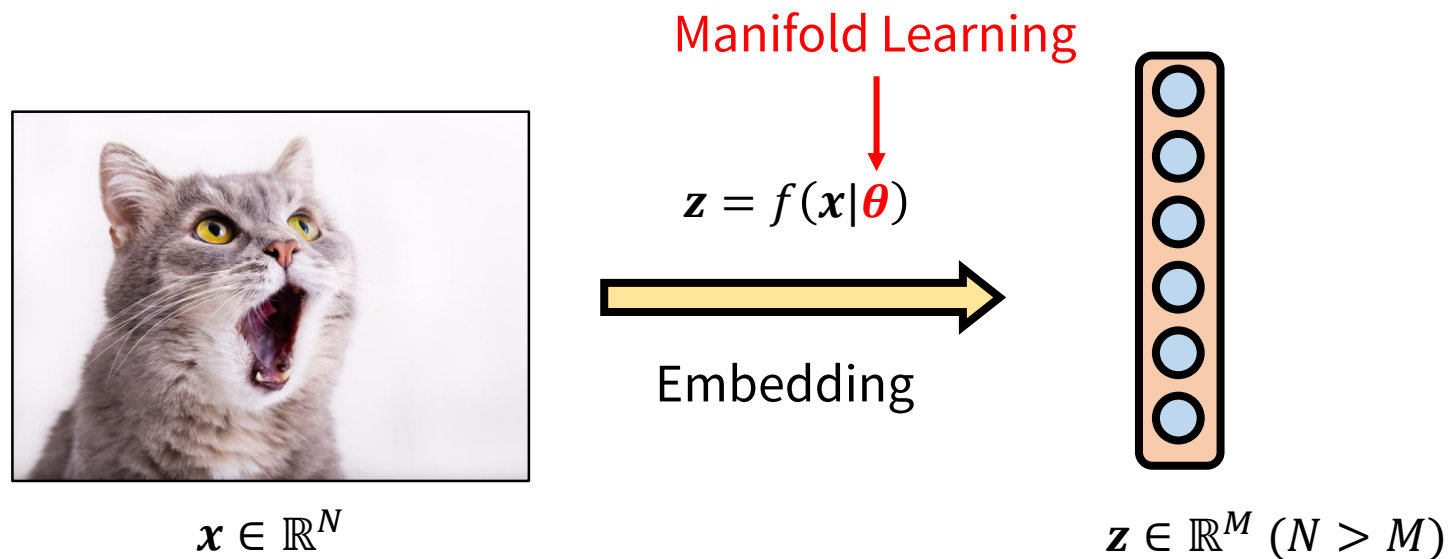
$$P(X = x_2)$$

$$P(X = x_1) = P(X = x_2)$$

임의로 영상을 생성할 경우, 좌측과 우측 영상이 발생할 확률은 동일하다.

Natural Image는 전체 영상 Domain 중에서 매우 Sparse하다.

# Manifold Learning



더 낮은 차원으로( $N > M$ ) 변환하는 것을 Embedding이라 하고,  
이 Embedding Function을 학습하는 것을 Manifold Learning이라 한다.

# Manifold Learning Example

Image Space

A selection from the 64-dimensional digits dataset



$$X \in \mathbb{R}^{8 \times 8}$$

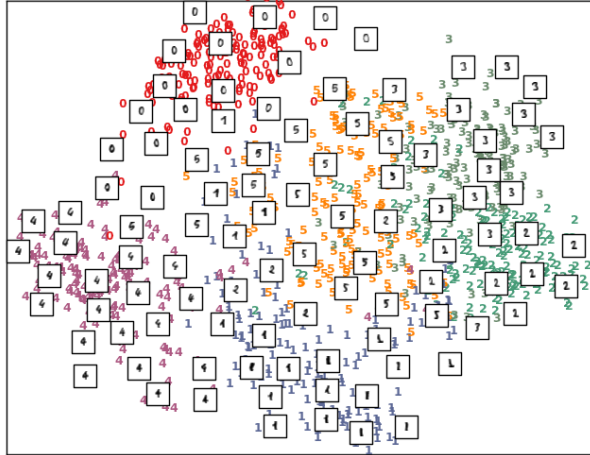
$$z = f(X|\theta)$$



Embedding

Latent Space

Principal Components projection of the digits (time 0.01s)

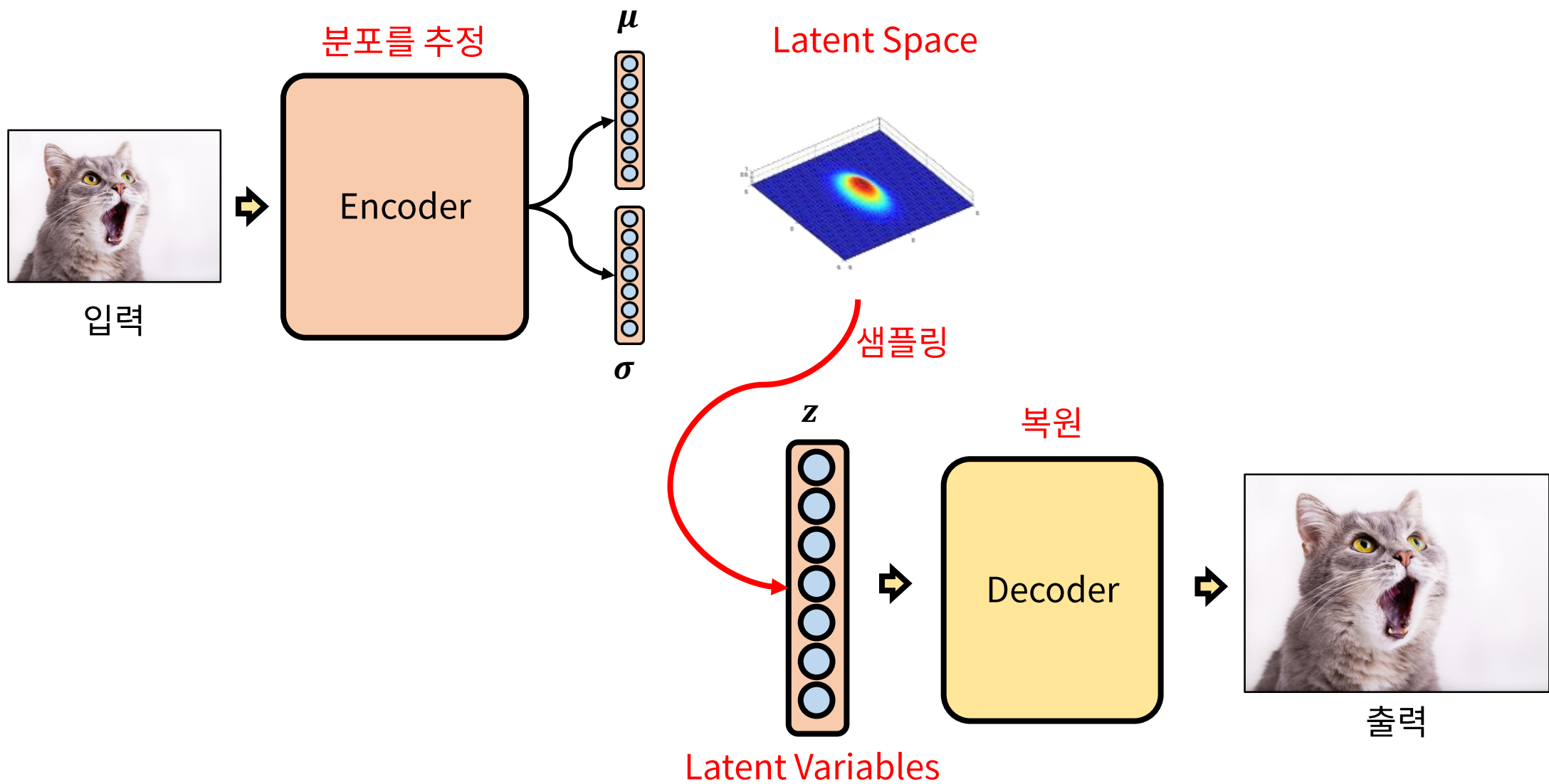


$$z \in \mathbb{R}^2$$

64차원 → 2차원의 Manifold Learning의 예 (PCA)

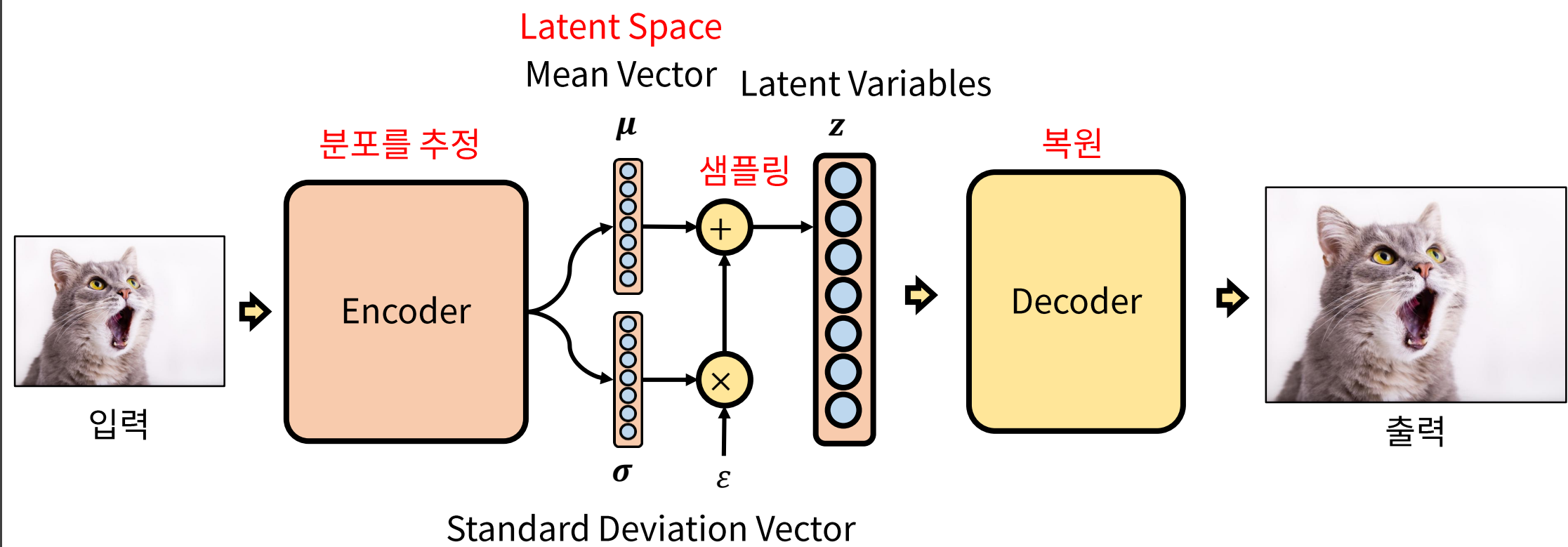
Image Space에서 생성하는 것 보다 Latent Space에서 생성하기가 훨씬 쉽다.

# VAE Structure (1/2)





# VAE Structure (2/2)



앞의 스토리를 보고 나니 좀 더 느낌이 오지 않는가?

# KL Divergence Kullback-Leibler Divergence

$$\begin{aligned} D_{KL}(p||q) &= \int p(x) \log \frac{p(x)}{q(x)} dx \\ &= \int p(x) \log p(x) dx - \int p(x) \log q(x) dx \\ &\neq D_{KL}(q||p) \end{aligned}$$

KL Divergence는 **두 확률 분포가 다른 정도**를 나타내는 것이 목적이다.  
'거리(Distance)'라고 부를 수는 없는데, 교환법칙이 성립하지 않기 때문이다.



# KLD of Gaussian Distribution

$$N(\mu, \sigma) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$

$$p(x) = N(\mu_1, \sigma_1)$$

$$q(x) = N(\mu_2, \sigma_2)$$

$$D_{KL}(p||q) = - \int p(x) \log q(x) dx + \int p(x) \log p(x) dx$$

$$= \frac{1}{2} \log(2\pi\sigma_1^2) + \frac{\sigma_1^2 + (\mu_1 - \mu_2)^2}{2\sigma_2^2} - \frac{1}{2} (1 + \log 2\pi\sigma_1^2)$$

$$= \log \frac{\sigma_2}{\sigma_1} + \frac{\sigma_1^2 + (\mu_1 - \mu_2)^2}{2\sigma_2^2} - \frac{1}{2}$$

두 Gaussian 분포의 비교는 KL Divergence로 쉽게 표현할 수 있다.

# Evidence Lower Bound (ELBO) (1/3)

$$\begin{aligned}\log P(x_i) &= \log \frac{P(x_i|z)p(z)}{P(z|x_i)} \\ &= \log P(x_i|z) + \log P(z) - \log P(z|x_i) - \log P(z|x_i)\end{aligned}$$

Log-Likelihood를 최대화 하고 싶다.  
(MLE)

$$\begin{aligned}\log P(x_i) &= \log P(x_i) \int q(z|x_i) dz \\ &= \int q(z|x_i) \log P(x_i) dz\end{aligned}$$

위 수식을 대입하고  
하나씩 전개해 보자.

# Evidence Lower Bound (ELBO) (2/3)

$$\begin{aligned}
 \log P(x_i) &= \int q(z|x_i) [\log P(x_i|z) + \log P(z) - \log P(z|x_i)] dz \\
 &= E_{q(z|x_i)} [\log P(x_i|z)] + \int q(z|x_i) \log p(z) dz - \int q(z|x_i) \log P(z|x_i) dz \pm \int q(z|x_i) \log q(z|x_i) dz \\
 &= E_{q(z|x_i)} [\log P(x_i|z)] - \int q(z|x_i) \log q(z|x_i) dz + \int q(z|x_i) \log p(z) dz \\
 &\quad + \int q(z|x_i) \log q(z|x_i) dz - \int q(z|x_i) \log P(z|x_i) dz \\
 &= E_{q(z|x_i)} [\log P(x_i|z)] - D_{KL}(q(z|x_i) || P(z)) + D_{KL}(q(z|x_i) || P(z|x_i)) \\
 &\geq E_{q(z|x_i)} [\log P(x_i|z)] - D_{KL}(q(z|x_i) || P(z))
 \end{aligned}$$

Decoder의 사후확률은 알기 어렵다.

# Evidence Lower Bound (ELBO) (3/3)

Loss Function

$$\log P(x_i) \geq \underbrace{E_{q(z|x_i)}[\log P(x_i|z)]}_{\text{Reconstruction Error}} - \underbrace{D_{KL}(q(z|x_i)||P(z))}_{\text{Regularization}}$$

Reconstruction Error

Regularization

둘 모두 Gaussian 분포를 따른다.

$$\underbrace{D_{KL}(q(z|x_i)||P(z))}_{\text{Regularization}} = D_{KL}\left(N(\mu_{q(x_i)}, \sigma_{q(x_i)}^2) || N(0,1)\right)$$

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

$$q(z|x_i) \sim N(\mu_{q(x_i)}, \sigma_{q(x_i)})$$

$$= \sum_i -\log \sigma_{q(x_i)} + \frac{1}{2} (\sigma_{q(x_i)}^2 + \mu_{q(x_i)}^2 - 1)$$

# Interesting Results

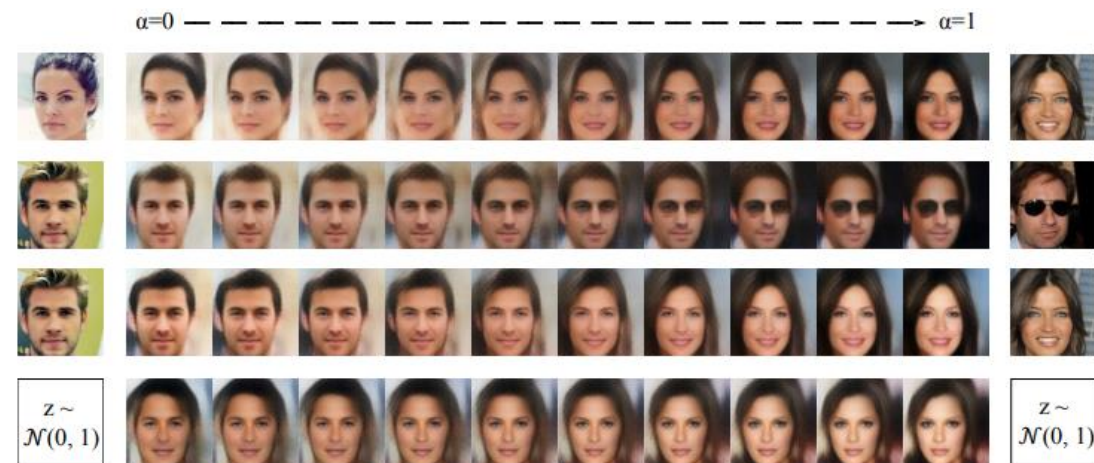


Figure 5. Linear interpolation for latent vector. Each row is the interpolation from left latent vector  $z_{left}$  to right latent vector  $z_{right}$ . e.g.  $(1 - \alpha)z_{left} + \alpha z_{right}$ . The first row is the transition from a non-smiling woman to a smiling woman, the second row is the transition from a man without eyeglass to a man with eyeglass, the third row is the transition from a man to a woman, and the last row is the transition between two fake faces decoded from  $z \sim \mathcal{N}(0, 1)$ .

# Interesting Results

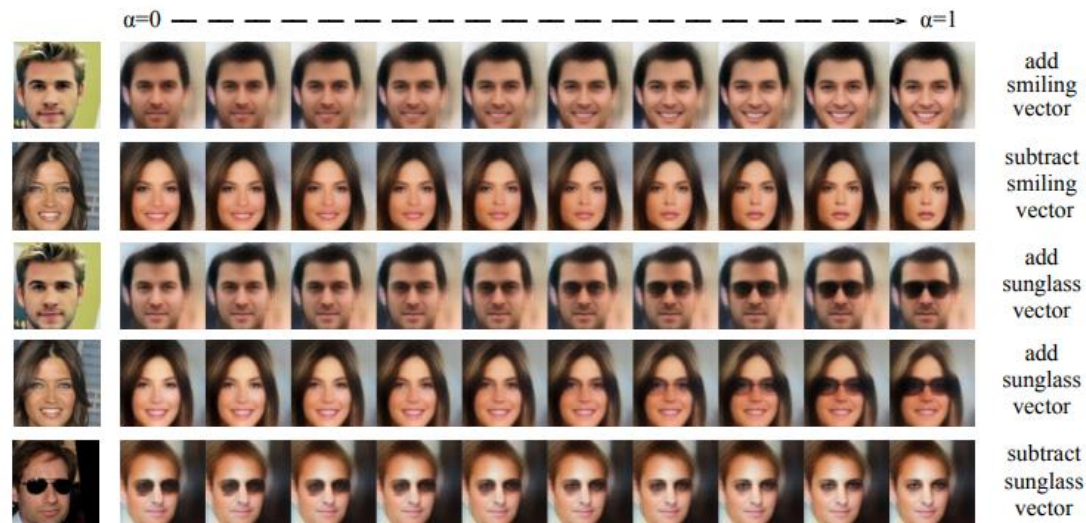


Figure 6. Vector arithmetic for visual attributes. Each row is the generated faces from latent vector  $z_{left}$  by adding or subtracting an attribute-specific vector, i.e.,  $z_{left} + \alpha z_{smiling}$ , where  $\alpha = 0, 0.1, \dots, 1$ . The first row is the transition by adding a smiling vector with a linear factor  $\alpha$  from left to right, the second row is the transition by subtracting a smiling vector, the third and fourth row are the results by adding a eyeglass vector to the latent representation for a man and women, and the last row shows results by subtracting an eyeglass vector.

<https://arxiv.org/pdf/1610.00291.pdf>