

Auto-Encoding Variational Bayes

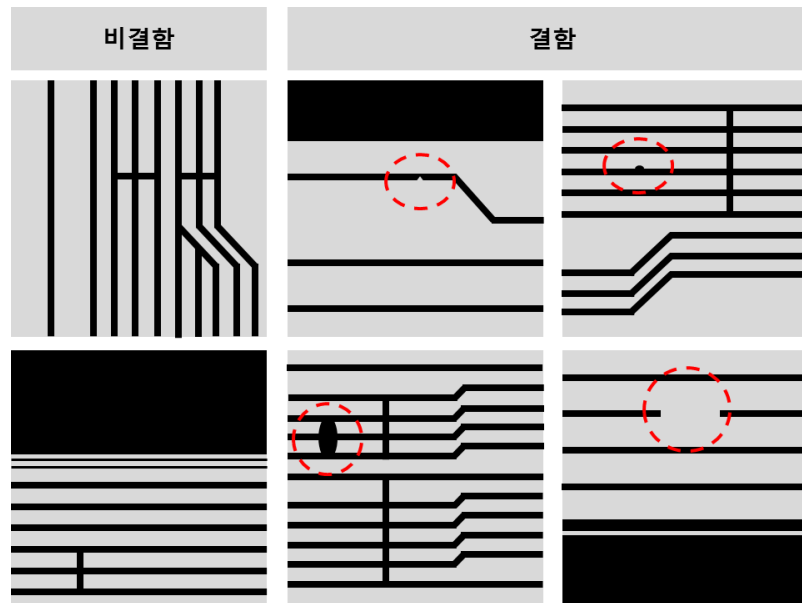
목차

- Generative Model
- Variational Inference
- Variational AutoEncoder (VAE)

Motivation

- 현업데이터(검사 도메인) controllable 한 EDA를 하기 위함
- 현업데이터 가능하다면 Generation
- 여러 논문들에서 중요한 개념이 되는 것을 공부하고자

현업데이터? Ambiguous? Model Capacity?



What is Deep Generative Model

Learning a Probability Distribution $P(X)$

- Generative Model의 궁극적 목적:
 - 데이터가 생성되는 과정, 즉 확률 분포 $P(X)$ 를 학습하고 싶은 것.
- 어떤 변수 X 의 확률 분포 $P(X)$ 를 안다는 것은 그 변수에 대한 모든 것을 안다는 의미다.
- 확률 분포 $P(X)$ 를 알고 안다면, 크게 2가지를 할 수 있다.
 1. (Inference) X 의 내부적 생성 구조, latent code z 혹은 class Y 정보에 관한 추론
 2. (Sampling) Data 생성 $X \sim P(X)$
- Generative Model을 학습하는 유일한 이유.

확률모형을 어떻게 학습 하지?

- 다양한 **Generative Model**

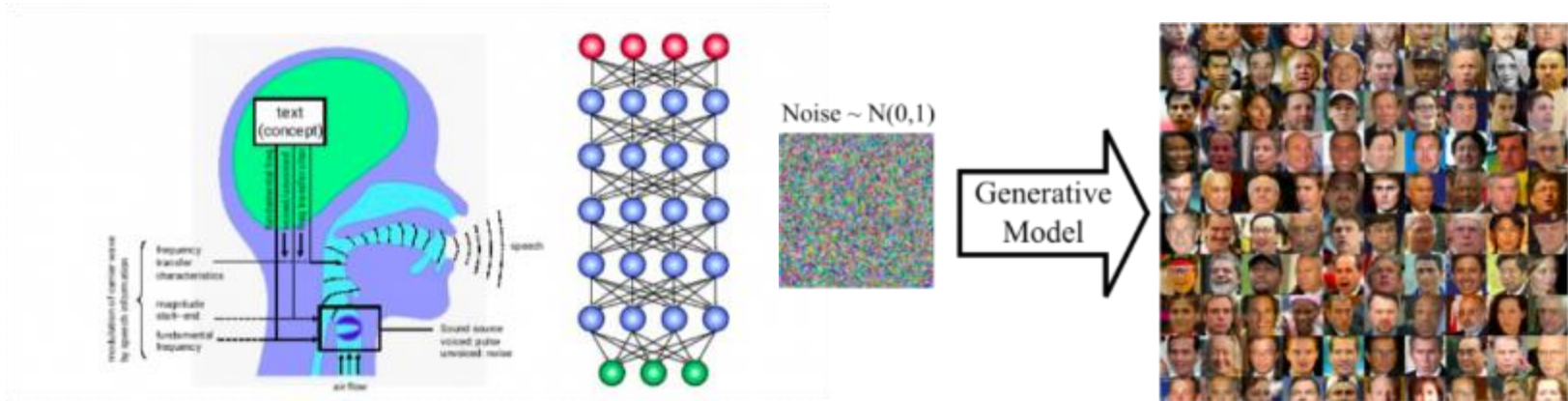
- Naïve Bayes, GDA, GMM, Beta-Bernoulli 모형...
- 이밖에 LDA (Topic Model), Hidden Markov Mode, non-parametric LDA, Markov Random Fields 등 다양한 전통적 Generative Model들이 있다.

- **Generative Model을 학습(parameter를추정)하는 몇가지 방법론 기법**

- MLE, MAP, EM, Analytical Posterior Inference, **(Stochastic) variational Inference w/ Re-parameterization Trick..**

- 대개 주어진데이터에 대한 확률분포 $P(X)$, 혹은 숨겨진구조 latent variable Z 에 대한 Prior $P(Z)$ 등을 가정하고, 주어진데이터를 잘 설명 하도록 모형을 학습했다.

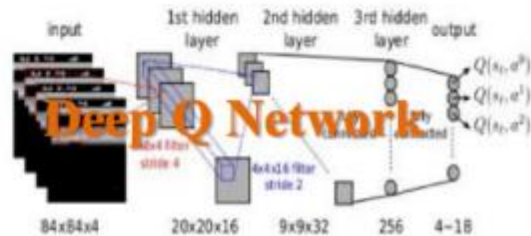
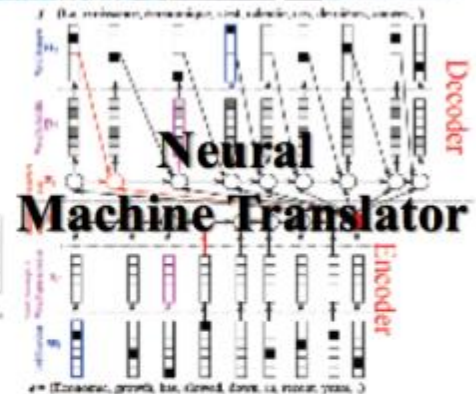
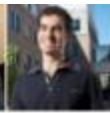
전통적인 Generative Model의 한계



- 기존인풋데이터에대한식별및추론에초점을두었던것과는달리, 딥러닝의시대에서는데이터의생성이중요한**task**다.
- 우리는인공지능이이미지,언어,소리등을만들어내길원한다.
- 하지만,현실세계의고차원적 데이터생성task를실현하기에는기존 확률모형들의표현력및학습방법등에분명한한계가있다.

Deep Learning & Supervised Learning

- 한편, Supervised Learning 영역에서는 인공지능과 같이 아주 복잡도가 높은 모형들에 대한 연구가 많이 이루어졌다.

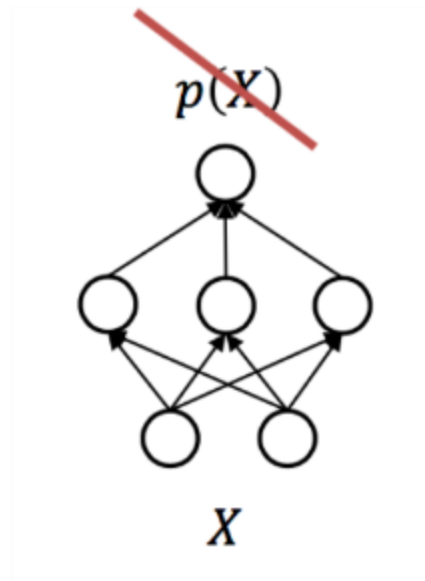


What is Deep Generative Model

“그렇다면, NN을 이용해서 확률모형을 학습할 수는 없을까?”

나이브한 접근

- 확률 X 의 확률밀도를 학습한다?



- 각 학습데이터에 대한 확률밀도값이 주어진다면 지도학습 (Supervised Learning)을 시도할 수 있다. 그러나, 일반적으로 데이터가 많아도 각 데이터의 정확한 확률밀도값을 알 수 없기 때문에 감독 학습만으로는 정확한 확률밀도값을 추정하기 힘들다.

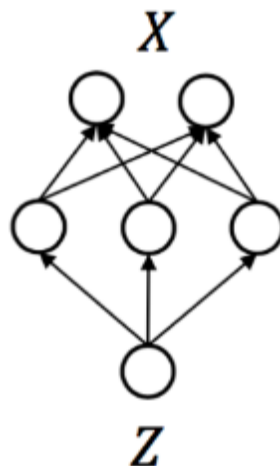
Deep Neural Network를 데이터 Sampling?

- 확률값 추론보다는 데이터를 잘 샘플링(생성)하는 모델을 만들어 보면 어떨까?
- 인공신경망은 한 벡터를 다른 벡터로 변환하는 것을 아주 잘 한다.

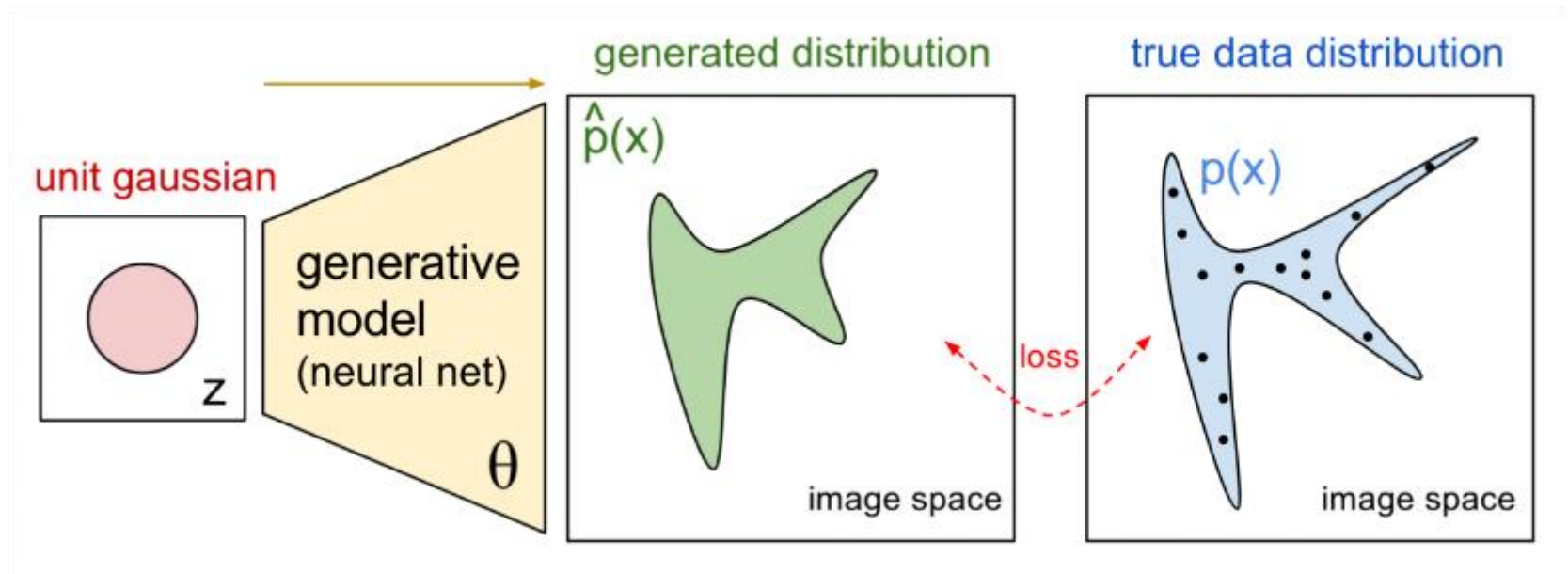
Key아이디어는

1. 간단한 확률 분포에서 얻어진 벡터를 $z \sim P(z)$ e.g. Gaussian
2. Neural Network를 통해서 실제와 비슷한 데이터 X 로 변환해보자!

$$Z \sim N(0, I)$$
$$X = f(Z)$$

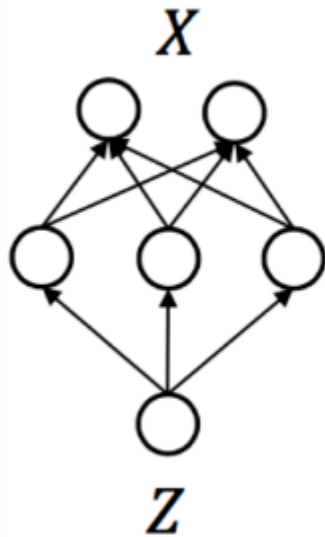


Deep Generative Model



Probability Distribution meets Deep Learning

- Random Variable의 함수를 적용하여 얻은 변수
== 새로운 Random Variable



$$Z \sim N(0, \mathbf{I})$$
$$p(X|Z) = f(Z)$$
~~$$X = f(Z)$$~~

- 엄밀하게 말하자면인공신경망을이용해서Conditional Probability $P(X|Z)$ 를모델링하자는것이다.

Latent Space? Z? Manifold?

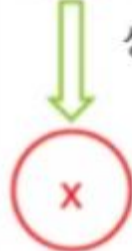
- 우리가 관측한 데이터 x 는 고차원이지만, 실제 의미가 표현되고 데이터의 생성에 관여하는 latent variable z 는 저차원으로 표현되고, 이 저차원상에 비슷한 데이터가 군집하는 공간이 있을 것이다.

잠재변수 z 의 예

: 물체의 형상, 카메라 좌표, 광원의 정보
(남자, [10, 2, -4], white)



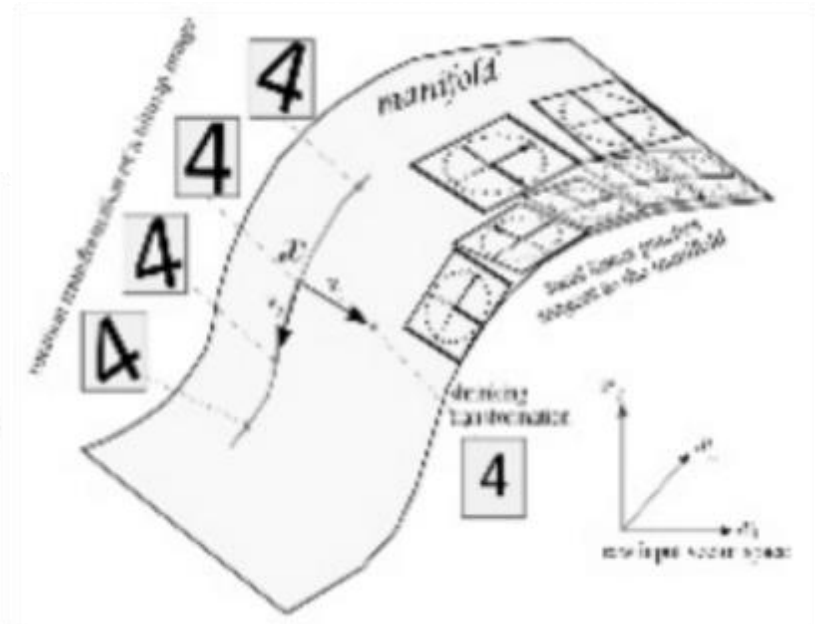
x : 이미지



생성



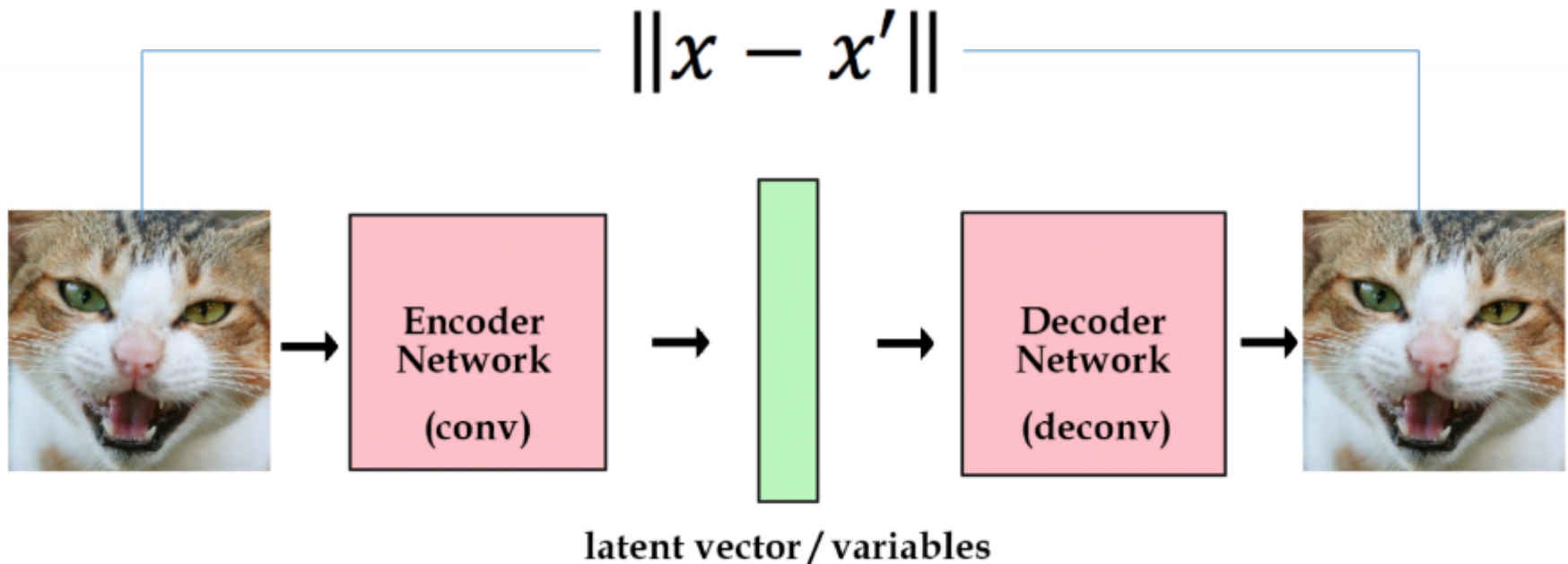
인식
(또는 추론)



Manifold hypothesis

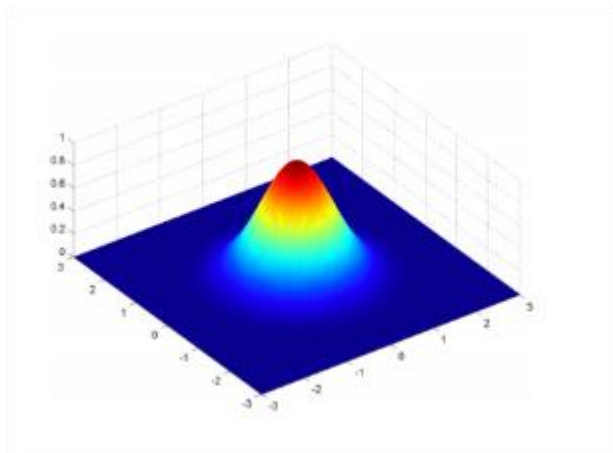
(vanilla) Auto-Encoder

- NN를 이용해서 각 영상에 대한 Latent Vector z 를 나타내는 방법 및 z 를 사용해서 영상을 복원하는 방법을 동시에 학습.
- 원본 데이터 x 자체가 일종의 레이블 역할을 하여 reconstructed된 데이터 또는 decoded 데이터와의 차이로 loss를 계산.
- L1 loss나 L2 loss를 많이 사용함(MSE)
- Unsupervised Learning, Feature Learning, Dimensionality Reduction

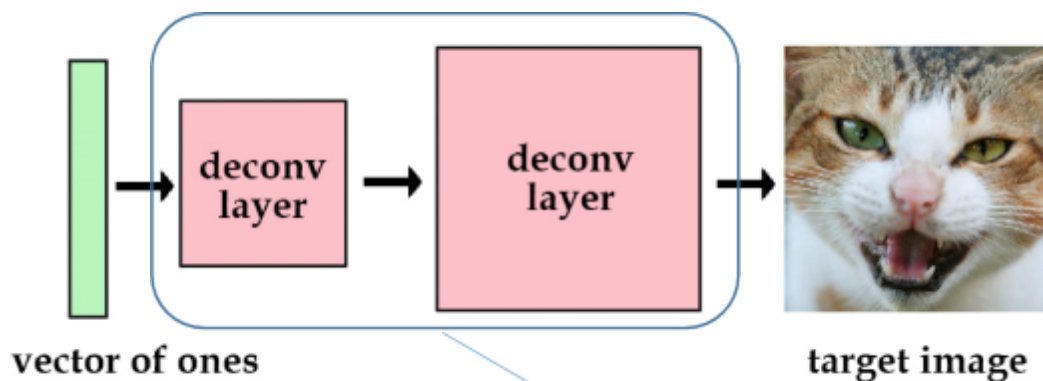


Variational Auto-Encoder의 목표

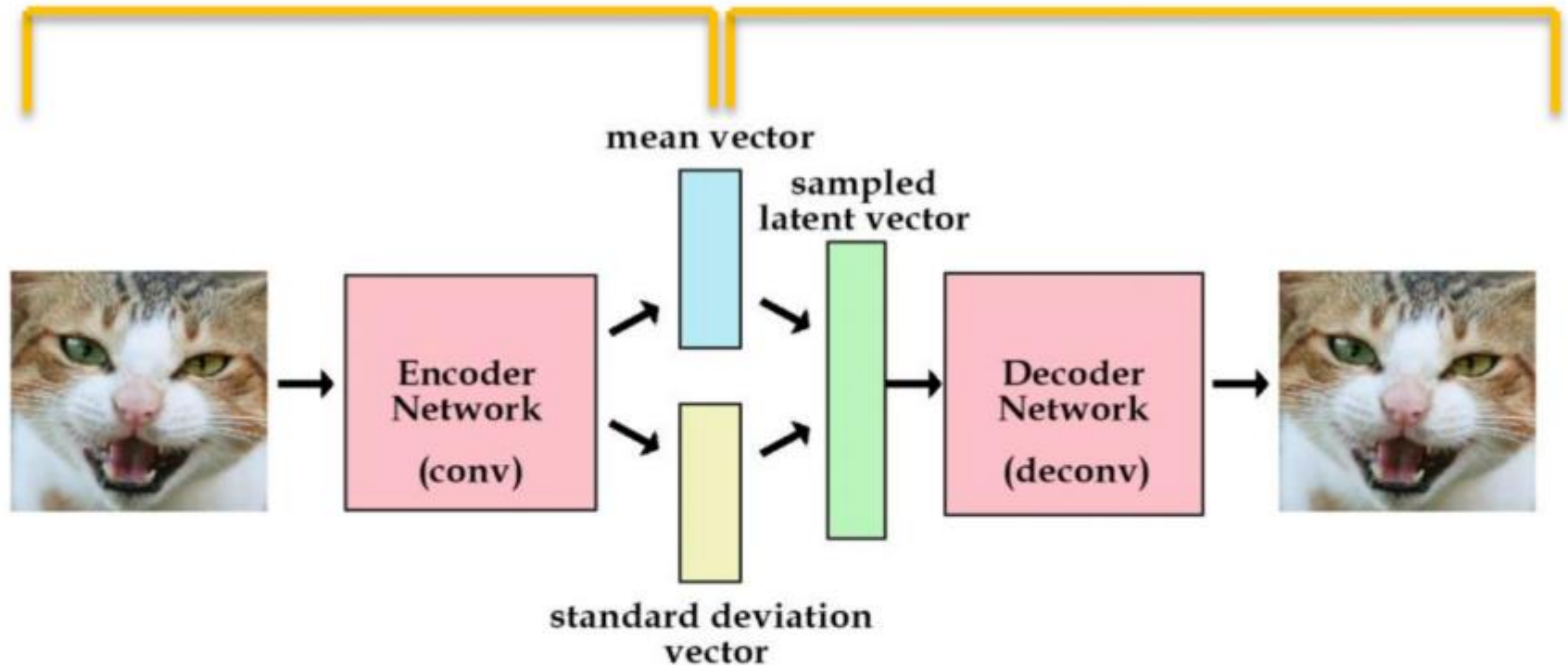
- (영상)데이터를 잘 생성하는 확률모형 $P(x|z;)$ 을 학습하는 것.
- $z \sim P(Z)$ 는 Gaussian 분포 사용.



$$P(z) = \mathcal{N}(z|0, I)$$



$$P(X|z; \theta) = \mathcal{N}(X|f(z; \theta), \hat{\sigma}^2 * I)$$

$q(z|x)$ $p(x|z)$ 

$$l_i(\theta, \phi) = \underbrace{-E_{z \sim q_\theta(z|x_i)}[\log p_\phi(x_i|z)]}_{\text{Reconstruction Loss}} + \underbrace{KL(q_\theta(z|x_i)||p(z))}_{\text{KL Divergence Regularizer}}$$

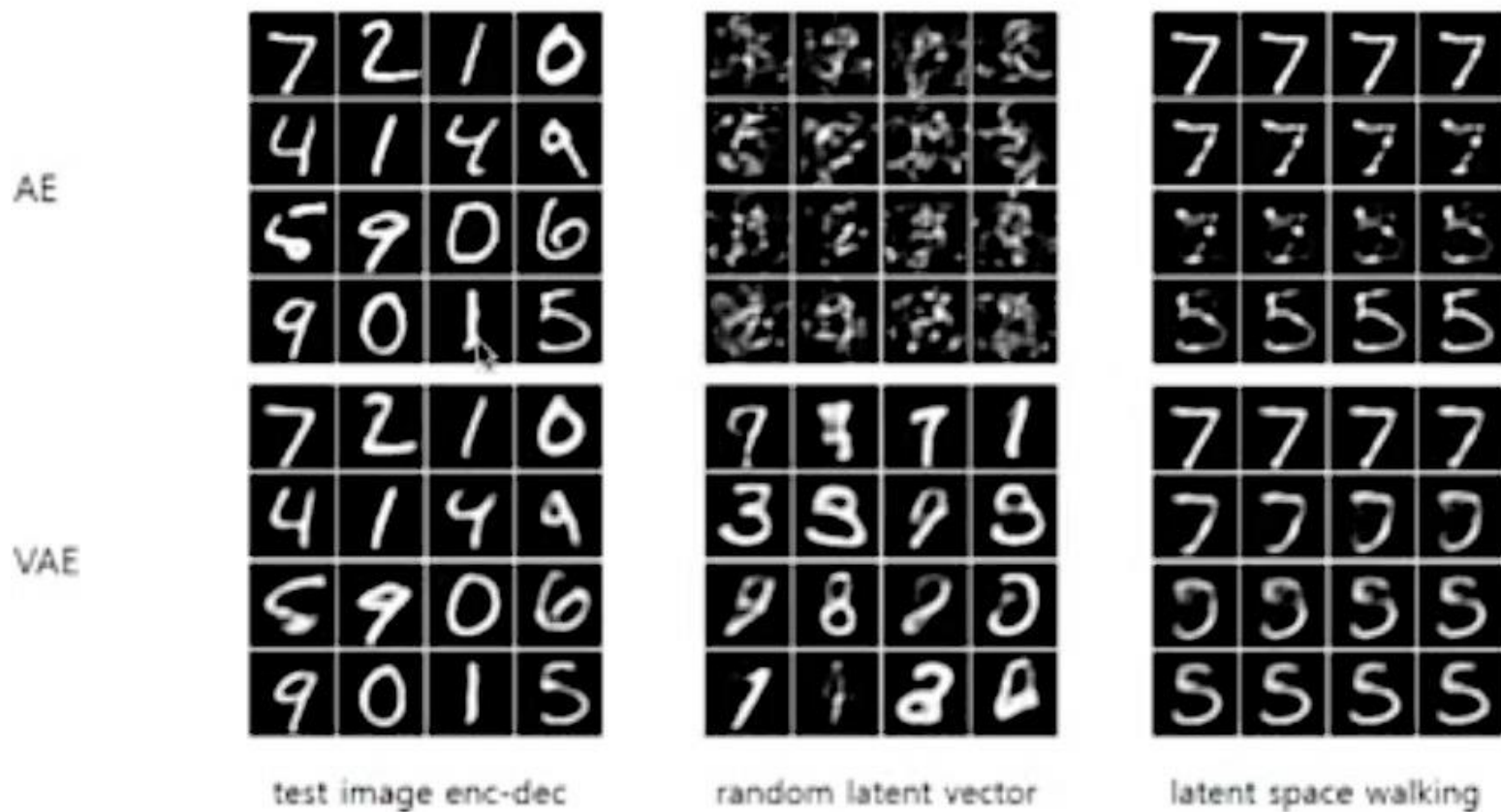
Reconstruction Loss
 $\|x - x'\|$

KL Divergence
Regularizer

AE vs VAE

- VAE는 AE와 다르게 latent vector z 에 대한 KL divergence **Regularization Term**이 추가적으로 붙는다.
- VAE는 AE와 다르게 **z 를 중간에 Sampling**하여 데이터 x 를 복원하는 방식으로 학습한다. 이러한 방식은 **Test-time**에 새로운 데이터 x 를 생성(샘플링)할 수 있게 한다(= x 에 대한 확률 분포를 학습한 것).
- 데이터 x 를 latent space z 의 차원으로 Mapping하는 방법, Encoder $q(z|x)$ 에 대한 확률 모형을 명시적으로 얻을 수 있다(= Posterior Distribution).

AE vs VAE 차이 비교



Background Keywords

Variational Inference (Bayesian Inference)

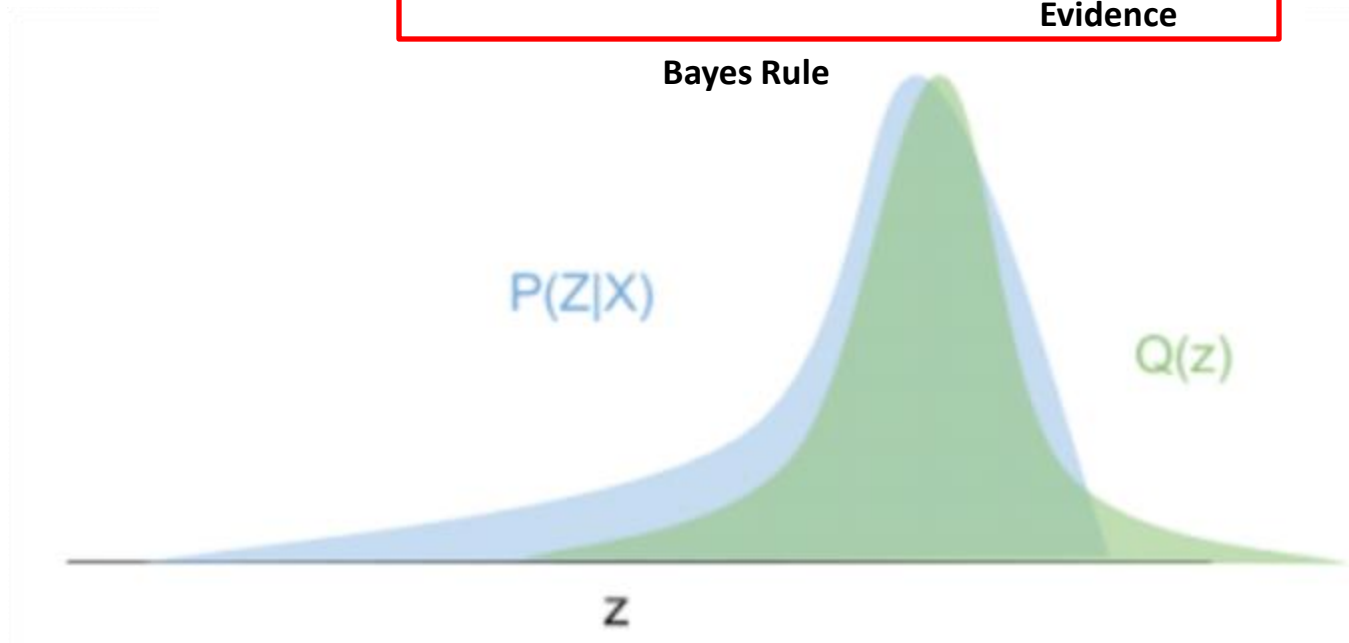
변분추론(Variational Inference)의 목적

- Posterior를 다루기 쉬운 분포 $Q(z)$ 로 근사하는 것임 (KLD 이용)

$$Q(z) \approx P(z|x) = \frac{P(x|z)P(z)}{P(x)} = \frac{P(x|z)P(z)}{\sum_z P(x, z)}$$

Evidence

Bayes Rule



V를 꺼놓으면 모델 어떻게 학습?

=> KL Divergence 관점

- KL divergence를 이용해서 $q(z)$ 와 $p(z|\mathbf{x})$ 와 오차를 측정!

$$D_{\text{KL}}(q(z) || p(z|\mathbf{x})) = \int q(z) \log \frac{q(z)}{p(z|\mathbf{x})} dz$$

gaussian

$$= \int q(z) \log \frac{q(z)p(\mathbf{x})}{p(\mathbf{x}|z)p(z)} dz$$

$$= \int q(z) \log \frac{q(z)}{p(z)} dz + \int q(z) \log p(\mathbf{x}) dz - \int q(z) \log p(\mathbf{x}|z) dz.$$

$$= \underline{D_{\text{KL}}(q(z) || p(z))} + \log p(\mathbf{x}) - \mathbb{E}_{z \sim q(z)} [\log p(\mathbf{x}|z)].$$

- KL divergence term이 줄어듦에 따라 $q(z)$ 를 학습하면 됨

- 결론: $p(z|\mathbf{x})$ 를 잘 근사한 $q^*(z)$ 를 얻음

원래 알고 싶었던 이놈(Posterior)
Latent effective knowledge set.

posterior

Evidence Lower Bound (ELBO) 관점

$$D_{\text{KL}}(q(z) \parallel p(z | \mathbf{x})) = D_{\text{KL}}(q(z) \parallel p(z)) + \log p(\mathbf{x}) - \mathbb{E}_{z \sim q(z)} [\log p(\mathbf{x} | z)]$$



$$\log p(\mathbf{x}) = \mathbb{E}_{z \sim q(z)} [\log p(\mathbf{x} | z)] - D_{\text{KL}}(q(z) \parallel p(z)) + \underbrace{D_{\text{KL}}(q(z) \parallel p(z | \mathbf{x}))}_{\geq 0}$$

$$\log p(\mathbf{x}) \geq \mathbb{E}_{z \sim q(z)} [\log p(\mathbf{x} | z)] - D_{\text{KL}}(q(z) \parallel p(z)) = \mathcal{L}(\theta_q)$$

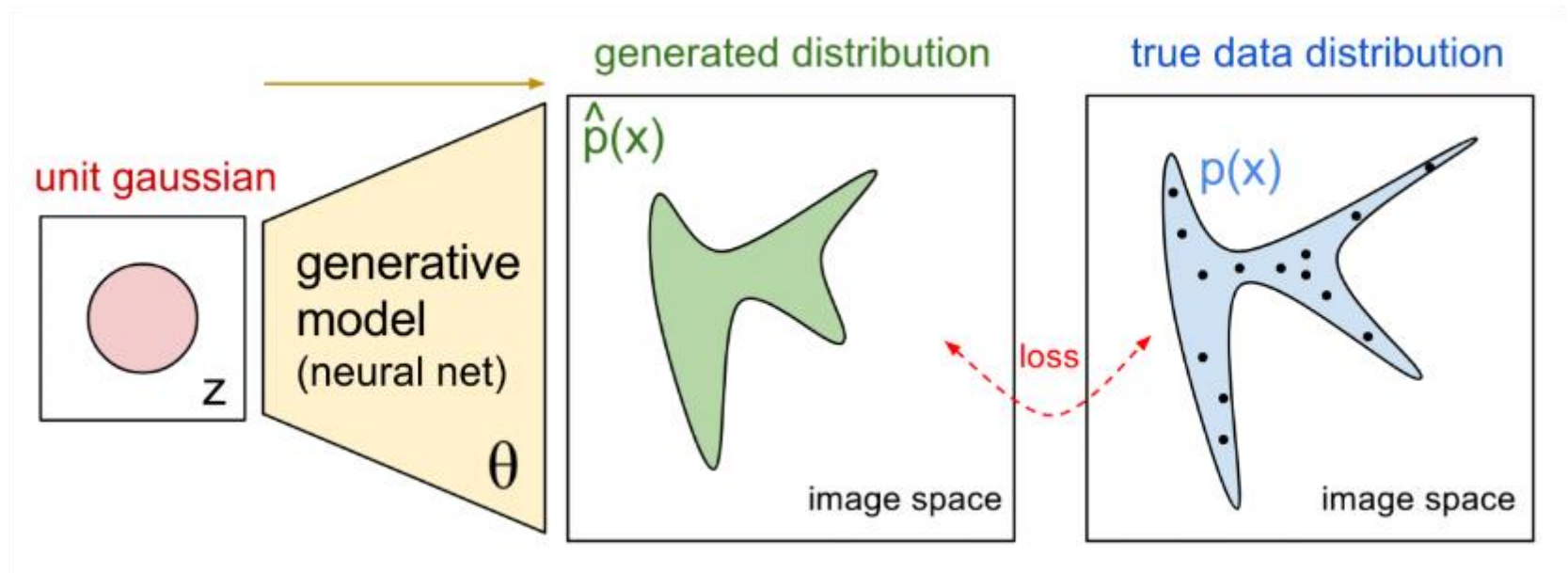
- $\mathcal{L}(\theta_q)$ = Evidence Lower Bound (ELBO)라고 부름
- $p(x|z)$ 에 추가적인 모델Parameter θ 가 붙으면 우리가 익숙한 marginal log-likelihood의 lower bound로 볼 수 있다.

$$\begin{aligned} \log p(x; \theta) &\geq \mathbb{E}_{z \sim q(z; \theta_q)} [\log p(x|z; \theta)] - D_{\text{KL}}(q(z; \theta_q) \parallel p(z)) \\ &= \mathcal{L}(\theta_q, \theta) = \sum_z q(z; \theta_q) \left[\log \frac{p(x|z; \theta)p(z)}{q(z; \theta_q)} \right] \end{aligned}$$

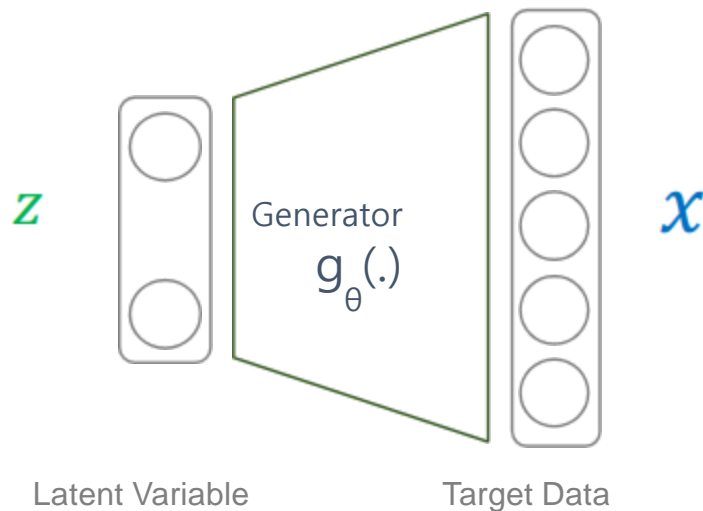
VAE

Variational Auto-Encoder

Deep Generative Model



Deep Generative Model



- Latent Variable z is target data x 의 특징을 결정짓는다.

- z 는또다른 latent code c 와 concatenation이될수있고, z 자체를 code로생각하여 $[0, \dots, 9]$ 을인풋으로 사용하는것도가능

$$p(x|g_{\theta}(z)) = p_{\theta}(x|z)$$

$z \sim p(z)$ Random variable

$g_{\theta}(\cdot)$ Deterministic function parameterized by θ

$x = g_{\theta}(z)$ Random variable

- Generative Model의 학습방법은 data x 에 대한 likelihood를 최대화하는 것이다. **MLE**

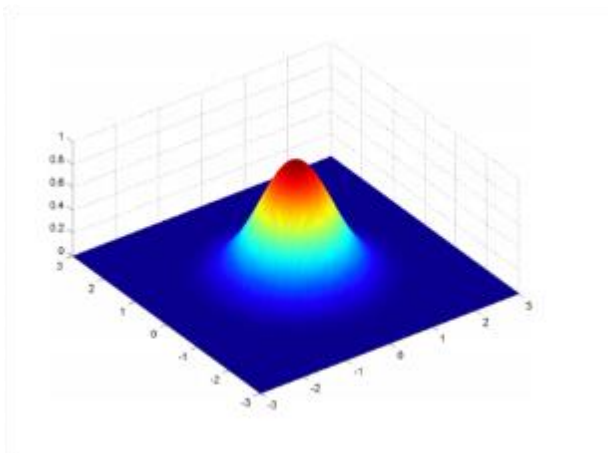
$$\int p(x|g_{\theta}(z))p(z)dz = p(x)$$

Probability density

VAE 목표

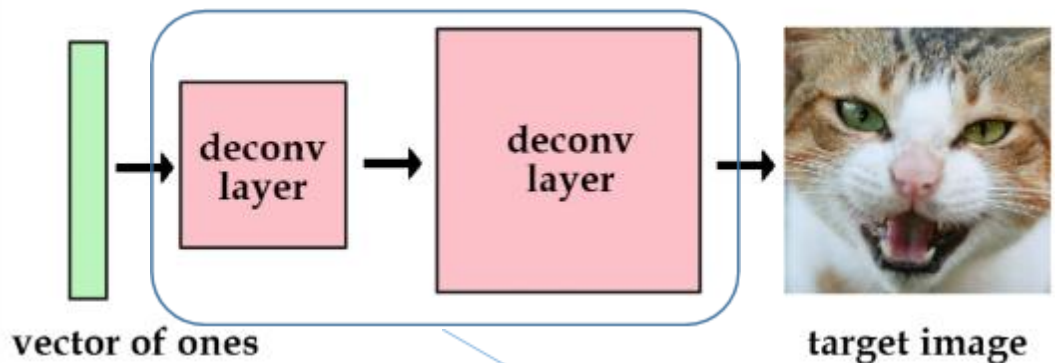
디코더

(영상)데이터를 잘 생성하는 확률 모형 $P(X|Z; \theta)$ 학습하는 것임
 $z \sim P(Z)$ 는 Gaussian 분포(N)사용.

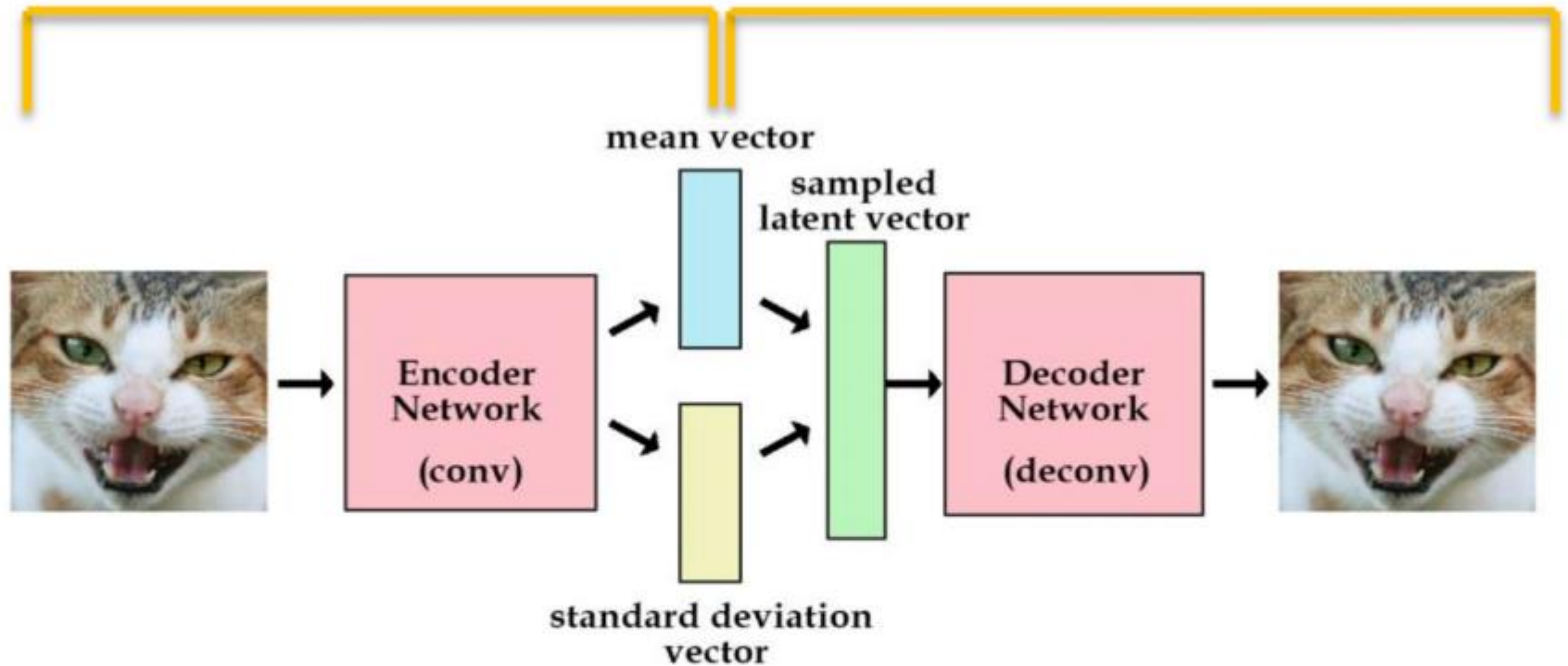


$$P(z) = \mathcal{N}(z|0, I)$$

가정: prior은 가우시안



$$P(X|z; \theta) = \mathcal{N}(X|f(z; \theta), \sigma^2 * I)$$

$q(z|x)$ $p(x|z)$ 

$$l_i(\theta, \phi) = \underbrace{-E_{z \sim q_\theta(z|x_i)}[\log p_\phi(x_i|z)]}_{\text{Reconstruction Loss}} + \underbrace{KL(q_\theta(z|x_i)||p(z))}_{\text{KL Divergence Regularizer}}$$

Reconstruction Loss
 $\|x - x'\|$

KL Divergence
Regularizer

모형(VAE)을 어떻게 학습 시킴?

- **MLE**을 통해서 Generative Model의 Parameter를 추정하려면 Marginal-log-likelihood $\log P(X; \theta)$ 를 maximize하면 되겠다!?

$$\log p(x; \theta) = \log \sum_z p(x|z; f_{\theta}(z), \sigma^2 * I) p(z) \quad \text{P}(X|Z; \theta)$$

- 보통 Generative Model (e.g. GMM)에서도 marginal log-Likelihood를 직접 최대화하기 어려움 $\pi\pi$ 계산불가, No Analytic
- VAE는 최적화 대상인 parameter (θ 가 확률모형 $p(x|z)$ 속 $f_{\theta}(z)$, 즉, neural network의 parameter θ 임으로 analytical 하게 최적의 Parameter를 추정하는 것은 어려움 $\pi\pi$ (문제발생)

⇒ 그럼

ELBO(VI=베이지안 접근)로 추정해볼까?

$$P(\theta|x) = \frac{\overset{\text{Likelihood}}{P(x|\theta)} \overset{\text{Prior}}{P(\theta)}}{\int d\theta P(x|\theta) P(\theta)}$$

사후 정보 = 데이터로 얻어진 정보 + 사전 정보

ELBO and Variational Inference

VI를 통해서 임의의 z 분포 $q(z; \theta_q)$ 를 만들고,
log-likelihood의 **Lower Bound**를 정할 수 있음.

$$\begin{aligned} \log p(x; \theta) &= \log \sum_z \left[\overset{\text{q 깁겨넣기}}{q(z; \theta_q)} \frac{p(x|z; \theta)p(z)}{q(z; \theta_q)} \right] \\ &\geq \sum_z q(z; \theta_q) \left[\log \frac{p(x|z; \theta)p(z)}{q(z; \theta_q)} \right] \\ &\text{By Jensen's inequality} \\ &= \mathbb{E}_{z \sim q(z; \theta_q)} [\log p(x|z; \theta)] - D_{KL}(q(z; \theta_q) || p(z)) = \mathcal{L}(\theta_q, \theta) \end{aligned}$$

목적: log-likelihood
Theta 학습

ELBO

• 계산이 쉬운 $\mathcal{L}(\theta_q, \theta)$ 를 최대화하면!

원래 알고싶었던 이놈(Posterior)
Latent effective knowledge set.

KL divergence term이 줄어들때 까지

$q(z; \theta_q)$ 는 알고싶었던 posterior $P(z|x)$ 를 닮아가며 학습 됨 $q^*(z)$ 를 얻음
 $\log p(x; \theta)$ 는 최대화되도록 학습이 가능하다. ■ 결론: $p(z|x)$ 를 잘 근사한 posterior

VAE에서 일반적인 VI와 달리 **Amortize trick** 사용하는 이유

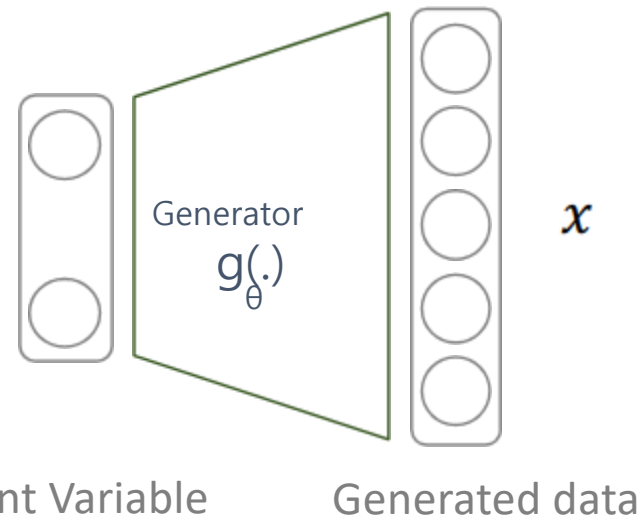
- 기존 Variational Inference에는 $q(z; \theta_q)$ 를 Gaussian으로 정했다.
- 그러나 x 의 차원이 클때는, $q(z; \theta_q)$ 가 Posterior $P(z|x)$ 에 수렴하는 속도가 아주 느리다.
- 왜? $q(z; \theta_q)$ 가 Gaussian 인것은, 너무 단순해서 z 와 x 의 복잡한 관계인 $P(z|x)$ 를 근사하는 것이 어렵다.
- $q(z; \theta_q)$ 가 Posterior에 수렴하지 않는다면, Lower Bound가 tight해 지지 않는다. $\rightarrow \log P(x; \theta)$ 의 최대화도 어렵다.

Amortize trick

- $q(z; \theta_q)$ 가 너무 단순할 때 생기는 문제점을 해결하기 위해서 Amortized Variational Inference 사용한다.
- 말은 어려워 보이지만.. $q(z; \theta_q)$ 의 Parameter θ_q 를 x 에 대한 함수로 만드는 것이다.
- 즉, $q(z) = N(\mu_q, \sigma_q)$ 을 $q(z|x) = N(\mu(x), \sigma(x))$ 로 바꾼다.
- 이때, 복잡한 x 와 z 의 관계를 고려해서, $\mu(x), \sigma(x)$ 를 또 다른 Neural Network로 정하였다.

VAE의 최종 목적 함수

$$p(z|x) \approx q_\phi(z|x) \sim z$$



Optimization Problem 1 on ϕ : Variational Inference

$$\log(p(x)) \geq \mathbb{E}_{q_\phi(z|x)}[\log(p(x|z))] - KL(q_\phi(z|x)||p(z)) = ELBO(\phi)$$

Optimization Problem 2 on θ : Maximum likelihood

$$-\sum_i \log(p(x_i)) \leq -\sum_i \mathbb{E}_{q_\phi(z|x_i)}[\log(p(x_i|g_\theta(z)))] - KL(q_\phi(z|x_i)||p(z))$$

Final Optimization Problem

$$\arg \min_{\phi, \theta} \sum_i -\mathbb{E}_{q_\phi(z|x_i)}[\log(p(x_i|g_\theta(z)))] + KL(q_\phi(z|x_i)||p(z))$$

ELBO: Evidence Lower Bound in VAE

$$\log(p(x)) = \log\left(\int p(x|z)p(z)dz\right) \geq \int \log(p(x|z))p(z)dz \quad \leftarrow \text{Jensen's Inequality}$$

[Jensen's Inequality]
For concave functions $f(\cdot)$
 $f(E[x]) \geq E[f(x)]$

$$\log(p(x)) = \log\left(\int p(x|z) \frac{p(z)}{q_\phi(z|x)} q_\phi(z|x) dz\right) \geq \int \log\left(p(x|z) \frac{p(z)}{q_\phi(z|x)}\right) q_\phi(z|x) dz$$

$$\log(p(x)) \geq \int \log(p(x|z)) q_\phi(z|x) dz - \int \log\left(\frac{q_\phi(z|x)}{p(z)}\right) q_\phi(z|x) dz$$

$$= \underbrace{\mathbb{E}_{q_\phi(z|x)}[\log(p(x|z))]}_{ELBO(\phi)} - KL(q_\phi(z|x) || p(z))$$

Variational lower bound
Evidence lower bound (ELBO)

**x에 대하여 잘 모사
& z에 대하여 잘 모사**

ELBO를 최대화하는 ϕ 값을 찾는다는 것은 $q_\phi(z|x)$ 와 $p(z|x)$ 간의 KL divergence를 줄이는 것과 같은 의미

ELBO: Evidence Lower Bound in VAE

$$\begin{aligned}
 \log(p(x)) &= \int \log(p(x)) q_{\phi}(z|x) dz \quad \leftarrow \int q_{\phi}(z|x) dz = 1 \\
 &= \int \log\left(\frac{p(x, z)}{p(z|x)}\right) q_{\phi}(z|x) dz \quad \leftarrow p(x) = \frac{p(x, z)}{p(z|x)} \\
 &= \int \log\left(\frac{p(x, z)}{q_{\phi}(z|x)} \cdot \frac{q_{\phi}(z|x)}{p(z|x)}\right) q_{\phi}(z|x) dz \\
 &= \underbrace{\int \log\left(\frac{p(x, z)}{q_{\phi}(z|x)}\right) q_{\phi}(z|x) dz}_{ELBO(\phi)} + \underbrace{\int \log\left(\frac{q_{\phi}(z|x)}{p(z|x)}\right) q_{\phi}(z|x) dz}_{KL(q_{\phi}(z|x) \parallel p(z|x))}
 \end{aligned}$$

두 확률분포 간의 거리 ≥ 0

ELBO를 최대화하는 ϕ 값을 찾는다는 것은 $q_{\phi}(z|x)$ 와 $p(z|x)$ 간의 KL divergence를 줄이는 것과 같은 의미

Variational Inference Objective and VAE

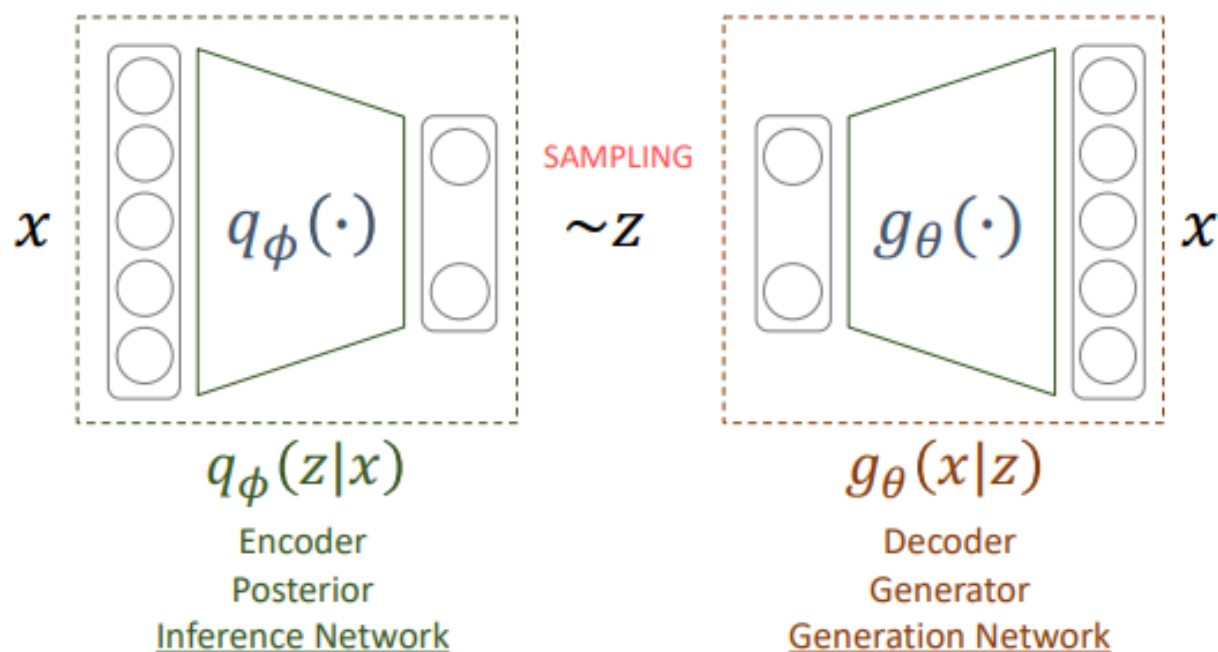
$$\log p(x) - D_{KL}(q_{\theta}(z|x)||p(z|x)) = \underbrace{\mathbb{E}_{z \sim q_{\theta}(z|x)} [\log p_{\phi}(x|z)]}_{\text{Reconstruction term}} - \underbrace{D_{KL}(q_{\theta}(z|x)||p(z))}_{\text{Regularizing Term}}$$

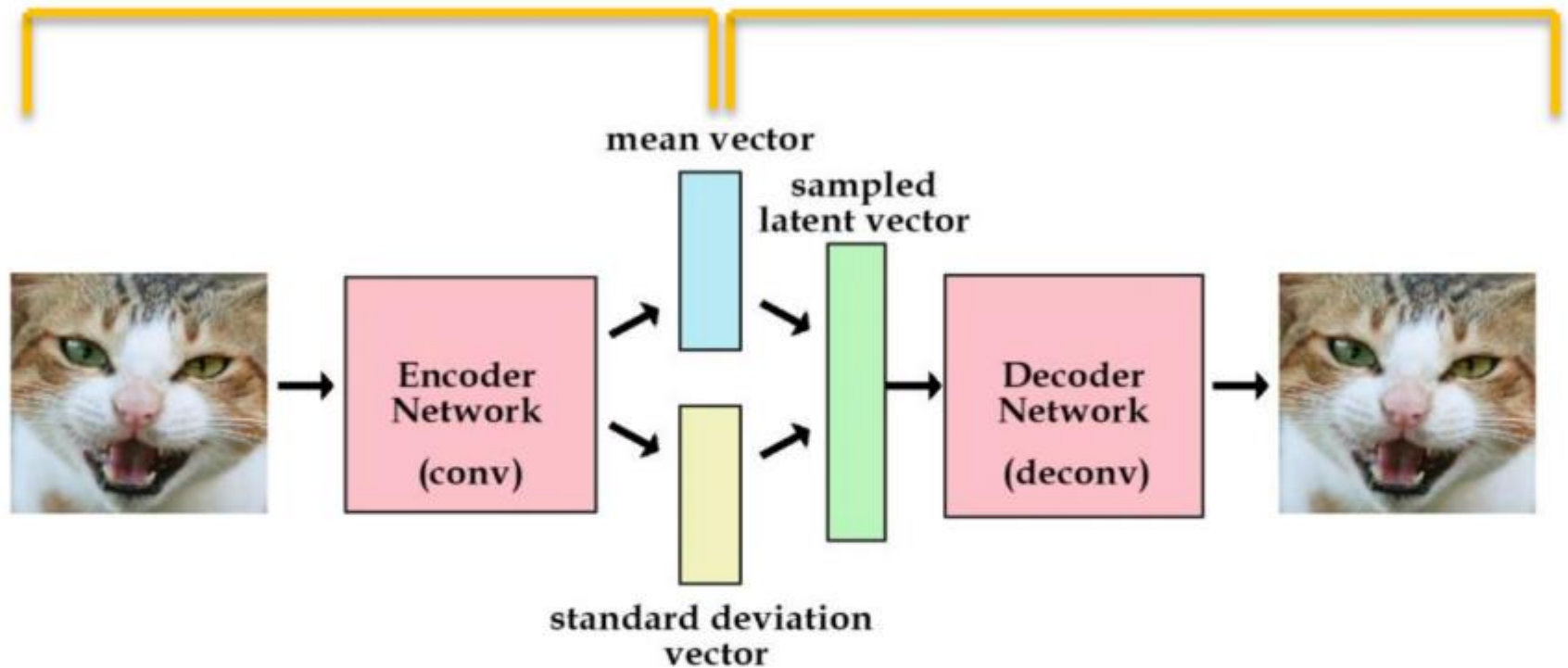
$$\log p(x) \geq \mathbb{E}_{z \sim q_{\theta}(z|x)} [\log p_{\phi}(x|z)] - D_{KL}(q_{\theta}(z|x)||p(z))$$

x에 대하여 잘 모사
& z에 대하여 잘 모사

Neural Network의 관점

$$\arg \min_{\phi, \theta} \sum_i -\mathbb{E}_{q_{\phi}(z|x_i)} [\log(p(x_i|g_{\theta}(z)))] + KL(q_{\phi}(z|x_i)||p(z))$$



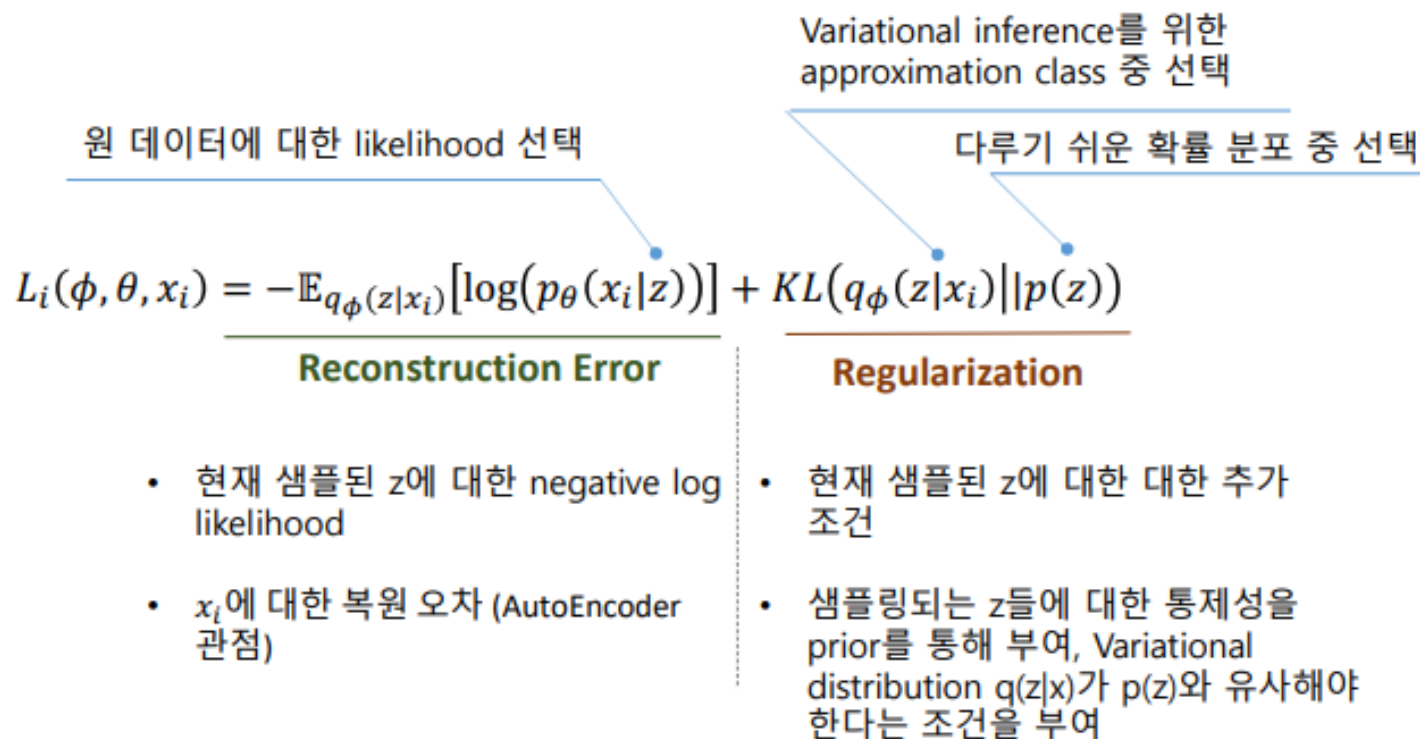
$q(z|x)$ $p(x|z)$ 

$$l_i(\theta, \phi) = \underbrace{-E_{z \sim q_\theta(z|x_i)}[\log p_\phi(x_i|z)]}_{\text{Reconstruction Loss}} + \underbrace{KL(q_\theta(z|x_i)||p(z))}_{\text{KL Divergence Regularizer}}$$

Reconstruction Loss
 $\|x - x'\|$

KL Divergence
Regularizer

VAE Loss 해석 (1)



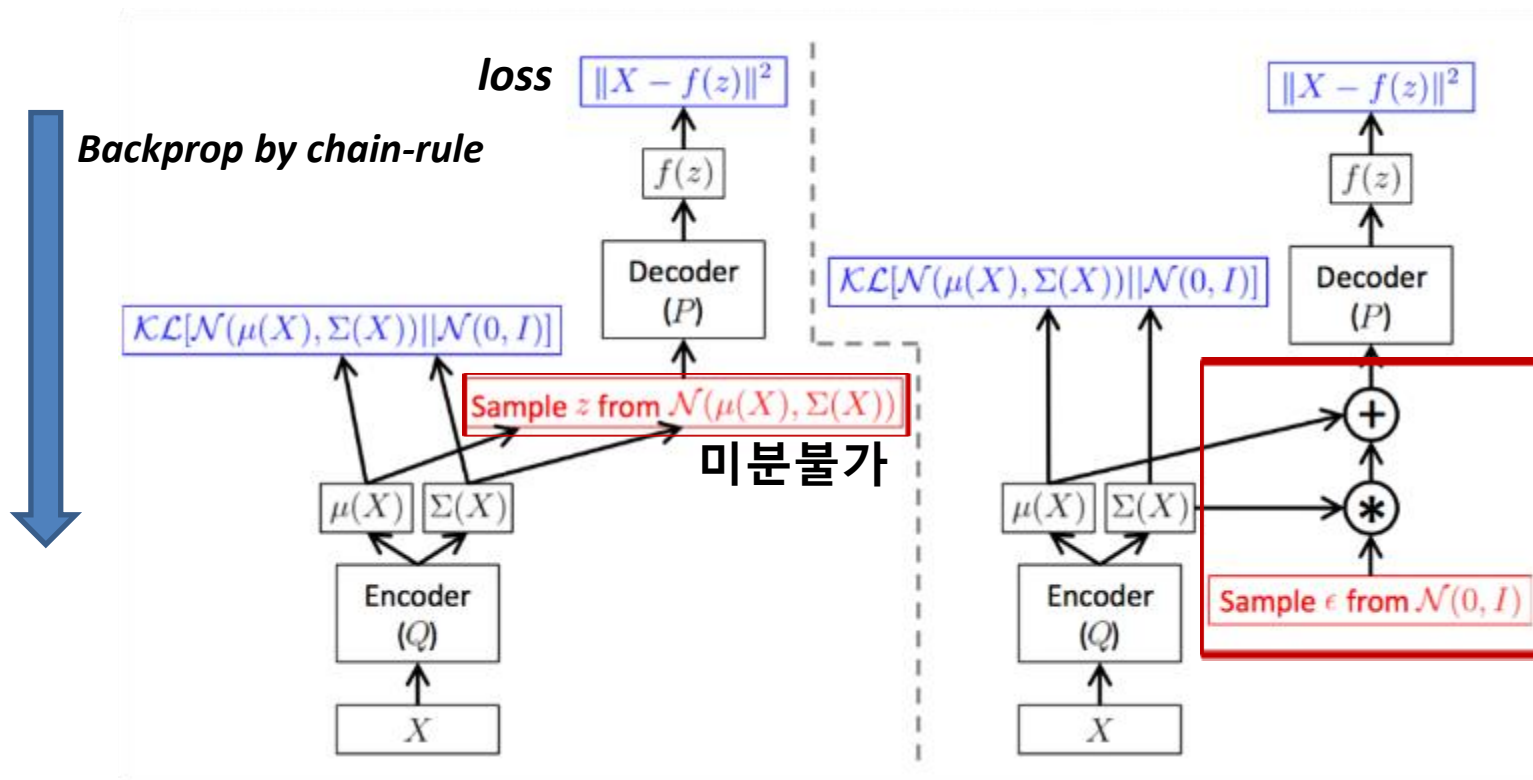
VAE Loss 해석 (2)

$$\begin{aligned} L_i(\phi, \theta, x_i) &= -\mathbb{E}_{q_\phi(z|x_i)}[\log(p_\theta(x_i|z))] + KL(q_\phi(z|x_i)||p(z)) \\ &= \underbrace{-\mathbb{E}_{q_\phi(z|x_i)}[\log(p_\theta(x_i|z))]}_{\text{Reconstruction Error}} - \underbrace{H(q_\phi(z|x_i))}_{\text{Posterior Entropy}} + \underbrace{H(q_\phi(z|x_i), p(z))}_{\text{Cross Entropy}} \end{aligned}$$

Posterior에서 샘플링 된 z 는 최대한
다양해야 한다
(mode collapse 방지 효과)

Posterior와 Prior의
정보량은 유사해야 한다

Re-parameterization Trick ?

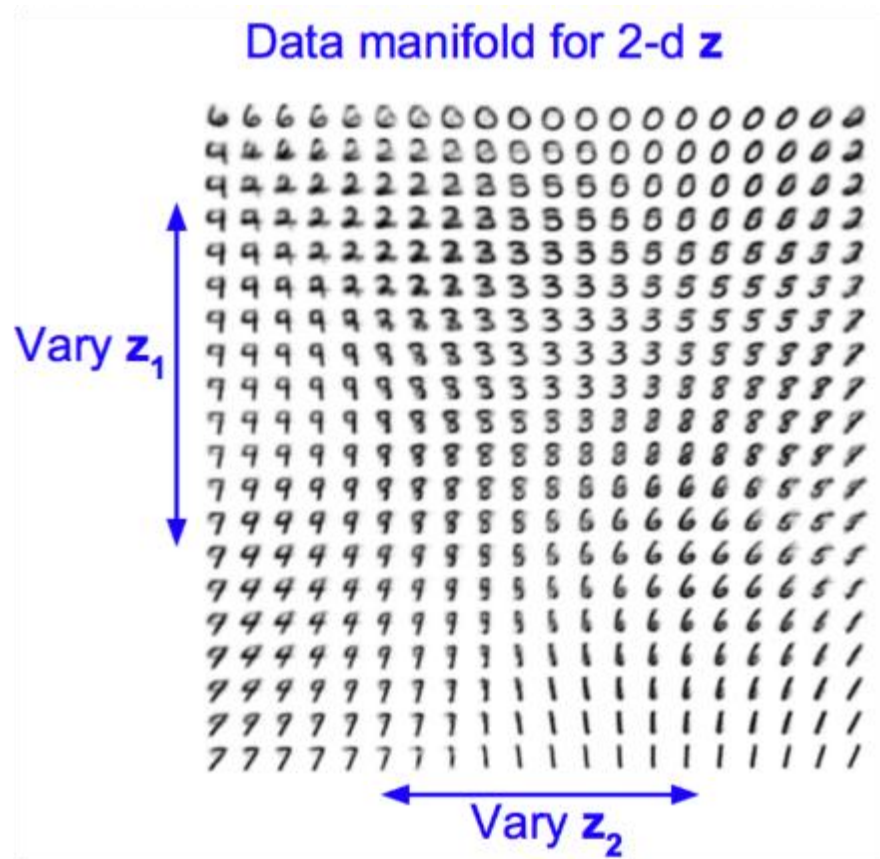
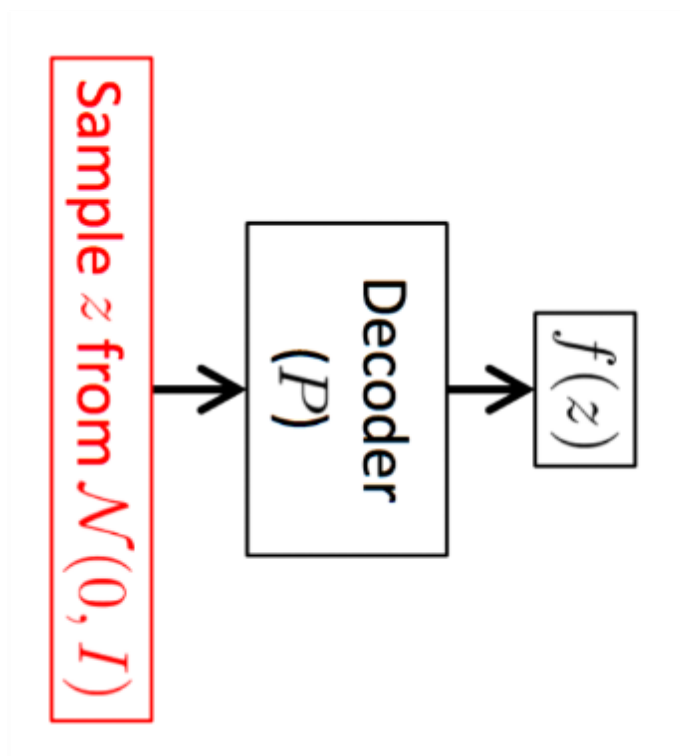


Sampling process $z^{i,l} \sim N(\mu_i, \sigma_i^2 I) \quad \Rightarrow \quad z^{i,l} = \mu_i + \sigma_i^2 \odot \epsilon$
 $\epsilon \sim N(0, I)$

같은 파이프라인/결과 이지만 Backprop을 가능하게 함

Generator

- Data를 생성할 때는 z 를 Decoder에 넣으면 된다.



MNIST Result



Input image



$J = |z| = 2$

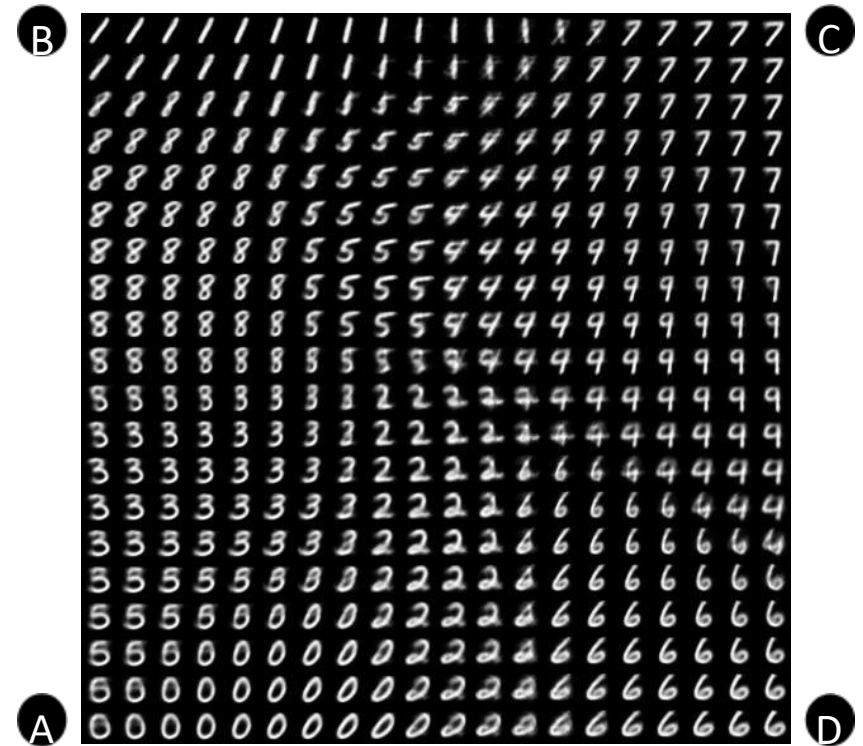
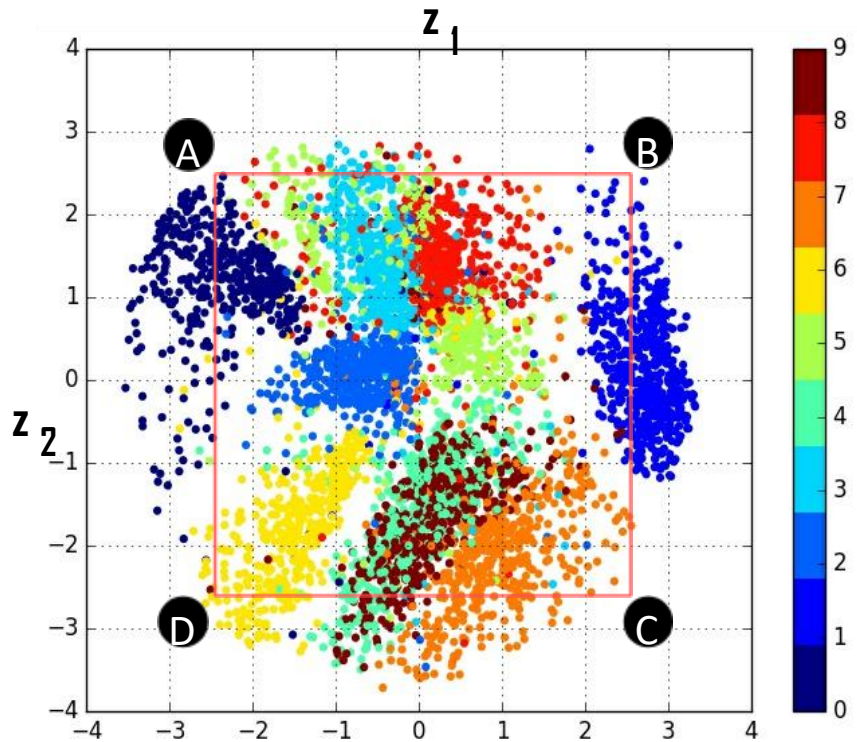


$J = |z| = 5$



$J = |z| = 20$

Learned Manifold (of latent space z)

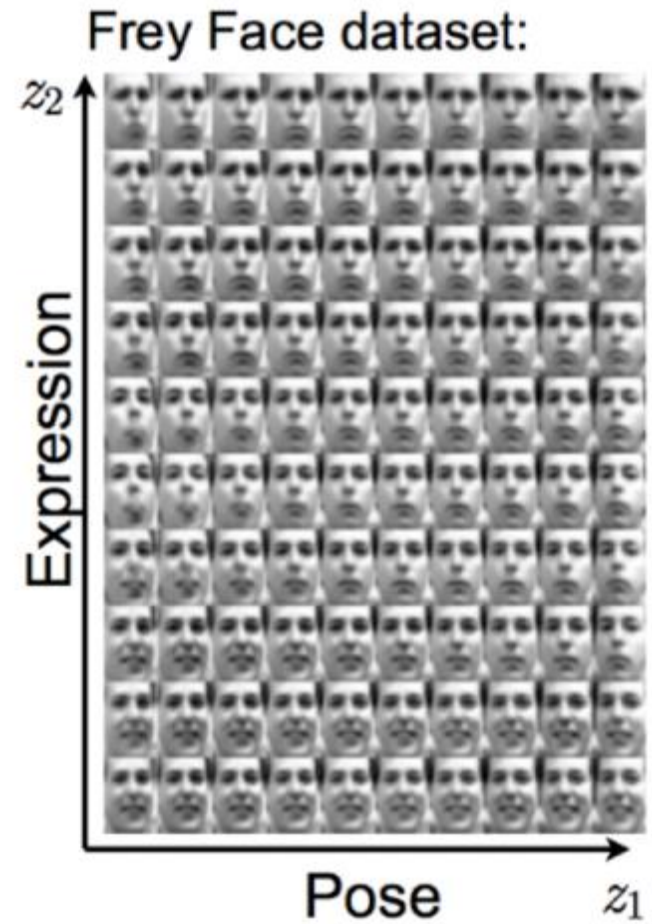


학습이 잘 되었을수록 2D 공간에서 같은 숫자들을 생성하는 z 들은 뭉쳐 있고, 다른 숫자들은 생성하는 z 들은 떨어져 있어야 한다.

Face Data 생성 예시



Labeled Faces in the Wild



VAE 장단점

- 학습이 안정적인편
 - 영상Generation뿐만아니라,주어진영상에서latent space의 posterior $q(z|x)$ 도함께학습시킨다. -> feature learning
 - 평가기준이명확하다 Reconstruction Error, lower bound (estimate likelihood of data)
-
- GAN에 비해서 출력이선명하지않고 평균값 형태로 표시되는문제
Reconstruction Loss -> Perceptual Loss
Complex Posterior Inference -> Normalizing Flow
 - Re-parameterization Trick이 Continuous한 함수에대해서만적용되는 문제
Gumbel Categorical Distribution

Thanks.