# 인공지능 기반 설계 이론 및 사례 연구

# 9차/10차) Variational AutoEncoder (VAE)
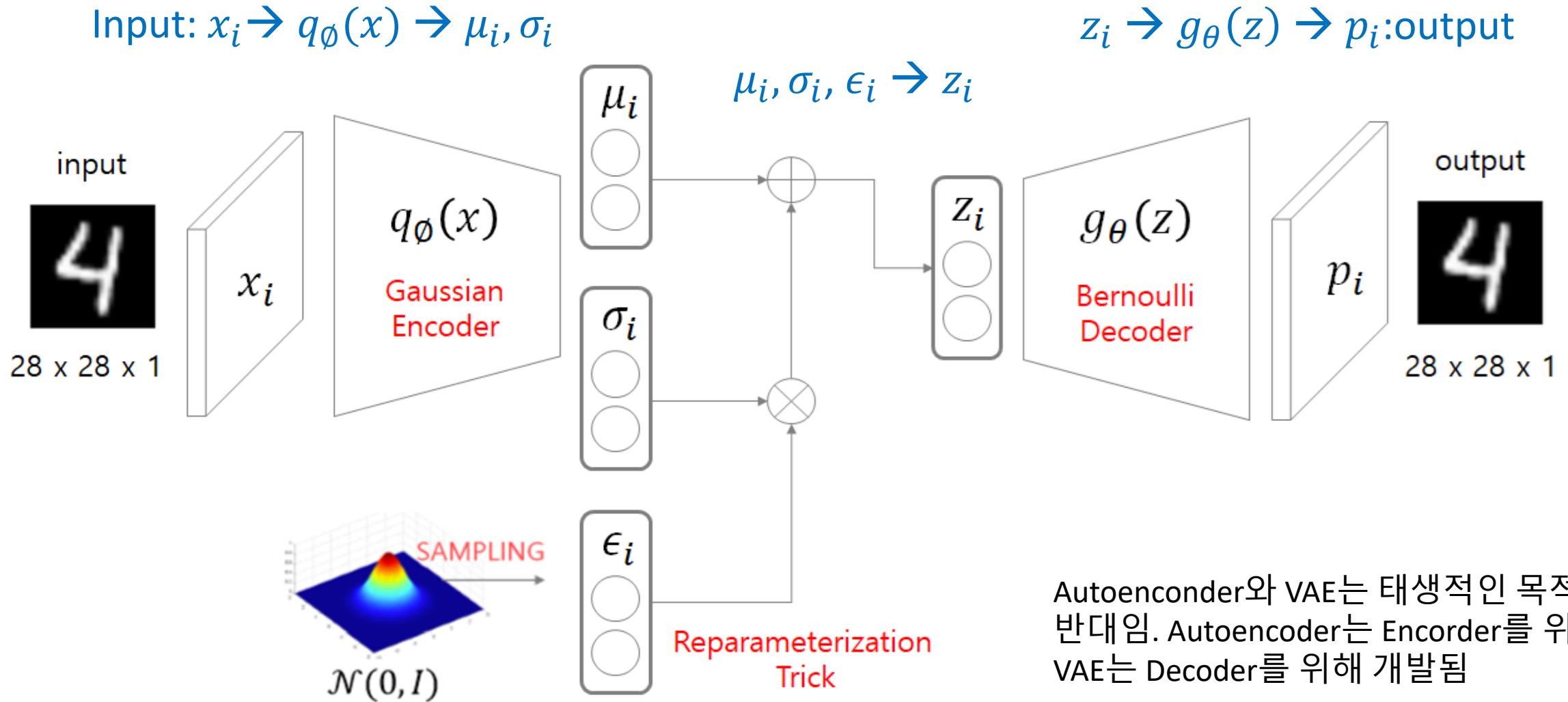
2020년 10월

## 강남우

기계시스템학부
숙명여자대학교

SDL SMART DESIGN LAB

# Variational Autoencoders (VAE) – How to work

Input: $x_i \rightarrow q_\emptyset(x) \rightarrow \mu_i, \sigma_i$

$z_i \rightarrow g_\theta(z) \rightarrow p_i$:output

$\mu_i, \sigma_i, \epsilon_i \rightarrow z_i$



input

$x_i$

$q_\emptyset(x)$

Gaussian Encoder

28 x 28 x 1

$\mu_i$

$\sigma_i$

SAMPLING

$\epsilon_i$

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

Reparameterization Trick

$z_i$

$g_\theta(z)$

Bernoulli Decoder

$p_i$

output

28 x 28 x 1

Autoenconder와 VAE는 태생적인 목적이 반대임. Autoencoder는 Encorder를 위해 VAE는 Decoder를 위해 개발됨
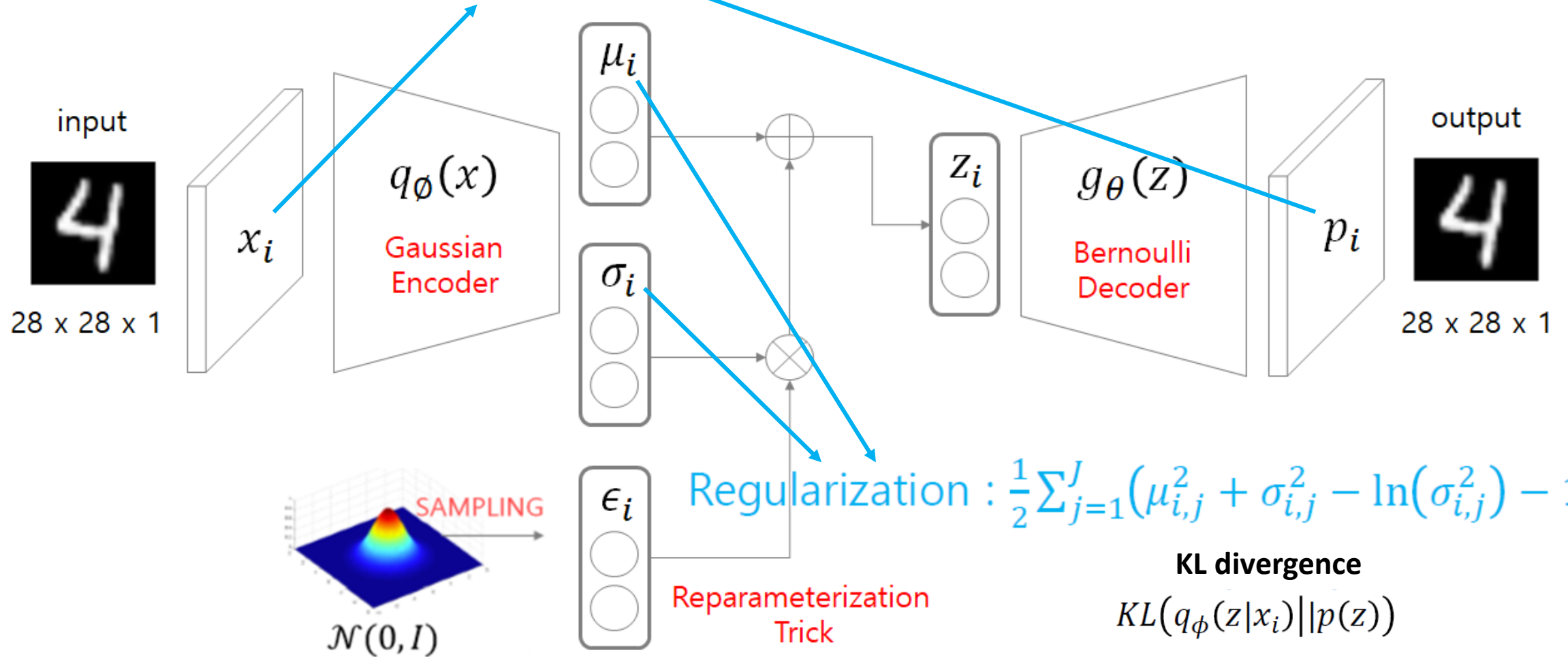
https://taeu.github.io/paper/deeplearning-paper-vae/

# VAE – How to work

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

**Cross entropy**

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

input

$q_\phi(x)$

Gaussian Encoder

$\mu_i$

$\sigma_i$

$z_i$

$g_\theta(z)$

Bernoulli Decoder

$p_i$

output

$x_i$

28 x 28 x 1

28 x 28 x 1

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

SAMPLING

$\epsilon_i$

Reparameterization Trick

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

**KL divergence**

$$KL(q_\phi(z|x_i)||p(z))$$

3  https://taeu.github.io/paper/deeplearning-paper-vae/

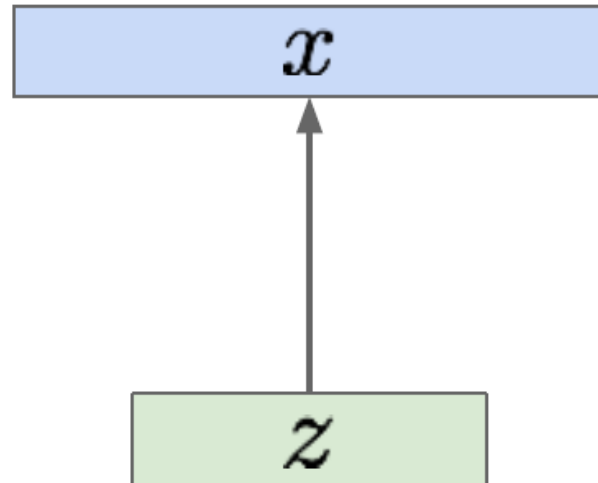Probabilistic spin on autoencoders - will let us sample from the model to generate data!

Assume training data $\{x^{(i)}\}_{i=1}^{N}$ is generated from underlying unobserved (latent) representation **z**

Sample from
true conditional
$p_{\theta^*}(x \mid z^{(i)})$

$x$

Sample from
true prior
$p_{\theta^*}(z)$

$z$

**Intuition** (remember from autoencoders!):
**x** is an image, **z** is latent factors used to
generate **x:** attributes, orientation, etc.

We want to estimate the true parameters $\theta*$ of this generative model.

How should we represent this model?

Sample from true conditional
$$p_{\theta*}(x \mid z^{(i)})$$

Sample from true prior
$$p_{\theta*}(z)$$

| $x$ |
| --- |

Decoder network

| $z$ |
| --- |

Choose prior $p(z)$ to be simple, e.g. Gaussian. Reasonable for latent attributes, e.g. pose, how much smile.

Conditional $p(x|z)$ is complex (generates image) => represent with neural network

We want to estimate the true parameters $\theta^*$ of this generative model.

How to train the model?
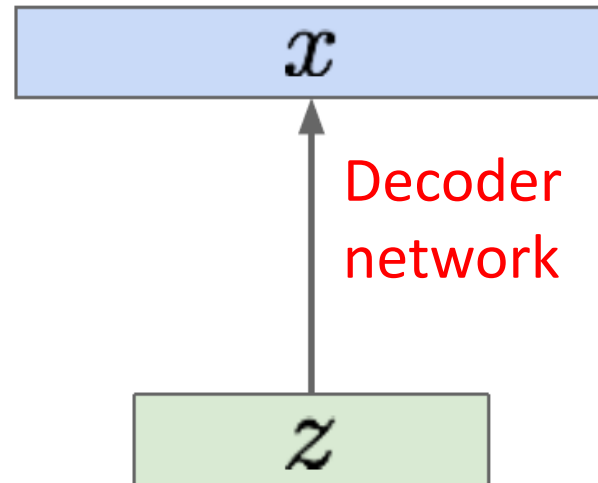
Sample from true conditional
$p_{\theta^*}(x \mid z^{(i)})$

Learn model parameters to maximize likelihood of training data



Decoder network

Sample from true prior
$p_{\theta^*}(z)$

$$p_\theta(x) = \int p_\theta(z) p_\theta(x|z) dz$$

Q: What is the problem with this?
Intractable!

Now with latent z

6

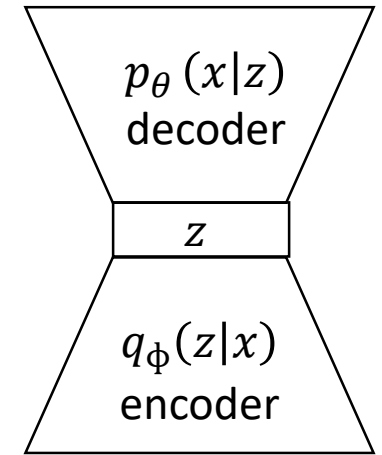# VAE - Loss Function

Intractible to compute
$p(x|z)$ for every z!

Data likelihood: $p_\theta(x) = \int p_\theta(z)p_\theta(x|z)dz$

Simple Gaussian prior          Decoder neural network

$p_\theta(x|z)$
decoder

$z$

$q_\phi(z|x)$
encoder

Posterior density also intractable: $p_\theta(z|x) = p_\theta(x|z)p_\theta(z)/p_\theta(x)$

Intractable data likelihood

Solution: In addition to decoder network modeling $p_\theta(x|z)$, define additional encoder network $q_\phi(z|x)$ that approximates $p_\theta(z|x)$

Will see that this allows us to derive a lower bound on the data likelihood that is tractable, which we can optimize

SMART DESIGN LAB

Now equipped with our encoder and decoder networks, let's work out the (log) data likelihood:

$$\log p_\theta(x^{(i)}) = \mathbf{E}_{z \sim q_\phi(z|x^{(i)})}\left[log p_\theta(x^{(i)})\right] \qquad (p_\theta(x^{(i)}) \text{ Does not depend on } z)$$

Taking expectation wrt. z (using encoder network) will come in handy later

$$= \mathbf{E}_z\left[log \frac{p_\theta(x^{(i)}|z)p_\theta(z)}{p_\theta(z|x^{(i)})}\right]$$

(Bayes' Rule)

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

$$= \mathbf{E}_z\left[log \frac{p_\theta(x^{(i)}|z)p_\theta(z)}{p_\theta(z|x^{(i)})} \frac{q_\phi(z|x^{(i)})}{q_\phi(z|x^{(i)})}\right] \qquad \text{(Multiply by constant)}$$
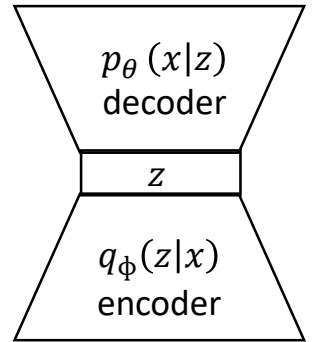
$p_\theta(x|z)$
decoder

$z$

$q_\phi(z|x)$
encoder

$$= \mathbf{E_z}\left[\log p_\theta(x^{(i)}|z)\right] - \mathbf{E_z}\left[log \frac{q_\phi(z|x^{(i)})}{p_\theta(z)}\right] + \mathbf{E}_z\left[log \frac{q_\phi(z|x^{(i)})}{p_\theta(z|x^{(i)})}\right] \qquad \text{(Logarithms)}$$

$$= \mathbf{E}_z\left[\log p_\theta(x^{(i)}|z)\right] - D_{kL}\left(q_\phi(z|x^{(i)}) \| p_\theta(z)\right) + D_{kL}\left(q_\phi(z|x^{(i)}) \| p_\theta(z|x^{(i)})\right)$$

참고: $E_{z \sim q_\phi(z|x^{(i)})}\left[log \frac{q_\phi(z|x^{(i)})}{p_\theta(z)}\right] = \int_z log \frac{q_\phi(z|x^{(i)})}{p_\theta(z)} q_\phi(z|x^{(i)}) dz$

$$KL(P\|Q) = \sum_x P(x) \log \frac{P(x)}{Q(x)}$$

The expectation wrt. z (using encoder network) let us write nice KL terms

SDL SMART DESIGN LAB

# VAE - Loss Function

Now equipped with our encoder and decoder networks, let's work out the (log) data likelihood:

$$\log p_\theta(x^{(i)}) = \mathbf{E}_{z \sim q_\phi(z|x^{(i)})}\left[\log p_\theta(x^{(i)})\right] \qquad (p_\theta(x^{(i)}) \text{ Does not depend on } z)$$
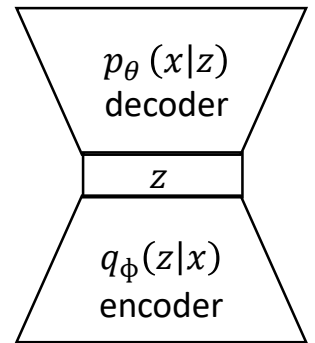
We want to maximize the data likelihood

$$= \mathbf{E}_z\left[\log \frac{p_\theta(x^{(i)}|z)p_\theta(z)}{p_\theta(z|x^{(i)})}\right] \text{ (Bayes' Rule)}$$

$$= \mathbf{E}_z\left[\log \frac{p_\theta(x^{(i)}|z)p_\theta(z)}{p_\theta(z|x^{(i)})} \frac{q_\phi(z|x^{(i)})}{q_\phi(z|x^{(i)})}\right] \text{ (Multiply by constant)}$$

$$= \mathbf{E}_z\left[\log p_\theta(x^{(i)}|z)\right] - \mathbf{E}_z\left[\log \frac{q_\phi(z|x^{(i)})}{p_\theta(z)}\right] + \mathbf{E}_z\left[\log \frac{q_\phi(z|x^{(i)})}{p_\theta(z|x^{(i)})}\right] \text{ (Logarithms)}$$

$$= \mathbf{E}_z\left[\log p_\theta(x^{(i)}|z)\right] - D_{kL}\left(q_\phi(z|x^{(i)}) \| p_\theta(z)\right) + D_{kL}\left(q_\phi(z|x^{(i)}) \| p_\theta(z|x^{(i)})\right)$$

$p_\theta(x|z)$ decoder

$z$

$q_\phi(z|x)$ encoder

Decoder network gives $p_\theta(x|z)$, can compute estimate of this term through sampling. (Sampling differentiable through reparam. trick, see paper.)

This KL term (between Gaussians for encoder and z prior) has nice closed-form solution!

$p_\theta(z|x)$ intractable (saw earlier), can't compute this KL term :( But we know KL divergence always >= 0.

**9**

SDL SMART DESIGN LAB

Now equipped with our encoder and decoder networks, let's work out the (log) data likelihood:

$$\log p_\theta(x^{(i)}) = \mathbf{E}_{z \sim q_\phi(z|x^{(i)})}\left[ log p_\theta(x^{(i)}) \right] \qquad (p_\theta(x^{(i)}) \; Does \; not \; depend \; on \; z)$$
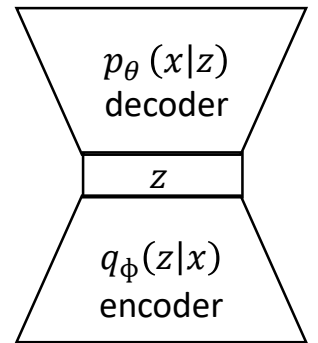
We want to maximize the data likelihood

$$= \mathbf{E}_z \left[ log \frac{p_\theta(x^{(i)}|z)p_\theta(z)}{p_\theta(z|x^{(i)})} \right] \quad \text{(Bayes' Rule)}$$

$$= \mathbf{E}_z \left[ log \frac{p_\theta(x^{(i)}|z)p_\theta(z)}{p_\theta(z|x^{(i)})} \frac{q_\phi(z|x^{(i)})}{q_\phi(z|x^{(i)})} \right] \quad \text{(Multiply by constant)}$$

$$= \mathbf{E}_z\left[ log \, p_\theta(x^{(i)}|z) \right] - \mathbf{E}_z\left[ log \frac{q_\phi(z|x^{(i)})}{p_\theta(z)} \right] + \mathbf{E}_z\left[ log \frac{q_\phi(z|x^{(i)})}{p_\theta(z|x^{(i)})} \right] \quad \text{(Logarithms)}$$

$$= \boxed{\mathbf{E}_z\left[ log \, p_\theta(x^{(i)}|z) \right] - D_{kL}\left( q_\phi(z|x^{(i)}) \, || \, p_\theta(z) \right)} + D_{kL}\left( q_\phi(z|x^{(i)}) \, || \, p_\theta(z|x^{(i)}) \right)$$

$$\mathcal{L}(x^{(i)}, \theta, \phi) \qquad\qquad\qquad \geq 0$$

**Tractable lower bound** which we can take gradient of and optimize! ($p_\theta(x|z)$ differentiable, KL term differentiable)

$p_\theta(x|z)$
decoder

$z$

$q_\phi(z|x)$
encoder

Now equipped with our encoder and decoder networks, let's work out the (log) data likelihood:

$$\log p_\theta(x^{(i)}) = \mathbf{E}_{z \sim q_\phi(z|x^{(i)})}\left[logp_\theta(x^{(i)})\right] \qquad (p_\theta(x^{(i)}) \ Does\ not\ depend\ on\ z)$$

We want to maximize the data likelihood

$$= \mathbf{E}_z\left[log\frac{p_\theta(x^{(i)}|z)p_\theta(z)}{p_\theta(z|x^{(i)})}\right] \text{ (Bayes' Rule)}$$

$$= \mathbf{E}_z\left[log\frac{p_\theta(x^{(i)}|z)p_\theta(z)}{p_\theta(z|x^{(i)})}\frac{q_\phi(z|x^{(i)})}{q_\phi(z|x^{(i)})}\right] \text{ (Multiply by constant)}$$

$p_\theta(x|z)$ decoder

$z$

$q_\phi(z|x)$ encoder

$$= \mathbf{E_z}\left[log\,p_\theta(x^{(i)}|z)\right] - \mathbf{E_z}\left[log\frac{q_\phi(z|x^{(i)})}{p_\theta(z)}\right] + \mathbf{E}_z\left[log\frac{q_\phi(z|x^{(i)})}{p_\theta(z|x^{(i)})}\right] \text{ (Logarithms)}$$

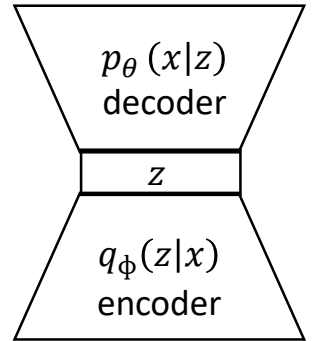$$= \underbrace{\mathbf{E_z}\left[log\,p_\theta(x^{(i)}|z)\right] - D_{kL}\left(q_\phi(z|x^{(i)})\,||p_\theta(z)\right)}_{\mathcal{L}(x^{(i)},\theta,\phi)} + \underbrace{D_{kL}\left(q_\phi(z|x^{(i)})\,||p_\theta(z|x^{(i)})\right)}_{\geq\,\mathbf{0}}$$

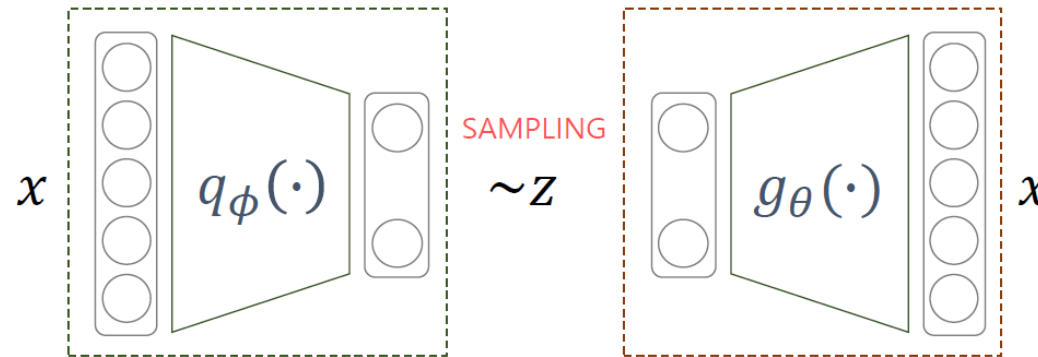$$log\,p_\theta(x^{(i)}) \geq \mathcal{L}(x^{(i)},\theta,\phi)$$
Variational lower bound ("ELBO")

$$\theta^*,\phi^* = arg\max_{\theta,\phi}\sum_{i=1}^{N}\mathcal{L}(x^{(i)},\theta,\phi)$$
Training: Maximize lower bound

SDL SMART DESIGN LAB

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

<span style="color:red">**Reconstruction Error**</span>  <span style="color:blue">**Regularization**</span>

- 현재 샘플링용 함수에 대한 negative log likelihood
- $x_i$에 대한 복원 오차 (Autoencoder 관점)

- 현재 샘플링용 함수에 대한 추가 조건
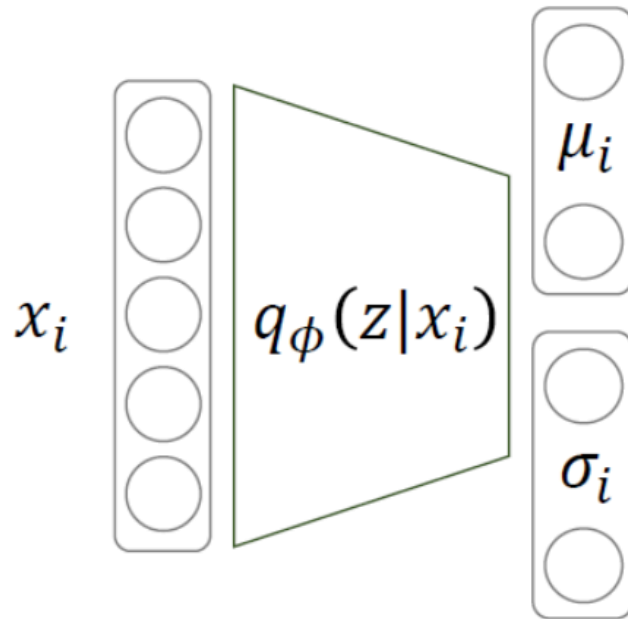- 샘플링의 용의성/생성 데이터에 대한 통제성을 위한 조건을 prior에 부여하고 이와 유사해야 한다는 조건을 부여

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

**Assumptions**

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

Regularization



$q_\phi(z|x_i)$

$x_i$

$\mu_i$

$\sigma_i$

## Assumption 1

[Encoder : approximation class]
multivariate gaussian distribution with a diagonal covariance

$$q_\phi(z|x_i) \sim N(\mu_i, \sigma_i^2 I)$$
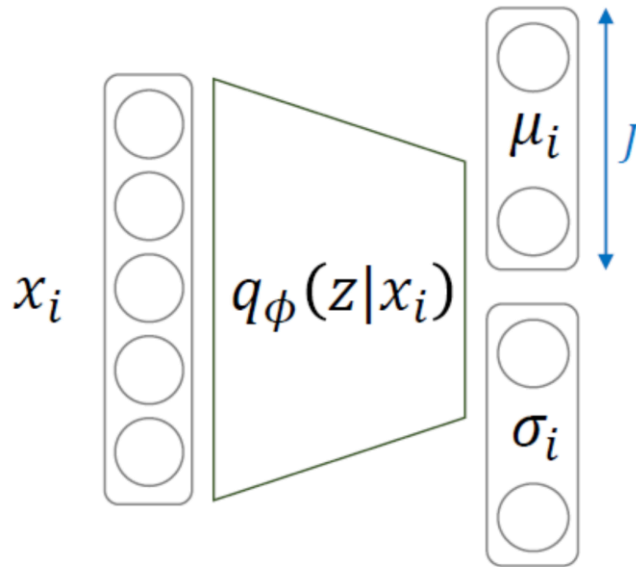
## Assumption 2

[prior] multivariate normal distribution

$$p(z) \sim N(0, I)$$

**KLD**

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

<p style="text-align:center; color:red;">Regularization</p>

$$KL(q_\phi(z|x_i)||p(z)) = \frac{1}{2}\left\{tr(\sigma_i^2 I) + \mu_i^T \mu_i - J + ln\frac{1}{\prod_{j=1}^{J}\sigma_{i,j}^2}\right\}$$

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

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

$x_i$  $q_\phi(z|x_i)$  $\mu_i$  $J$  $\sigma_i$

**KLD for multivariate normal distributions**

$$D_{KL}(\mathcal{N}_0 \| \mathcal{N}_1) = \frac{1}{2}\left(tr(\Sigma_1^{-1}\Sigma_0) + (\mu_1 - \mu_0)^\mathsf{T}\Sigma_1^{-1}(\mu_1 - \mu_0) - k + ln\left(\frac{\det \Sigma_1}{\det \Sigma_0}\right)\right)$$
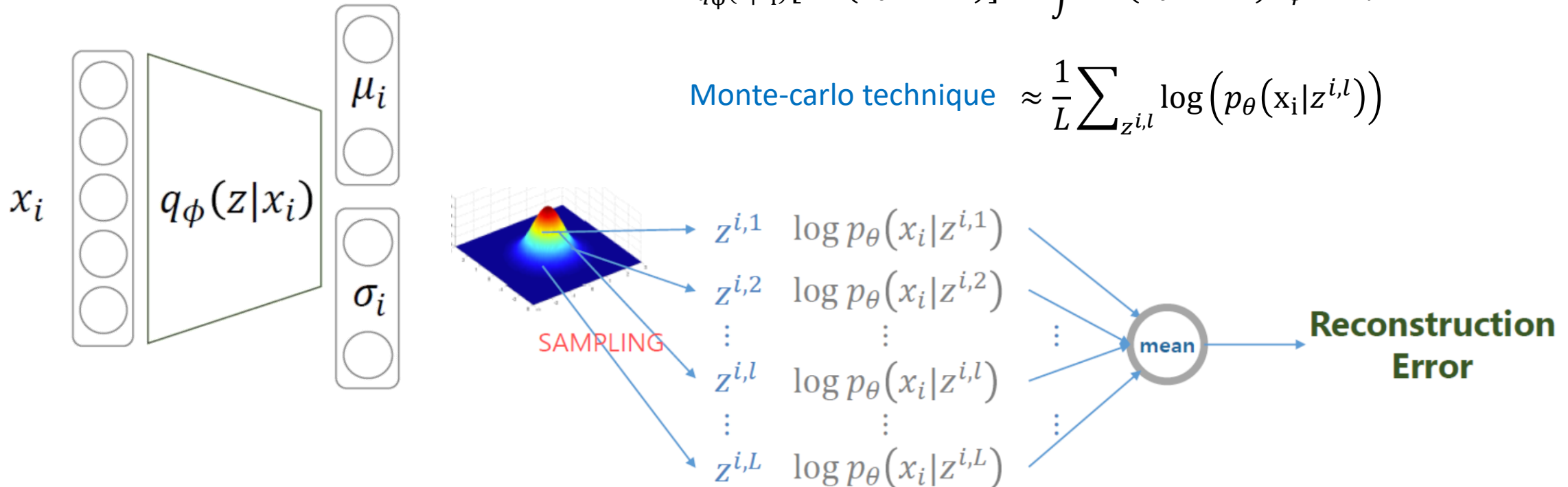
SDL SMART DESIGN LAB

**Sampling**

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

Reconstruction Error

$$\mathbb{E}_{q_\phi(z|x_i)}[\log(p_\theta(x_i|z))] = \int \log(p_\theta(x_i|z))q_\phi(z|x_i)dz$$

Monte-carlo technique $\approx \dfrac{1}{L}\sum_{z^{i,l}} \log\left(p_\theta(x_i|z^{i,l})\right)$



$x_i$   $q_\phi(z|x_i)$   $\mu_i$   $\sigma_i$

SAMPLING

$z^{i,1}$   $\log p_\theta(x_i|z^{i,1})$

$z^{i,2}$   $\log p_\theta(x_i|z^{i,2})$

$z^{i,l}$   $\log p_\theta(x_i|z^{i,l})$

$z^{i,L}$   $\log p_\theta(x_i|z^{i,L})$

mean → **Reconstruction Error**

- L is the number of samples for latent vector
- Usually L is set to 1 for convenience

SDL SMART DESIGN LAB

## *Reparameterization Trick*

**Sampling process**

$$z^{i,l} \sim N(\mu_i, \sigma_i^2 I)$$

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

Same distribution!
But it makes backpropagation possible!

**Assumptions**

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

<span style="color:red">Reconstruction Error</span>

$$\mathbb{E}_{q_\phi(z|x_i)}\big[log\big(p_\theta(x_i|z)\big)\big] = \int log\big(p_\theta(x_i|z)\big)\, q_\phi(z|x_i)dz \approx \frac{1}{L}\sum_{z^{i,l}} log\big(p_\theta(x_i|z^{i,l})\big) \approx log\big(p_\theta(x_i|z^i)\big)$$

Monte-carlo
technique

L=1

**Assumption 3-1**

[Decoder, likelihood]
multivariate bernoulli or gaussain distribution

$$z^i \quad g_\theta(\cdot) \quad p_i$$

$$D$$

$$log\big(p_\theta(x_i|z^i)\big) = log\prod_{j=1}^{D} p_\theta(x_{i,j}|z^i) = \sum_{j=1}^{D} log\, p_\theta(x_{i,j}|z^i)$$

$$= \sum_{j=1}^{D} log\, p_{i,j}^{x_{i,j}}(1-p_{i,j})^{1-x_{i,j}} \quad \longleftarrow \quad \color{red}{p_{i,j}: \text{network output}}$$

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

$$p_\theta(x_i|z^i) \sim Bernoulli(p_i)$$

<span style="color:red">**Cross entropy**</span>

**Assumptions**

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

<span style="color:red">Reconstruction Error</span>

$$\mathbb{E}_{q_\phi(z|x_i)}[log(p_\theta(x_i|z))] \approx log\big(p_\theta(x_i|z^i)\big)$$

## Assumption 3-2

[Decoder, likelihood]
multivariate bernoulli or gaussain distribution

$$log\big(p_\theta(x_i|z^i)\big) = log(N(x_i; \mu_i, \sigma_i^2 I))$$

$$= -\sum_{j=1}^{D} \frac{1}{2}log(\sigma_{i,j}^2) + \frac{(x_{i,j} - \mu_{i,j})^2}{2\sigma_{i,j}^2}$$
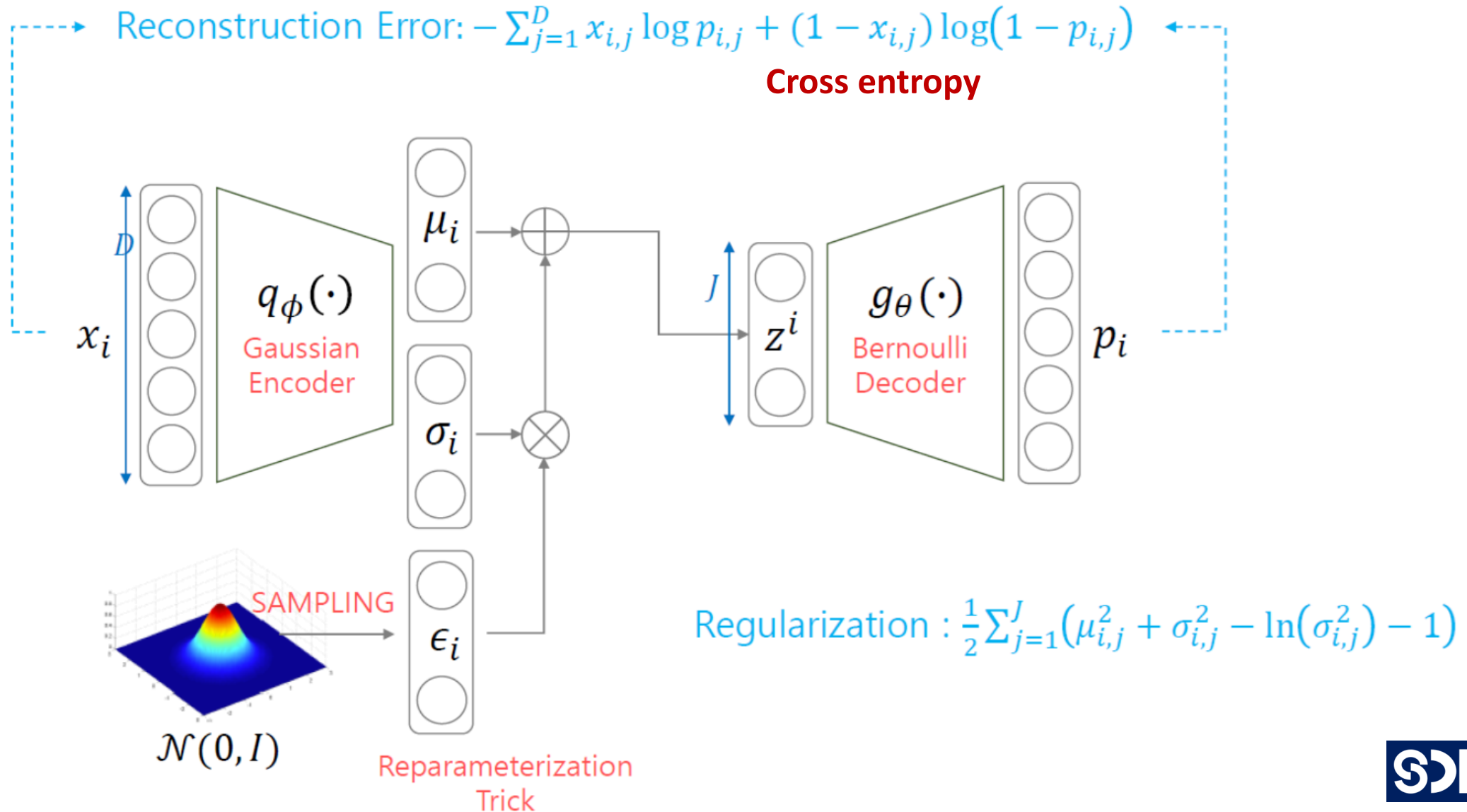
For gaussain distribution with identity covariance

$$log\big(p_\theta(x_i|z^i)\big) \propto -\sum_{j=1}^{D} (x_{i,j} - \mu_{i,j})^2 \qquad \text{\color{red}\textbf{Squared Error}}$$



$z^i \quad g_\theta(\cdot) \quad \mu_i \quad \sigma_i$

**Default : Gaussian Encoder + Bernoulli Decoder**

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

**Cross entropy**



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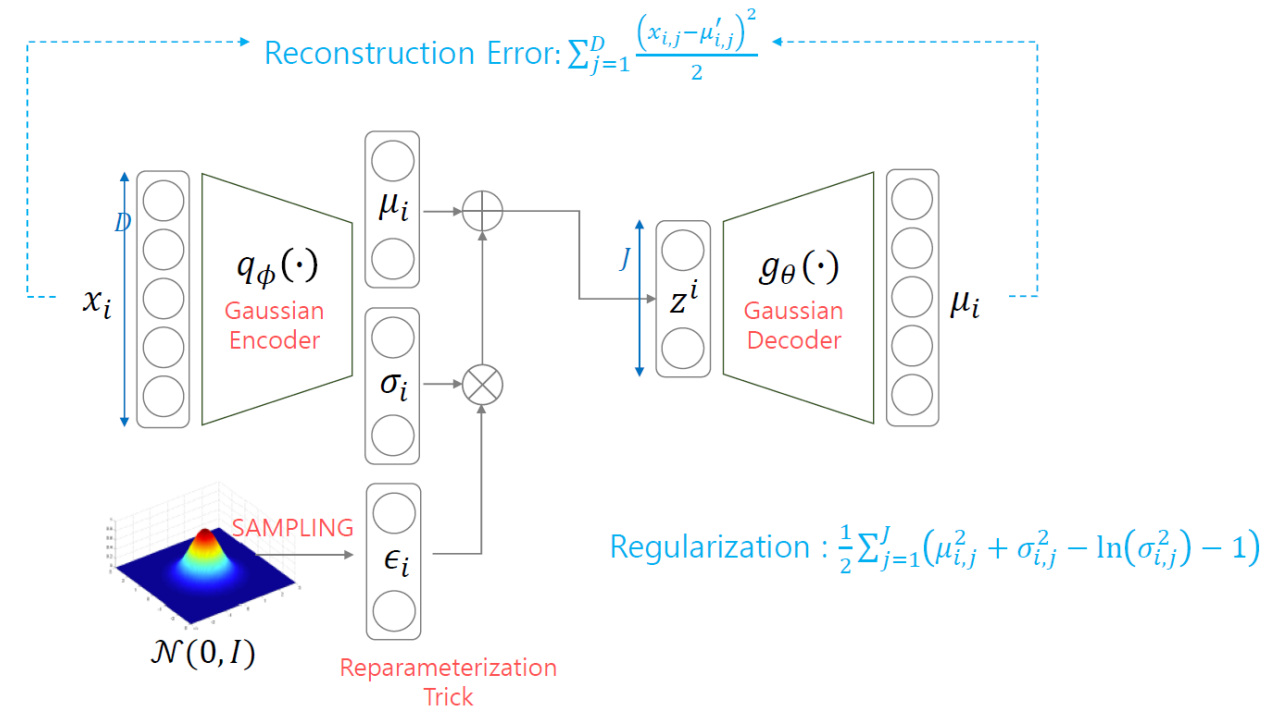
## Gaussian Encoder + Gaussian Decoder

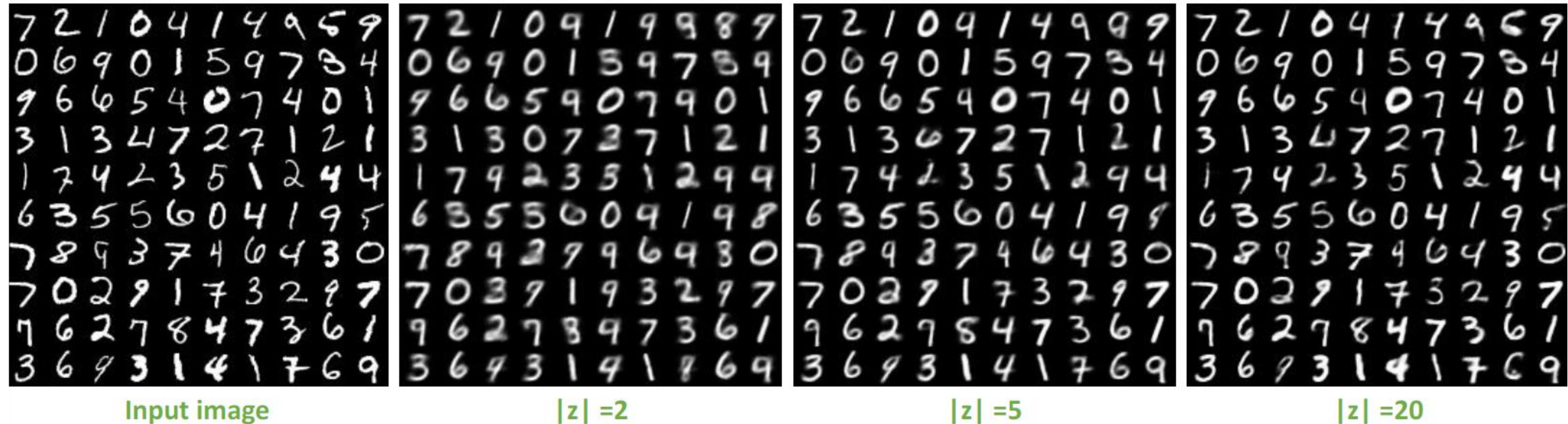$$\sum_{j=1}^{D} \frac{1}{2}\log(\sigma_{i,j}^2) + \frac{(x_{i,j} - \mu_{i,j})^2}{2\sigma_{i,j}^2}$$

## Gaussian Encoder + Gaussian Decoder with Identity Covariance

$$\sum_{j=1}^{D} (x_{i,j} - \mu_{i,j})^2 \quad \textbf{\textcolor{red}{Squared Error}}$$



Reconstruction Error: $\sum_{j=1}^{D} \frac{1}{2}\log(\sigma_{i,j}'^2) + \frac{(x_{i,j} - \mu_{i,j}')^2}{2\sigma_{i,j}'^2}$

$x_i$, $D$, $q_\phi(\cdot)$ Gaussian Encoder, $\mu_i$, $\sigma_i$, $z^i$, $J$, $g_\theta(\cdot)$ Gaussian Decoder, $\mu_i'$, $\sigma_i'$, $\epsilon_i$, $\mathcal{N}(0, I)$, SAMPLING, Reparameterization Trick

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



Reconstruction Error: $\sum_{j=1}^{D} \frac{(x_{i,j} - \mu_{i,j}')^2}{2}$

$x_i$, $D$, $q_\phi(\cdot)$ Gaussian Encoder, $\mu_i$, $\sigma_i$, $z^i$, $J$, $g_\theta(\cdot)$ Gaussian Decoder, $\mu_i$, $\epsilon_i$, $\mathcal{N}(0, I)$, SAMPLING, Reparameterization Trick

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

**Latent variable dimensions**



Input image

|z| =2

|z| =5

|z| =20

# VAE – Characteristics

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

**복원 오차**

입력과 출력 간의 *cross-entropy*

**제약 조건**

*Prior* 분포와의 다른 정도

- Probabilistic spin to traditional autoencoders ➔ allows generating data
- Defines an intractable density ➔ derive and optimize a (variational) lower bound

[ **VAE**의 **특징들** ]

1. **Decoder**가 <u>최소한</u> 학습 데이터는 생성해 낼 수 있게 된다.
➔ 생성된 데이터가 학습 데이터 좀 닮아 있다.

2. **Encoder**가 <u>최소한</u> 학습 데이터는 잘 latent vector로 표현할 수 있게 된다.
➔ 데이터의 추상화를 위해 많이 사용된다.

$y_i = x_i$

Decoder
Bernoulli
$p_\theta(x|z)$

$z_i$

$\mu_i$    $\sigma_i$

$q_\phi(z|x_i)$
Gaussian
Encoder

$\epsilon_i$

$\epsilon \sim N(0, I)$

$x_i$

# VAE coding

$$L_i(\phi, \theta, x_i) = -\mathbb{E}_{q_\phi(z|x_i)}[\log\left(p(x_i|g_\theta(z))\right)] + KL(q_\phi(z|x_i)||p(z)) \quad \Rightarrow \quad argmax \text{ ELBO}(\phi)$$

*Reconstruction Error*
원데이터에 대한 Log Likelihood

*Regularization*
다루기 쉬운 확률 분포 중 선택해서 변이추론을 위한 근사 class중 선택하여
유사해야 한다는 조건을 부여함.

**코딩에 적용된 수식**

[*Regularization* : **Kullback − leibler divergence**]

$$KL(q_\phi(z|x_i)||p(z)) = \frac{1}{2}\sum_{j=1}^{J} (\mu_{i,j}^2 + \sigma_{i,j}^2 - \ln(\sigma_{i,j}^2) - 1)$$

[*Reconstruction Error*]

$$-\mathbb{E}_{q_\phi(z|x_i)}\left[\log\left(p(x_i|g_\theta(z))\right)\right] = \int \log\left(p(x_i|g_\theta(z))\right) \approx \frac{1}{L}\sum_{z^{i,l}} log(p_\theta(x_i|z^{i,l})) \approx \log(p_\theta(x_i|z^{i,l})) = \sum_{j=1}^{D} x_{i,j}logp_{i,j} + (1 - x_{i,j})log(1 - p_{i,j})$$

Monte-carlo technique

For Bernoulli = cross-entropy

For Gaussian distibition
= mean square Error

# What Questions Do You Have?

nwkang@sm.ac.kr

www.smartdesignlab.org