# AUTOENCODER

정인호

# Bayes rule?



$$P(H \mid e) = \frac{P(e \mid H)\,P(H)}{P(e)}$$

**Likelihood**
How probable is the evidence given that our hypothesis is true?

**Prior**
How probable was our hypothesis *before* observing the evidence?
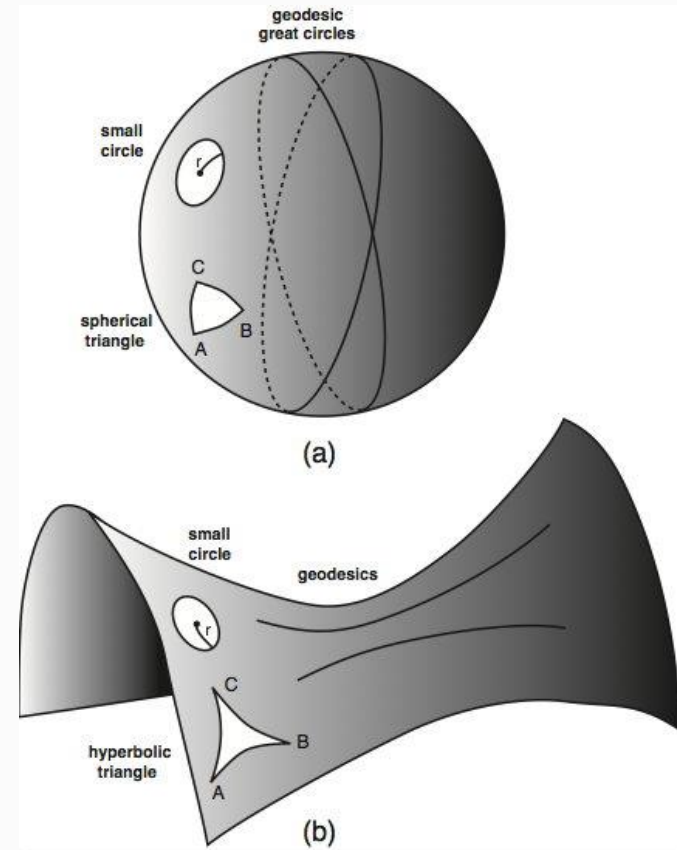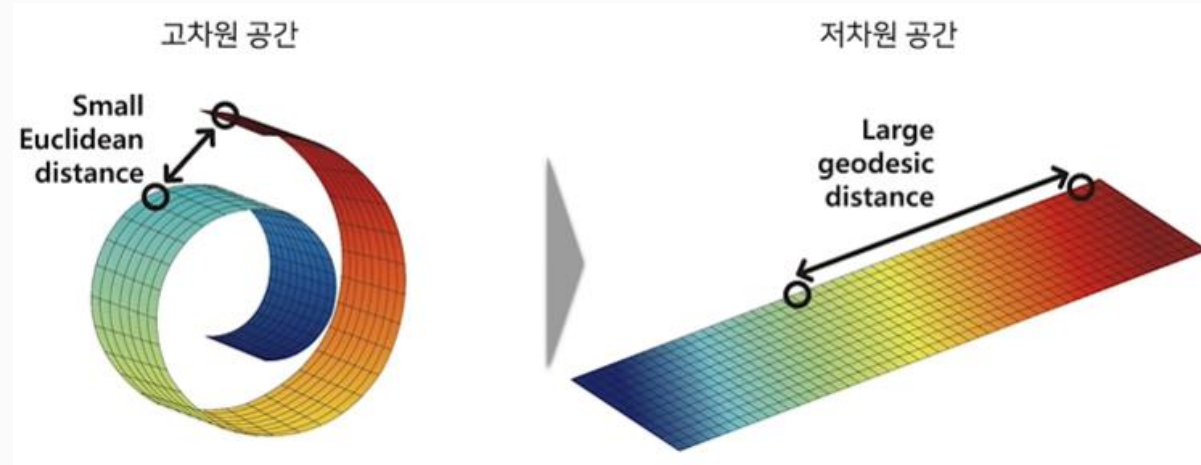
**Posterior**
How probable is our hypothesis given the observed evidence?
(Not directly computable)

**Marginal**
How probable is the new evidence under all possible hypotheses?
$P(e) = \sum P(e \mid H_i)\,P(H_i)$

# Manifold?

In mathematics, a manifold is a topological space that locally resembles Euclidean space near each point.

# KL-divergence

정의 (쿨백-라이블러 발산, Kullback-Leibler divergence)

두 확률밀도함수 $P(x)$, $g(x)$에 대해서

$$KL(P \| g) := \int_{\mathbb{R}} P(x) \log\left(\frac{P(x)}{g(x)}\right) dx$$
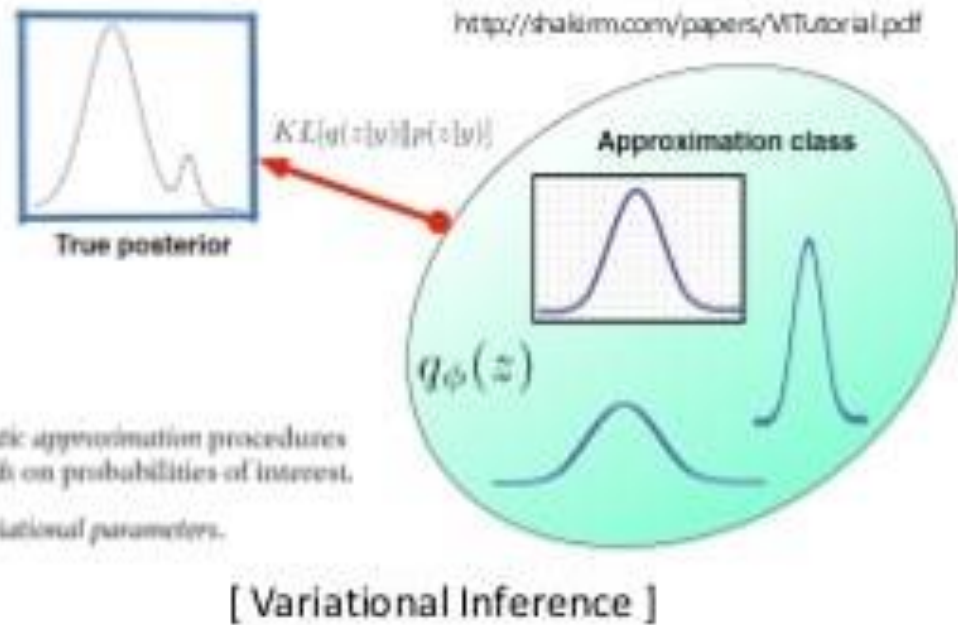
를 두 확률분포 간의 KL-divergence or relative entropy 라고 부른다.

KL-divergence의 성질
- 임의의 $P, g$에 대해서 $KL(P\|g) \geq 0$
- $KL(P\|g) = 0$ if and only if $P = g$.
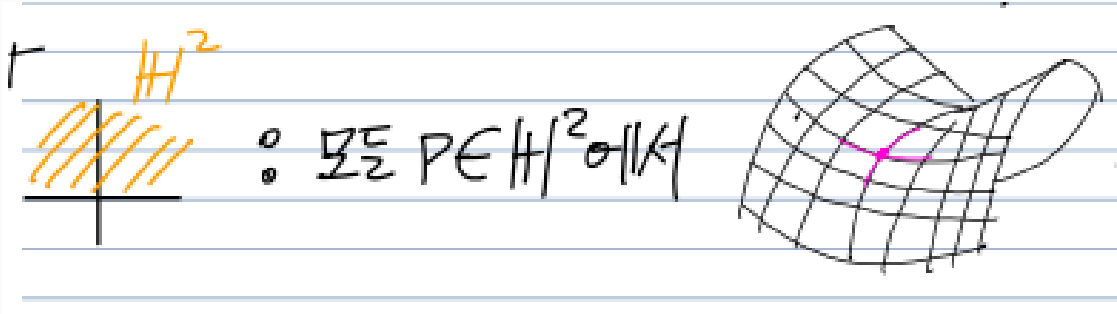
수학적으로 엄밀하게 말하면 metric은 아님. (삼각부등식 만족X)

# Information geometry
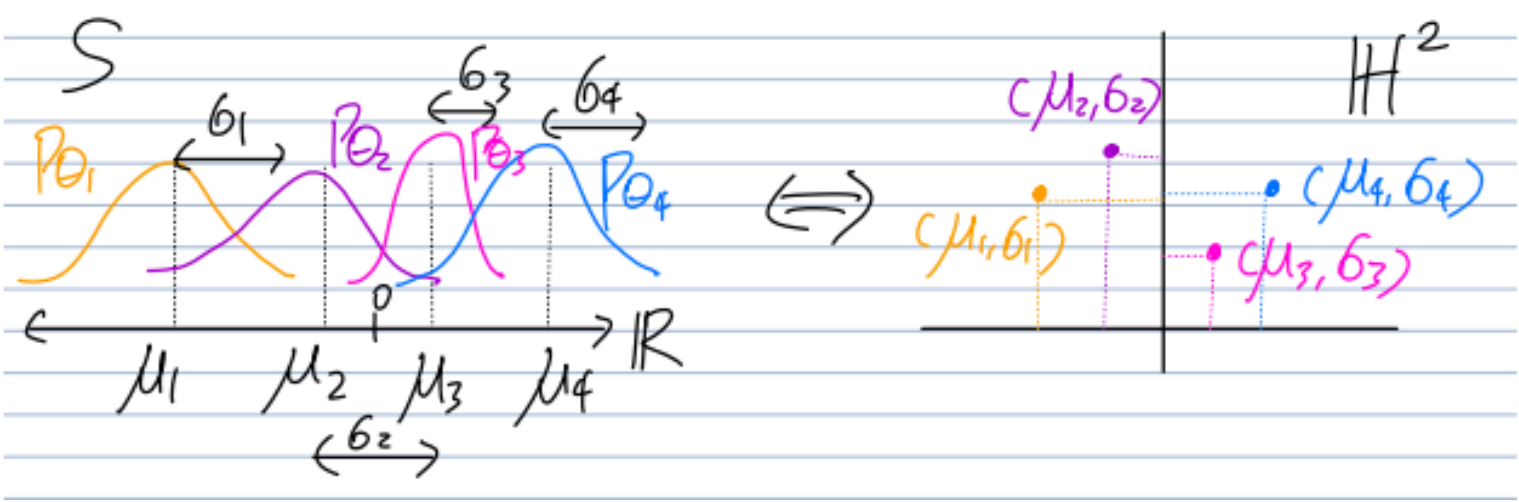
Q. The most reasonable Gaussian distribution?



[ Variational Inference ]

# Information geometry

# Information geometry

# Curse of dimensionality





200*200 RGB image

**차원이 증가할수록 데이터의 분포 분석**
**또는 모델추정에 필요한 샘플 게이터의 개수가**
**기하급수적으로 증가**

# Manifold Hypothesis

**Natural data in high dimensional spaces concentrates close to lower dimensional manifolds.**

# Manifold Hypothesis



Distance in high dimension          Distance in manifold



MNIST Data → 2D manifold

Entangled manifold          Disentangled manifold

Q. 4차원 이상에서 Disentangled 확인하는 방법?

# Autoencoder

Trained to reconstruct input $x$ as output $\tilde{x}$



Bottleneck Hidden Layer

Input layer

$x$

Output layer recostruct input

$y$

$z$

Encoding

Decoding

Undercomplete

Overcomplete

- Code
- Latent Variable
- Feature
- Hidden representation

http://videolectures.net/deeplearning2015_vincent_autoencoders/?q=vinc

Unsupervised Learning → Supervised Learning (Self Learning)

( + minimum 성능 보장)

# Linear Autoencoder = PCA



**LINEAR AUTOENCODER** | Multi-Layer Perceptron

**General Autoencoder**

latent vector $z \in \mathbb{R}^{d_z}$

$z = h(x)$

Encoder $h(\cdot)$    $g(\cdot)$ Decoder

input $x \in \mathbb{R}^d$    output $y \in \mathbb{R}^d$

$y = g(h(x))$

reconstruction error $L(x, y)$

Minimize $L_{AE} = \sum_{x \in D} L(x, g(h(x)))$

**Linear Autoencoder**

$$h(x) = W_e x + b_e$$

$$g(h(x)) = W_d z + b_d$$

$\|x - y\|^2$ or cross-entropy

Hidden layer 1개이고 레이어 간
fully-connected로 연결된 구조

http://videolectures.net/deeplearning2015_vincent_autoencoders/?q=vincent%20autoencoder

# Pretraining



Tied weighted.

W' W = I

# Denoising Autoencoder



Denoising AutoEnocder

latent vector $z \in \mathbb{R}^{d_z}$

$z = h(\tilde{x})$

Encoder $h(\cdot)$     $g(\cdot)$ Decoder

$y = g(h(\tilde{x}))$

corrupted input $\tilde{x} \in \mathbb{R}^d$

output $y \in \mathbb{R}^d$

add random noise $q(\tilde{x}|x)$

input $x \in \mathbb{R}^d$

$L(x, y)$ reconstruction error

Noise가 끼어 있지만 의미적으로 같아야 한다
(manifold 위에서 같은 점에 위치)

DAE | Performance – Visualization of learned filters

Natural image patches (12x12 pixels) : 100 hidden units

랜덤값으로 초기화하였기 때문에
노이즈처럼 보이는 필터일 수록 학습이
잘 안 된 것이고 edge filter와 같은 모습
일 수록 학습이 잘 된 것이다.

- Mean Squared Error
- 100 hidden units
- Salt-and-pepper noise

AE

AE with weight decay

DAE

10% salt-and-pepper noise

http://videolectures.net/deeplearning2015_vincent_autoencoders/?q=vincent%20autoencoder

# Variational Autoencoders

KEY WORD : Generative model
      (Autoencoder : manifold learning)



Latent Variable      Target Data

$z \sim p(z)$      Random variable

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

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

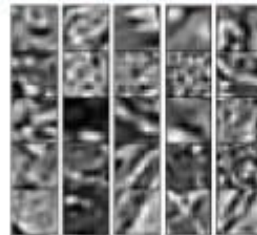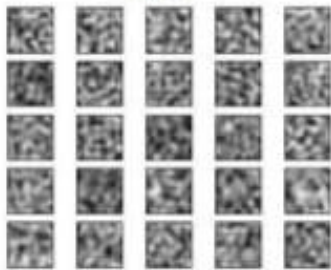Latent variable can be seen as a set of control parameters for target data (generated data)

For MNIST example, our model can be trained to generate image which match a digit value z randomly sampled from the set [0, ..., 9].

그래서, p(z)는 샘플링 하기 용이해야 편하다.

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

We are aiming maximize the probability of each x in the training set, under the entire generative process, according to:

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

**Question: Is it enough to model p(z) with simple distribution like normal distribution?**

**Yes**



Figure 2: Given a random variable $z$ with one distribution, we can create another random variable $X = g(z)$ with a completely different distribution. Left: samples from a gaussian distribution. Right: those same samples mapped through the function $g(z) = z/10 + z/||z||$ to form a ring. This is the strategy that VAEs use to create arbitrary distributions: the deterministic function $g$ is learned from data.

Tutorial on Variational Autoencoders : https://arxiv.org/pdf/1606.05908

## Question: Why don't we use maximum likelihood estimation directly?

$$p(x) \approx \sum_i p(x|g_\theta(z_i))p(z_i)$$

If $p(x|g_\theta(z)) = \mathcal{N}(x|g_\theta(z), \sigma^2 * I)$, the negative log probability of X is proportional squared Euclidean distance between $g_\theta(z)$ and $x$.

$x$ : Figure 3(a)

$z_{bad} \rightarrow g_\theta(z_{bad})$ : Figure 3(b)

$z_{bad} \rightarrow g_\theta(z_{good})$: Figure 3(c) (identical to x but shifted down and to the right by half a pixel)

$$\|x - z_{bad}\|^2 < \|x - z_{good}\|^2 \rightarrow p(x|g_\theta(z_{bad})) > p(x|g_\theta(z_{good}))$$

Solution 1: we should set the $\sigma$ hyperparameter of our Gaussian distribution such that this kind of erroroneous digit does not contribute to p(X) → hard..

Solution 2: we would likely need to sample many thousands of digits from $z_{good}$ → hard..



(a)          (b)          (c)

Figure 3: It's hard to measure the likelihood of images under a model using only sampling. Given an image X (a), the middle sample (b) is much closer in Euclidean distance than the one on the right (c). Because pixel distance is so different from perceptual distance, a sample needs to be extremely close in pixel distance to a datapoint X before it can be considered evidence that X is likely under the model.

생성기에 대한 확률모델을 가우시안으로 할 경우, MSE관점에서 가까운 것이 더 p(x)에 기여하는 바가 크다.
MSE가 더 작은 이미지가 의미적으로도 더 가까운 경우가 아닌 이미지들이 많기 때문에 현실적으로 올바른 확률값을 구하기가 어렵다.

Tutorial on Variational Autoencoders : https://arxiv.org/pdf/1606.05908

# Loss Function of VAE

**Relationship among** $p(x), p(z|x), q_\phi(z|x)$

**LOSS FUNCTION** | NeuralNet Perspective

$$\underset{\phi,\theta}{\arg\min} \sum_i \underbrace{-\mathbb{E}_{q_\phi(z|x_i)}\left[\log\left(p(x_i|g_\theta(z))\right)\right] + KL\left(q_\phi(z|x_i)||p(z)\right)}_{L_i(\phi,\theta,x_i)}$$



$x$    $q_\phi(\cdot)$    SAMPLING   $\sim z$    $g_\theta(\cdot)$    $x$

$q_\phi(z|x)$               $g_\theta(x|z)$

Encoder                 Decoder
Posterior               Generator
Inference Network      Generation Network

The mathematical basis of VAEs actually has relatively little to do with classical autoencoders

Tutorial on Variational Autoencoders : https://arxiv.org/pdf/1606.05908

# Loss Function of VAE

**LOSS FUNCTION** | Explanation

$$\underset{\phi,\theta}{\arg\min} \sum_i \underbrace{-\mathbb{E}_{q_\phi(z|x_i)}\left[\log(p(x_i|g_\theta(z)))\right] + KL(q_\phi(z|x_i)||p(z))}_{L_i(\phi,\theta,x_i)}$$

Variational inference를 위한
approximation class 중 선택

원 데이터에 대한 likelihood

다루기 쉬운 확률 분포 중 선택

$$L_i(\phi,\theta,x_i) = \underbrace{-\mathbb{E}_{q_\phi(z|x_i)}\left[\log(p(x_i|g_\theta(z)))\right]}_{\textbf{Reconstruction Error}} + \underbrace{KL(q_\phi(z|x_i)||p(z))}_{\textbf{Regularization}}$$

**Reconstruction Error**

- 현재 샘플링 용 함수에 대한 negative log likelihood
- $x_i$에 대한 복원 오차 (AutoEncoder 관점)
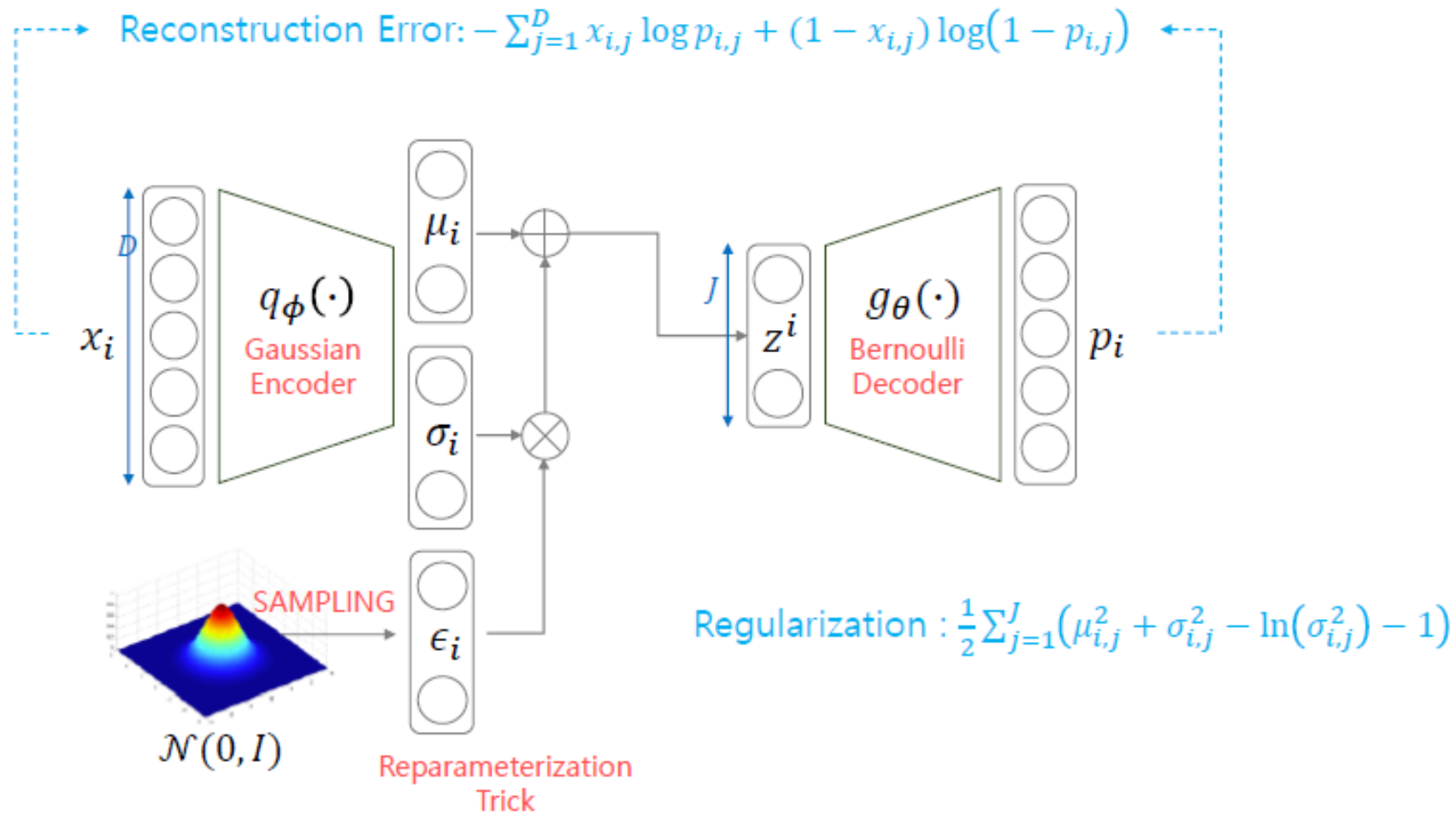
**Regularization**

- 현재 샘플링 용 함수에 대한 추가 조건
- 샘플링의 용의성/생성 데이터에 대한 통제성을 위한 조건을 prior에 부여 하고 이와 유사해야 한다는 조건을 부여

# Loss Function of VAE

**STRUCTURE** | Default : Gaussian Encoder + Bernoulli Decoder



Reconstruction Error: $-\sum_{j=1}^{D} x_{i,j} \log p_{i,j} + (1 - x_{i,j}) \log(1 - p_{i,j})$

$q_{\phi}(\cdot)$
Gaussian Encoder

$\mu_i$

$\sigma_i$

$x_i$

$D$

$\mathcal{N}(0,I)$

SAMPLING

$\epsilon_i$

Reparameterization Trick

$J$

$z^i$

$g_{\theta}(\cdot)$
Bernoulli Decoder

$p_i$

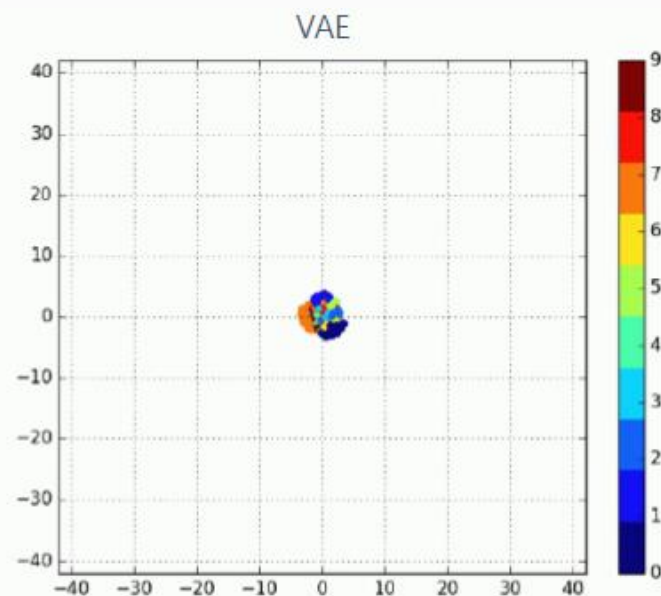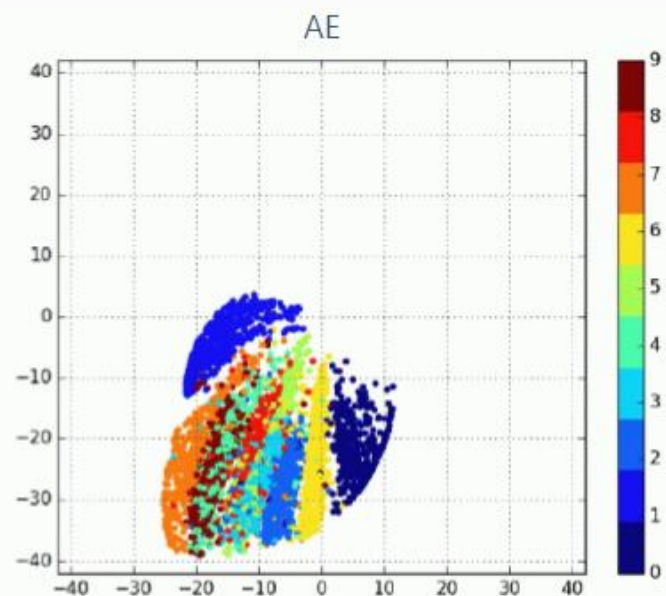Regularization : $\frac{1}{2}\sum_{j=1}^{J}\left(\mu_{i,j}^2 + \sigma_{i,j}^2 - \ln(\sigma_{i,j}^2) - 1\right)$

# AE vs VAE

AE & VAE 코드 관점에서 한 줄 다름.
(Loss function 에서 KL_divergence term 추가)

AE – for dimensionally reduction. Generating에는 적합하지 않음.
VAE – for generating. Manifold 위치가 안정적임.

Reconstruct만 할거면 AE Loss가 훨씬 저렴함.
AE는 text, image, sound, …, multi_model까지 domain을 가리지 않고 안정적이라고 함.



영상 자료 : https://www.youtube.com/watch?v=rNh2CrTFpm4
43분 15초~

# Adversarial Autoencoder (AAE)

## Adversarial Autoencoder

$$L_i(\phi, \theta, x_i) = -\mathbb{E}_{q_\phi(z|x_i)}\big[\log(p_\theta(x_i|z))\big] + \boxed{KL(q_\phi(z|x_i) \parallel p(z))}$$

**Regularization**

Conditions for $q_\phi(z|x_i)$, $p(z)$
1. Easily draw samples from distribution
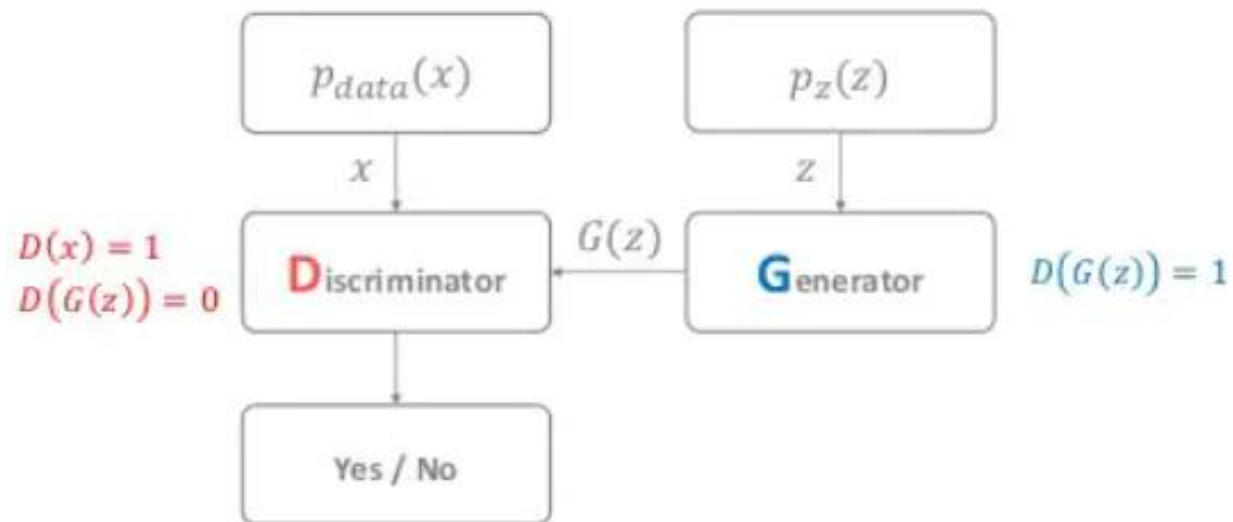2. KL divergence can be calculated

**Adversarial Autoencoder (AAE)**

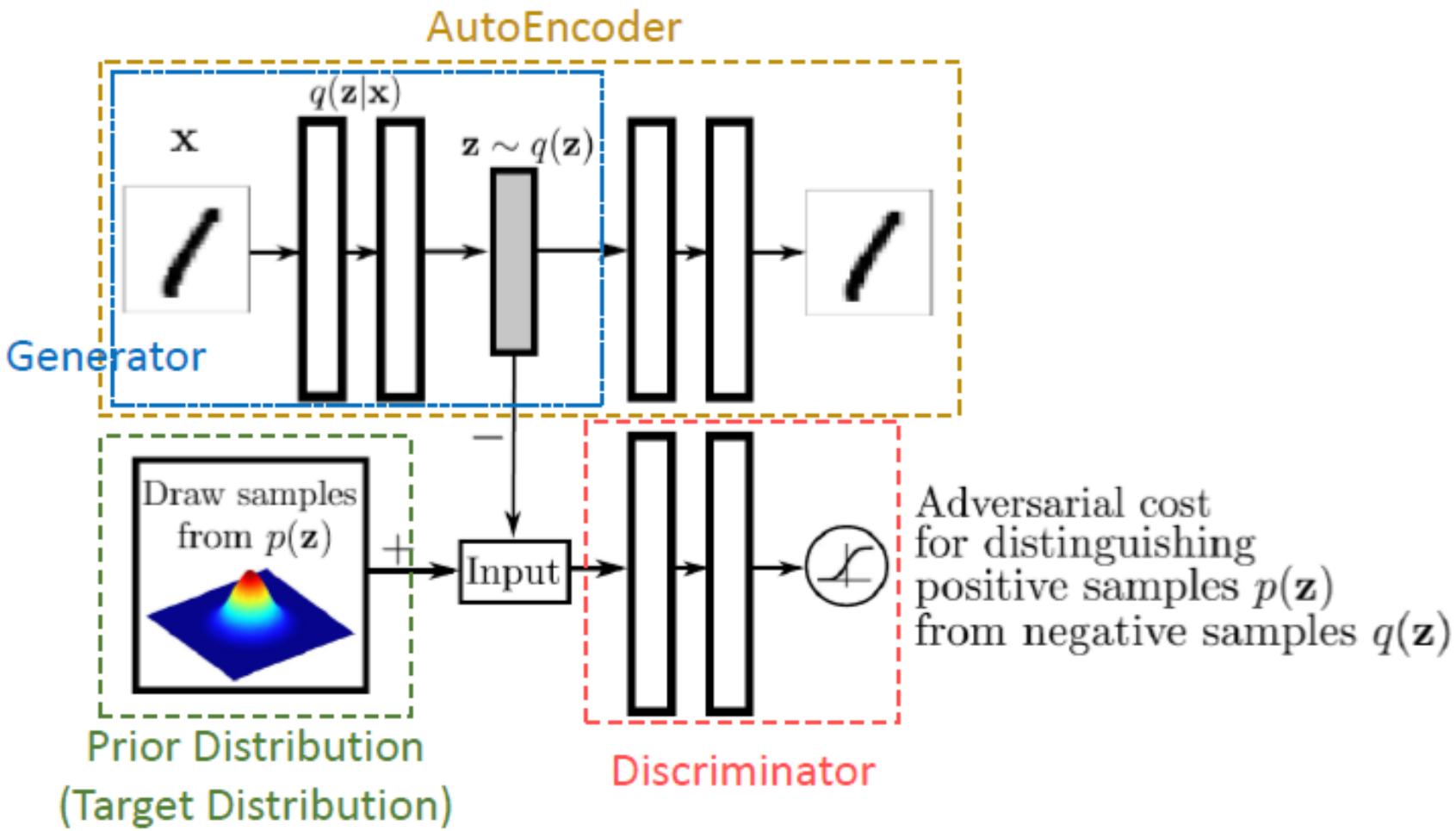| Conditions | $q_\phi(z|x_i)$ | $p(z)$ |
|---|---|---|
| Easily draw samples from distribution | O | O |
| KL divergence can be calculated | X | X |

**KL divergence is replaced by discriminator in GAN**

# AAE

## Generative Adversarial Network

$p_{data}(x)$ → $x$ → **D**iscriminator

$p_z(z)$ → $z$ → **G**enerator → $G(z)$ → **D**iscriminator

$D(x) = 1$
$D(G(z)) = 0$

$D(G(z)) = 1$

**D**iscriminator → Yes / No

Value function of GAN : $V(D, G) = \mathbb{E}_{x \sim p_{data}(x)}[\log D(x)] + \mathbb{E}_{z \sim p_z(z)}\left[\log\left(1 - D(G(z))\right)\right]$

Goal : $D^*, G^* = \min_{G} \max_{D} V(D, G)$    GAN은 $G(z) \sim p_{data}(x)$로 만드는 것이 목적이다

# AAE

# AAE

## Loss Function

**GAN loss**

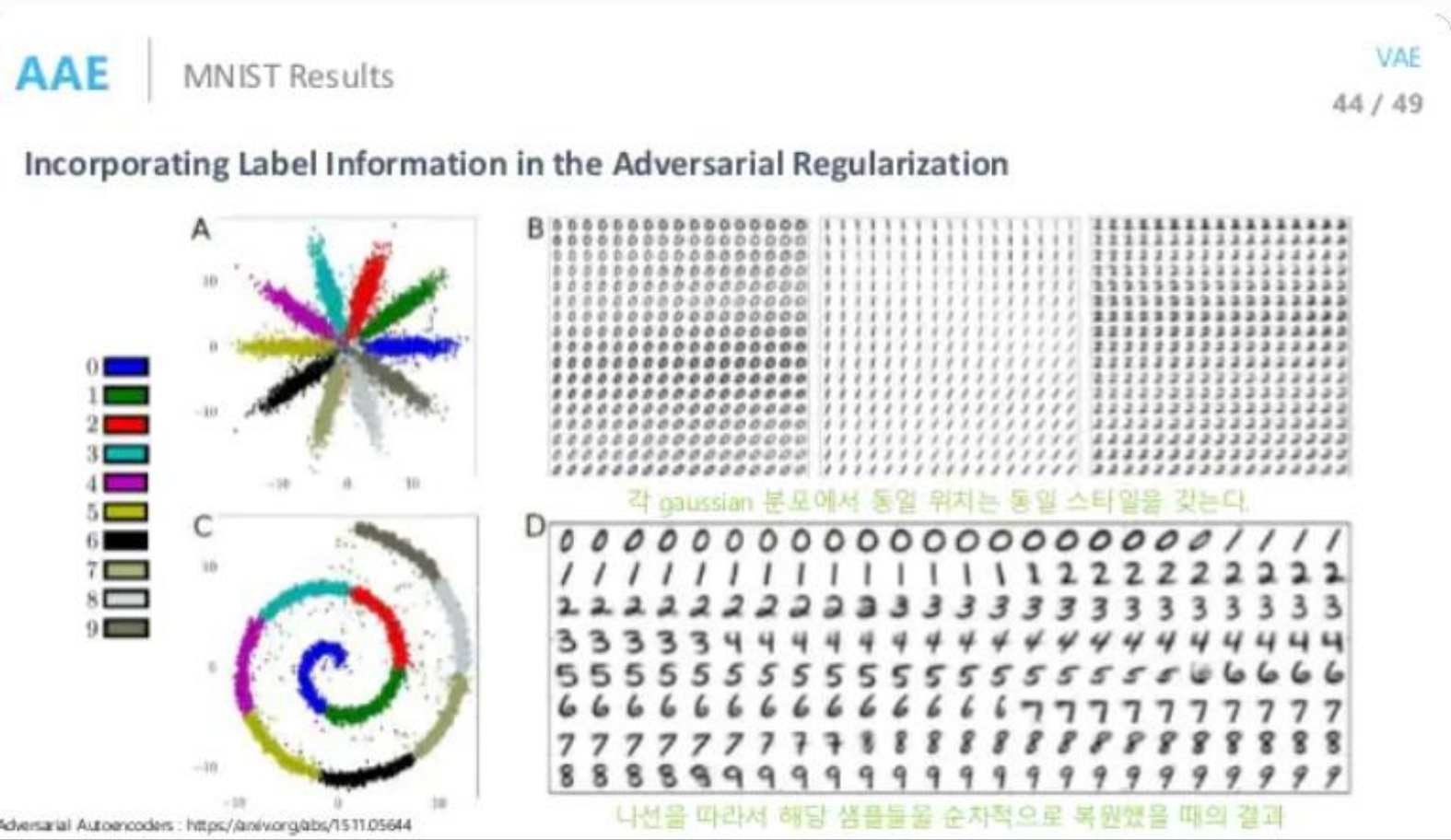$$V(D, G) = \mathbb{E}_{z \sim p(z)}[\log D(z)] + \mathbb{E}_{x \sim p(x)}\left[\log\left(1 - D(q_\phi(x))\right)\right]$$

Let's say G is defined by $q_\phi(\cdot)$ and D is defined by $d_\lambda(\cdot)$

$$V_i(\phi, \lambda, x_i, z_i) = \log d_\lambda(z_i) + \log\left(1 - d_\lambda(q_\phi(x_i))\right)$$
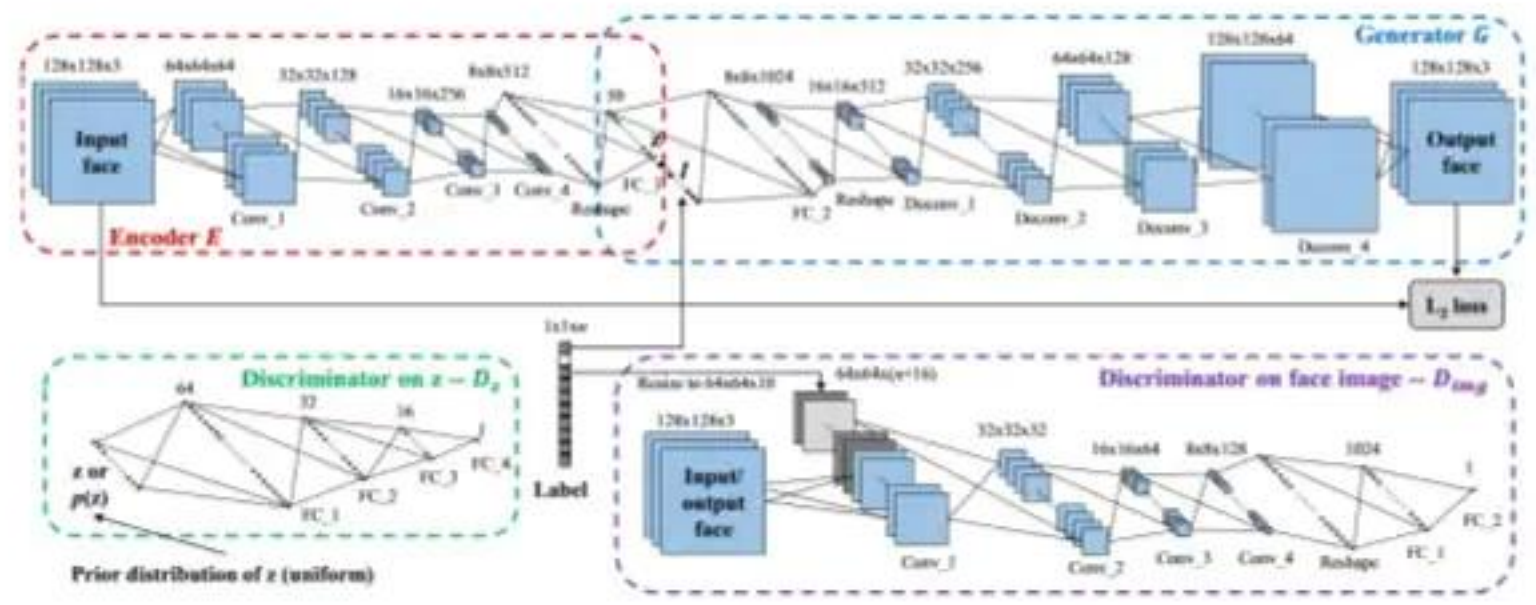
*논문에는 로스 정의가 제시되어 있지 않아 새로 정리한 내용

**VAE loss**

$$L_i(\phi, \theta, x_i) = -\mathbb{E}_{q_\phi(z|x_i)}\left[\log(p_\theta(x_i|z))\right] + KL(q_\phi(z|x_i)||p(z))$$

# AAE

AAE는 prior distribution을 원하는 모양으로 만들 수 있음.

# GAN + VAE...?

# Q & A

Reference
-오토인코더의 모든 것(이활석 NAVER)
(https://www.slideshare.net/NaverEngineering/ss-96581209)
-수학의 즐거움 정보기하
(https://www.youtube.com/watch?v=4s06EgHHRrA)