

AE & GAN

Adversarial autoencoders

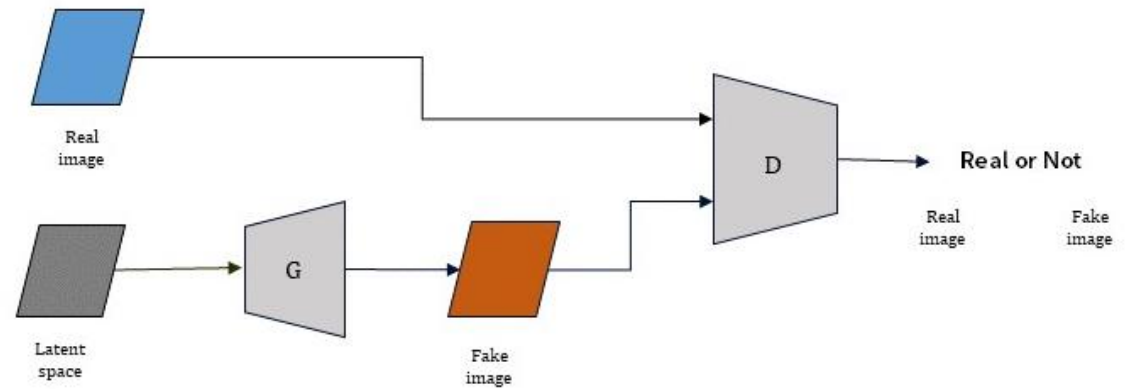
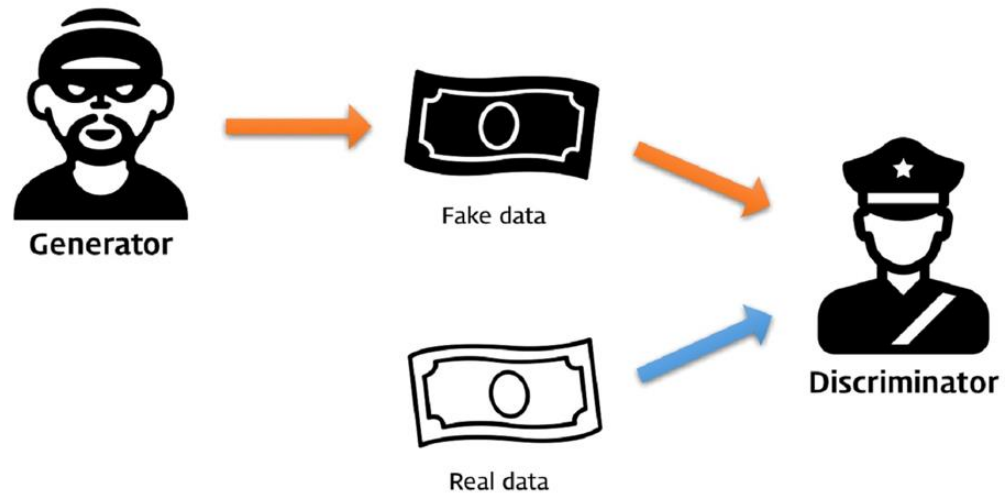
Makhzani, Alireza, et al. "Adversarial autoencoders."
arXiv preprint arXiv:1511.05644 (2015).

김수형

- What is GAN?
(feat. Generative model)
- What is AE?
(feat. manifold Hypothesis)
- VAE(variational autoencoder)
- AAE(adversarial autoencoders)

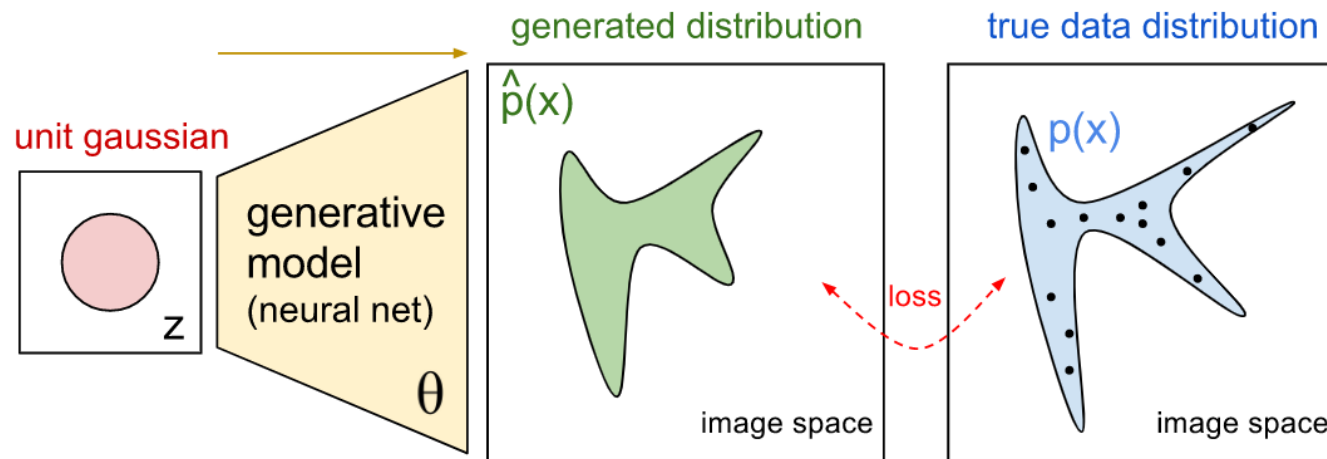
- **GAN(generative adversarial network)**

✓ 그럴듯한 가짜를 만드는 적대적 생성 모델



- **Generative model**

- ✓ Target data의 분포를 학습하는 모델

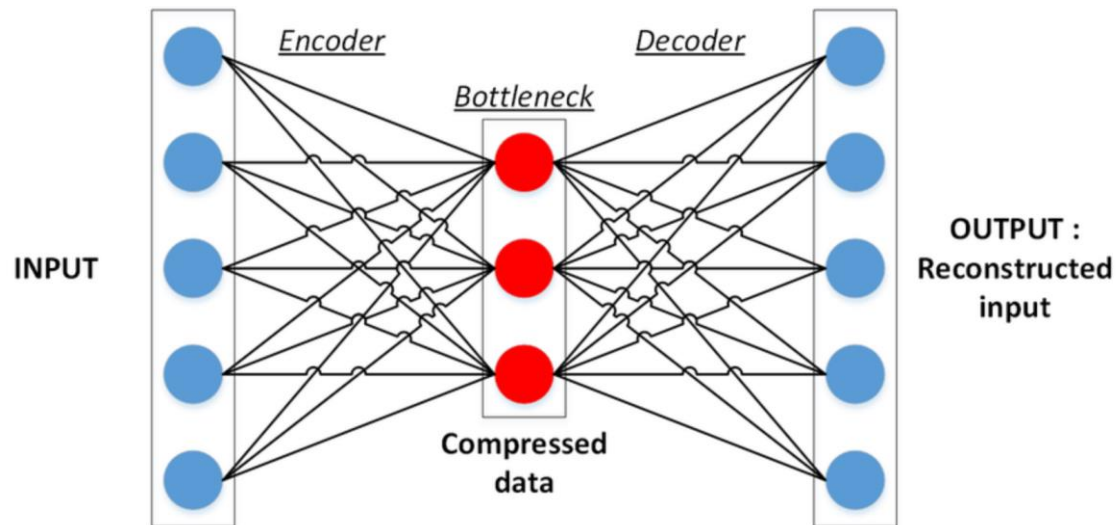


AE & GAN

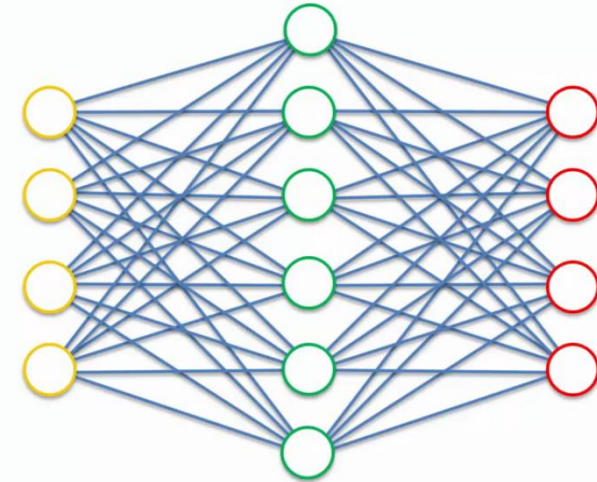
What is AE

- **AE(auto-encoder)**

✓ input을 넣었을 때 input과 똑같은 output을 출력하도록 만든 모델(신경망)



Overcomplete Hidden Layers



Machine Learning A-Z

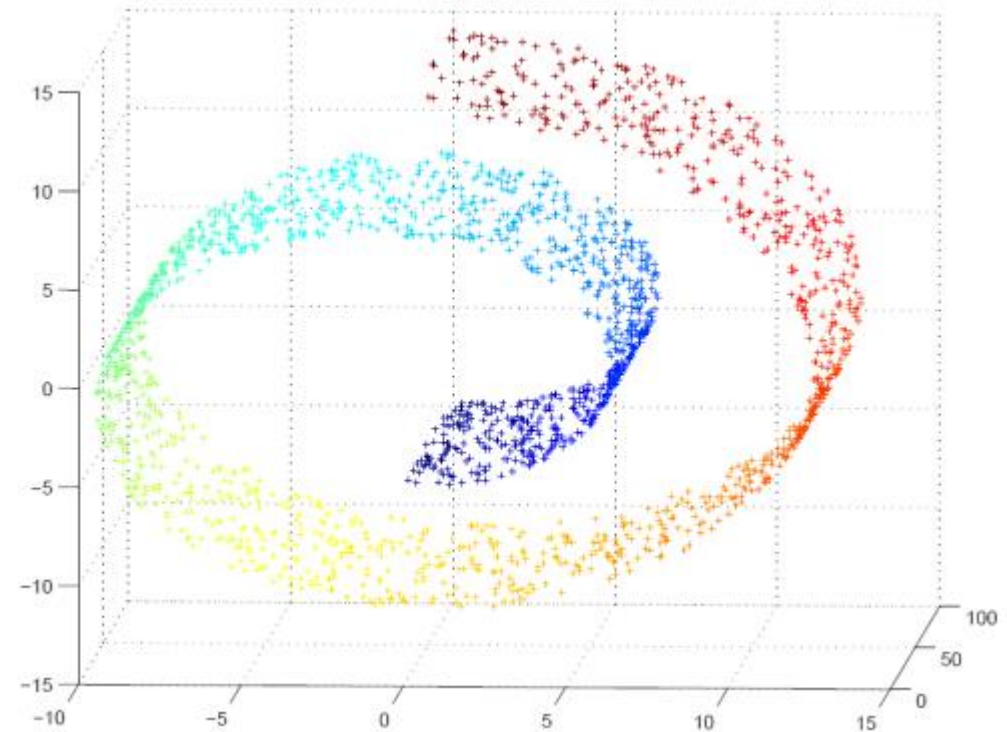
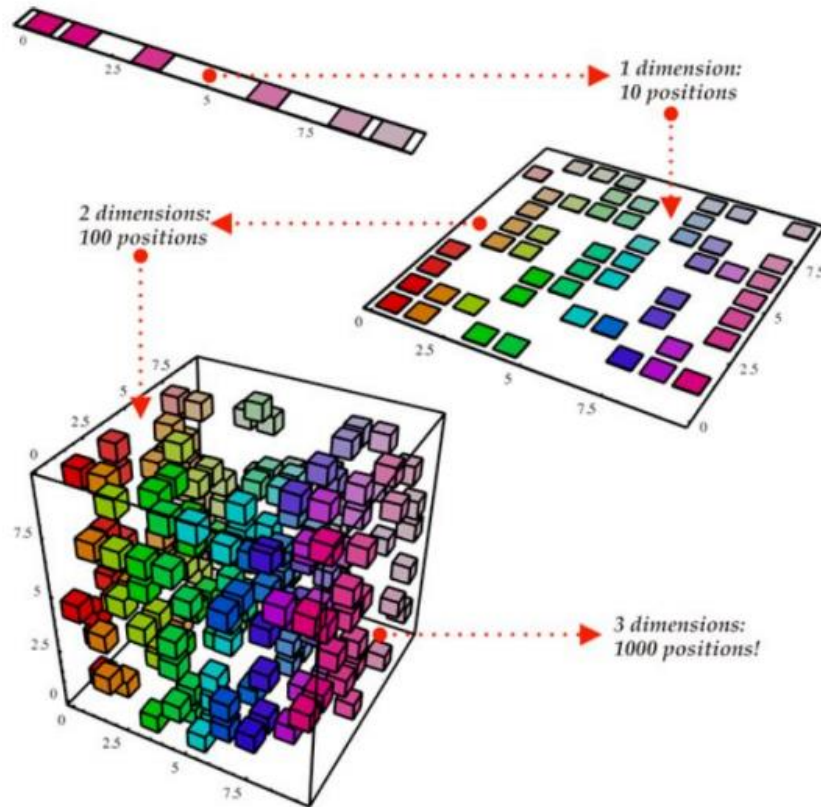
© SuperDataScience

AE & GAN

What is AE

- **Manifold Hypothesis**

- ✓ 고차원의 데이터의 밀도는 낮지만, 이들의 집합을 포함하는 저차원의 manifold(영역)이 존재한다. (차원이 높아질수록 밀도가 급격히 낮아진다.)

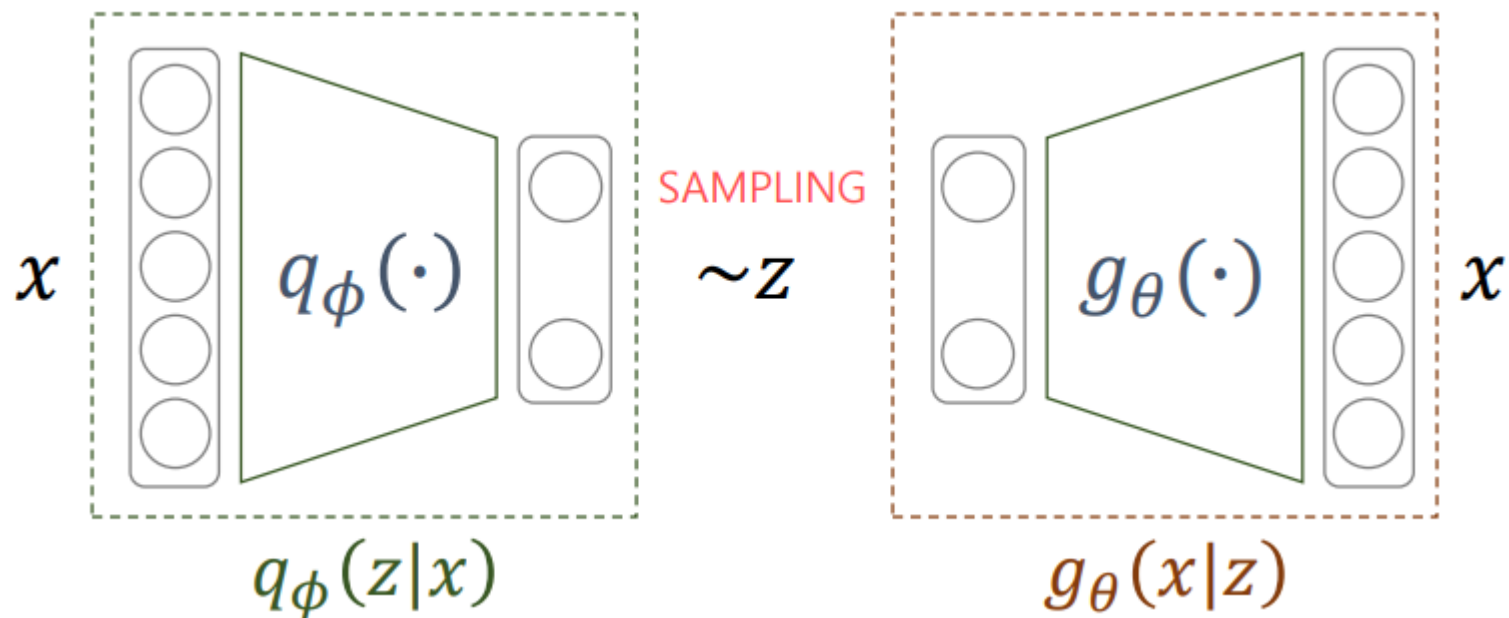


AE & GAN

VAE(variational autoencoder)

- **VAE(variational autoencoder)**

- ✓ 데이터의 분포를 다루기 쉬운 확률 분포로 정규화 하여 해당 분포로부터 input과 비슷한 다른 output을 출력(generative model)



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

- **VAE(variational autoencoder)**

- ✓ 데이터의 분포를 다루기 쉬운 확률 분포로 정규화 하여 해당 분포로부터 input과 비슷한 다른 output을 출력(generative model)

LOSS FUNCTION

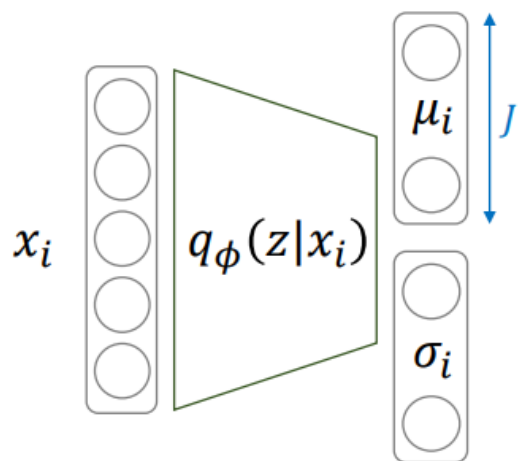
Regularization

VAE

13 / 49

KL divergence

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



$$\begin{aligned} KL(q_\phi(z|x_i)||p(z)) &= \frac{1}{2} \left\{ \text{tr}(\sigma_i^2 I) + \mu_i^T \mu_i - J + \ln \frac{1}{\prod_{j=1}^J \sigma_{i,j}^2} \right\} \\ &= \frac{1}{2} \left\{ \sum_{j=1}^J \sigma_{i,j}^2 + \sum_{j=1}^J \mu_{i,j}^2 - J - \sum_{j=1}^J \ln(\sigma_{i,j}^2) \right\} \\ &= \frac{1}{2} \sum_{j=1}^J (\mu_{i,j}^2 + \sigma_{i,j}^2 - \ln(\sigma_{i,j}^2) - 1) \quad \text{Easy to compute!} \end{aligned}$$

Kullback-Leibler divergence [edit]

The Kullback-Leibler divergence from $\mathcal{N}_0(\mu_0, \Sigma_0)$ to $\mathcal{N}_1(\mu_1, \Sigma_1)$, for non-singular matrices Σ_0 and Σ_1 , is:^[8]

$$D_{KL}(\mathcal{N}_0||\mathcal{N}_1) = \frac{1}{2} \left\{ \text{tr}(\Sigma_1^{-1} \Sigma_0) + (\mu_1 - \mu_0)^T \Sigma_1^{-1} (\mu_1 - \mu_0) - k + \ln \frac{|\Sigma_1|}{|\Sigma_0|} \right\},$$

where k is the dimension of the vector space.

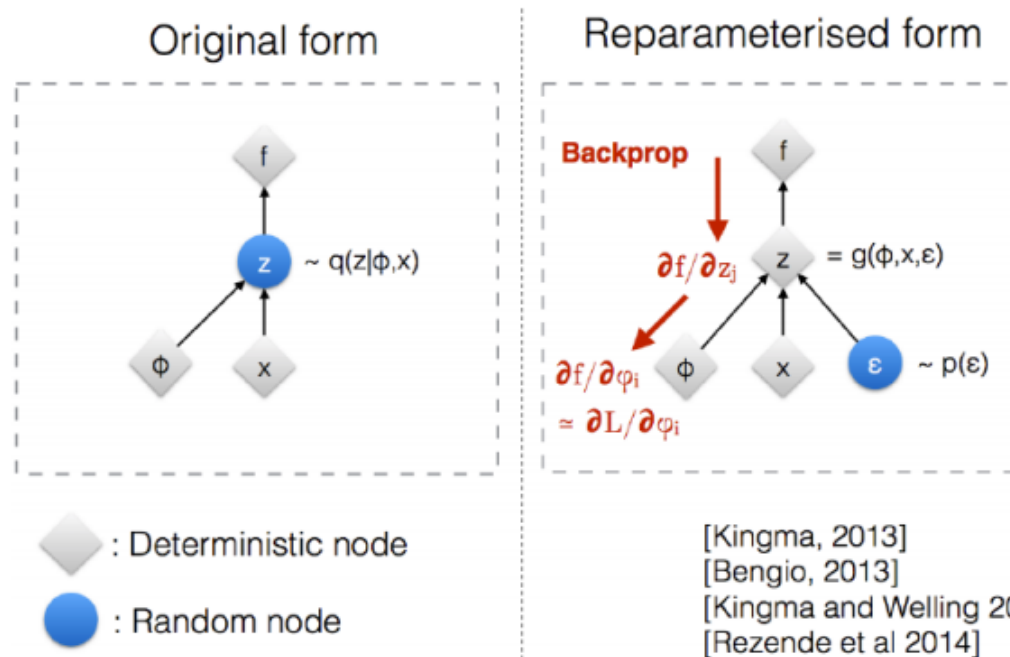
AE & GAN

VAE(variational autoencoder)

- **VAE(variational autoencoder)**

- ✓ 데이터의 분포를 다루기 쉬운 확률 분포로 정규화 하여 해당 분포로부터 input과 비슷한 다른 output을 출력(generative model)

Reparameterization Trick



Sampling Process

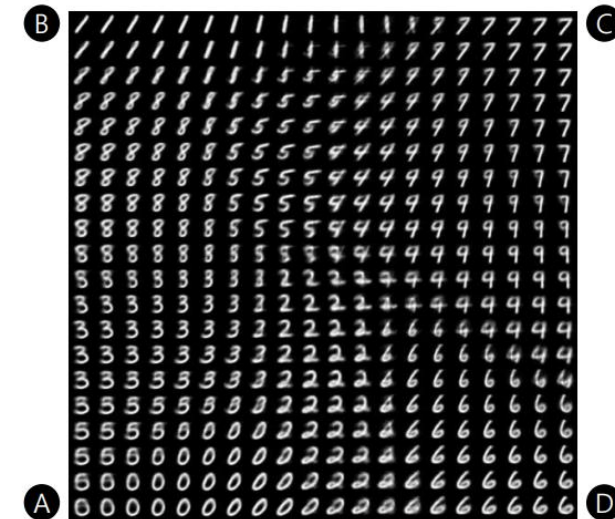
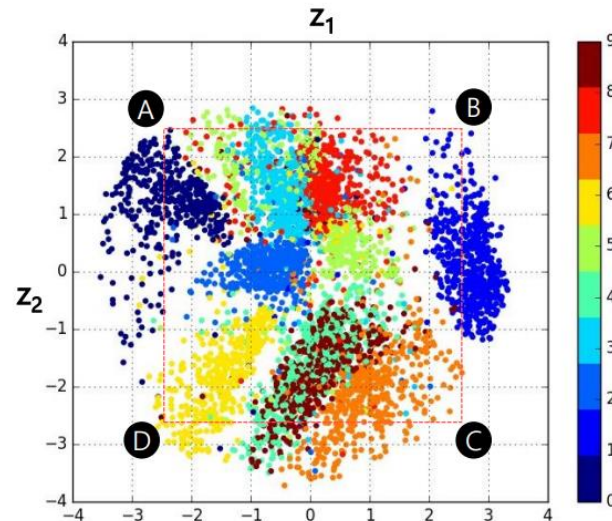
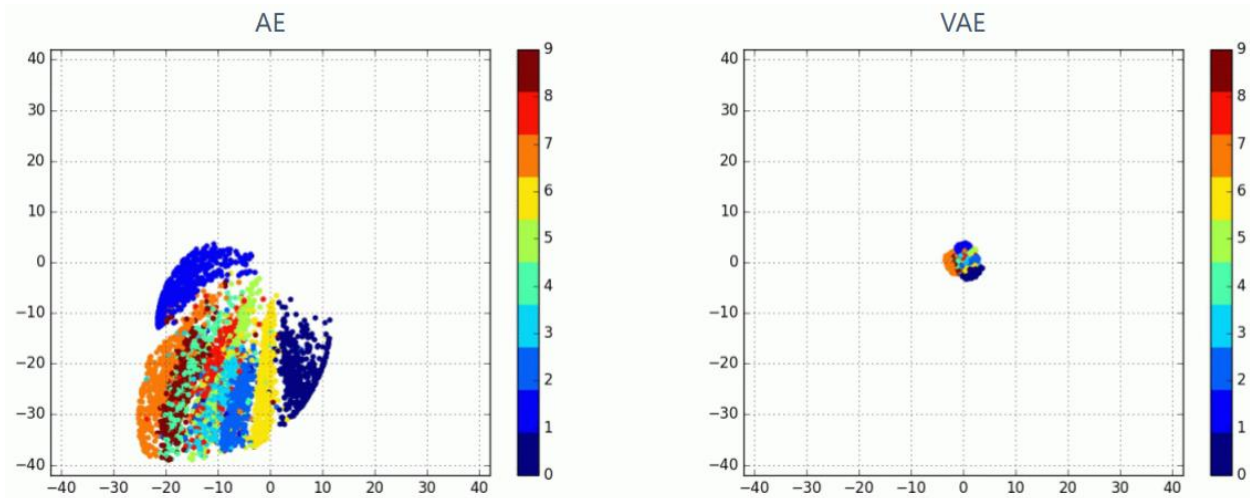
$$z^{i,l} \sim N(\mu_i, \sigma_i^2 I) \quad \Rightarrow \quad z^{i,l} = \mu_i + \sigma_i^2 \odot \epsilon$$
$$\epsilon \sim N(0, I)$$

AE & GAN

VAE(variational autoencoder)

- **VAE(variational autoencoder)**

- ✓ 데이터의 분포를 다루기 쉬운 확률 분포로 정규화 하여 해당 분포로부터 input과 비슷한 다른 output을 출력(generative model)

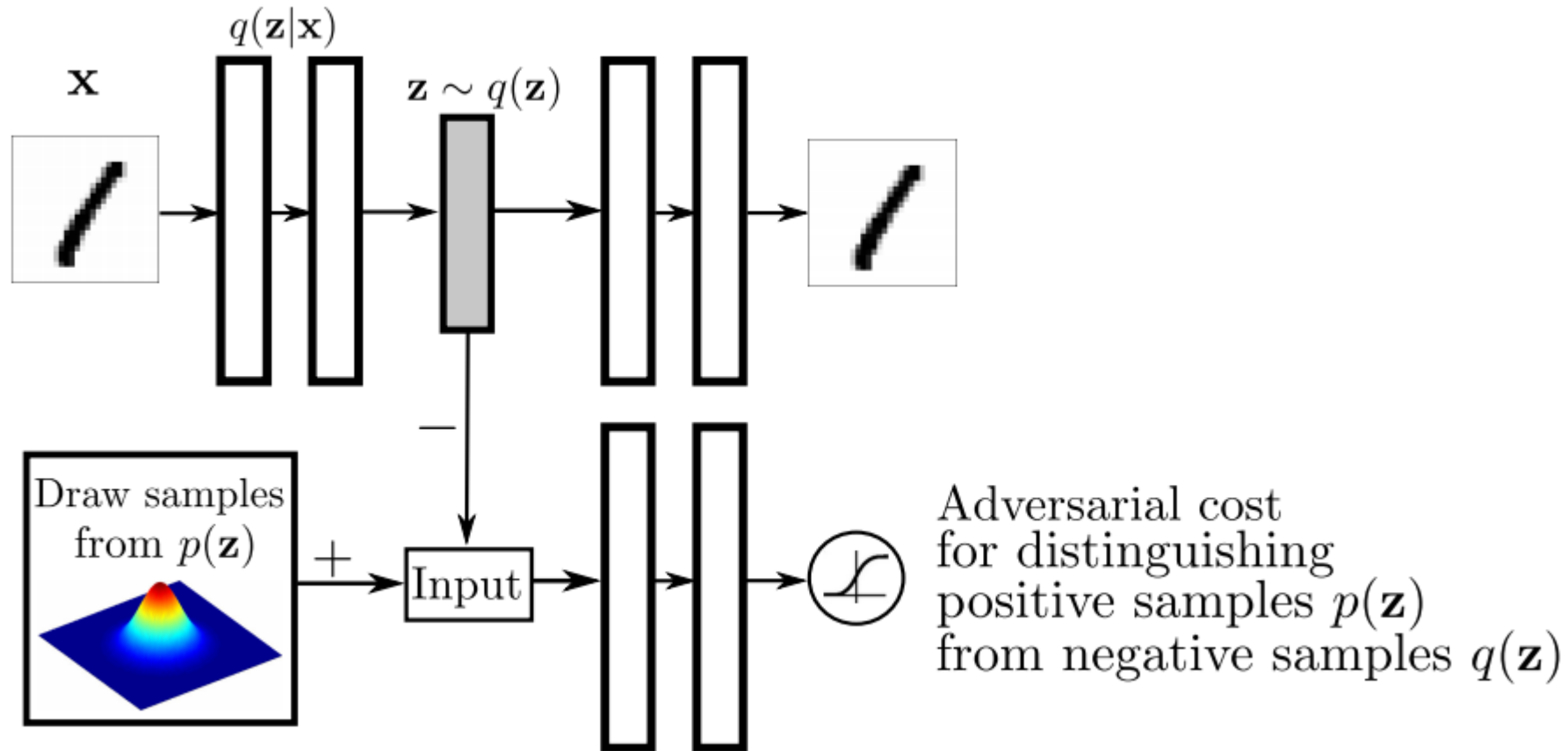


AE & GAN

VAE(variational autoencoder)

- **AAE(adversarial autoencoder)**

- ✓ 데이터의 분포를 원하는 확률 분포로 정규화 하여 해당 분포로부터 input과 비슷한 다른 output을 출력(generative model)



- **AAE(adversarial autoencoder)**

- ✓ 데이터의 분포를 원하는 확률 분포로 정규화 하여 해당 분포로부터 input과 비슷한 다른 output을 출력(generative model)

Loss Function

GAN loss

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

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

$$V_i(\phi, \lambda, x_i, z_i) = \log d_\lambda(z_i) + \log (1 - d_\lambda(q_\phi(x_i)))$$

*논문에는 로스 정의가 제시되어 있지 않아 새로 정리한 내용

VAE loss

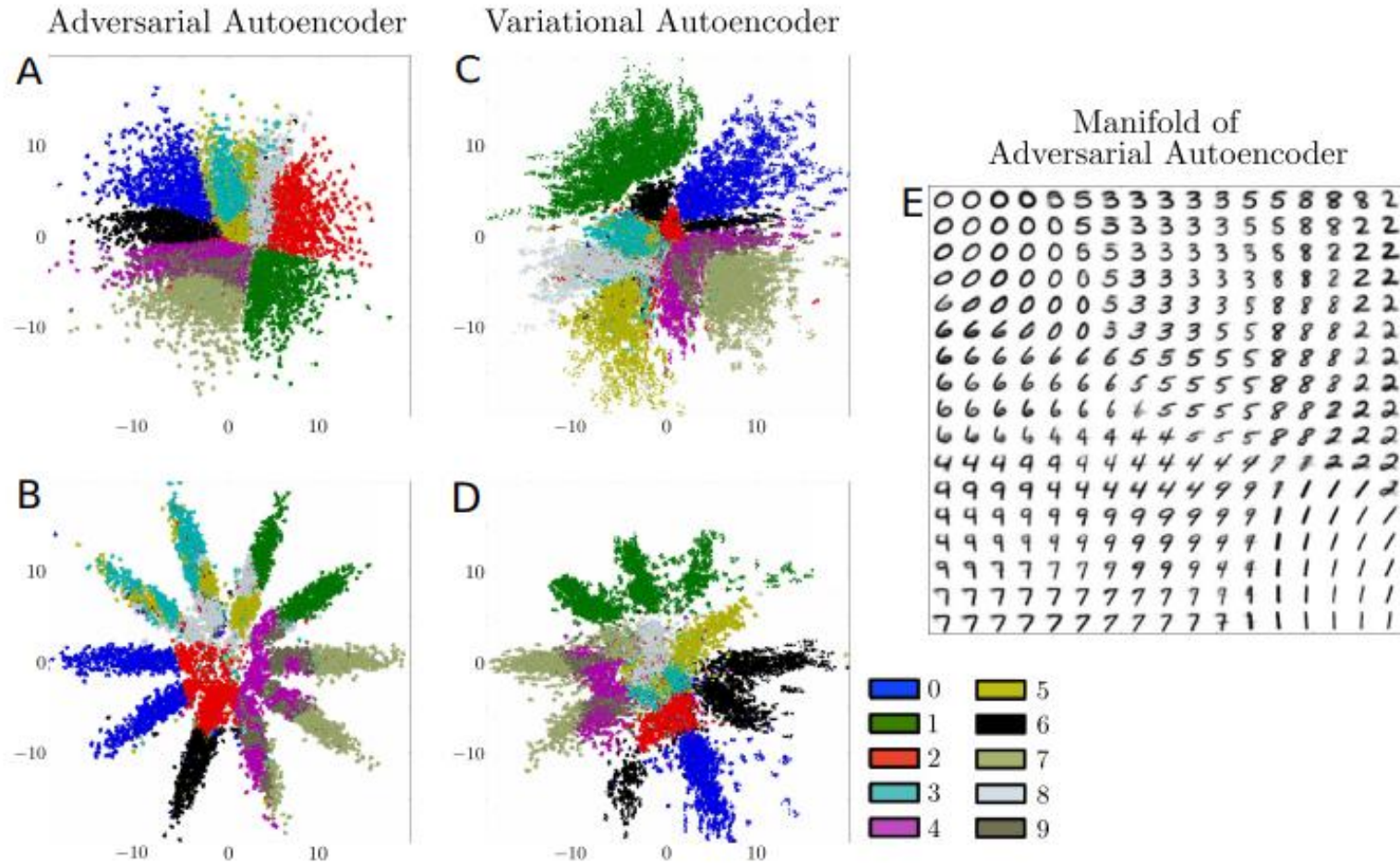
$$L_i(\phi, \theta, x_i) = -\mathbb{E}_{q_\phi(z|x_i)} [\log(p_\theta(x_i|z))] + \cancel{KL(q_\phi(z|x_i)||p(z))}$$

AE & GAN

VAE(variational autoencoder)

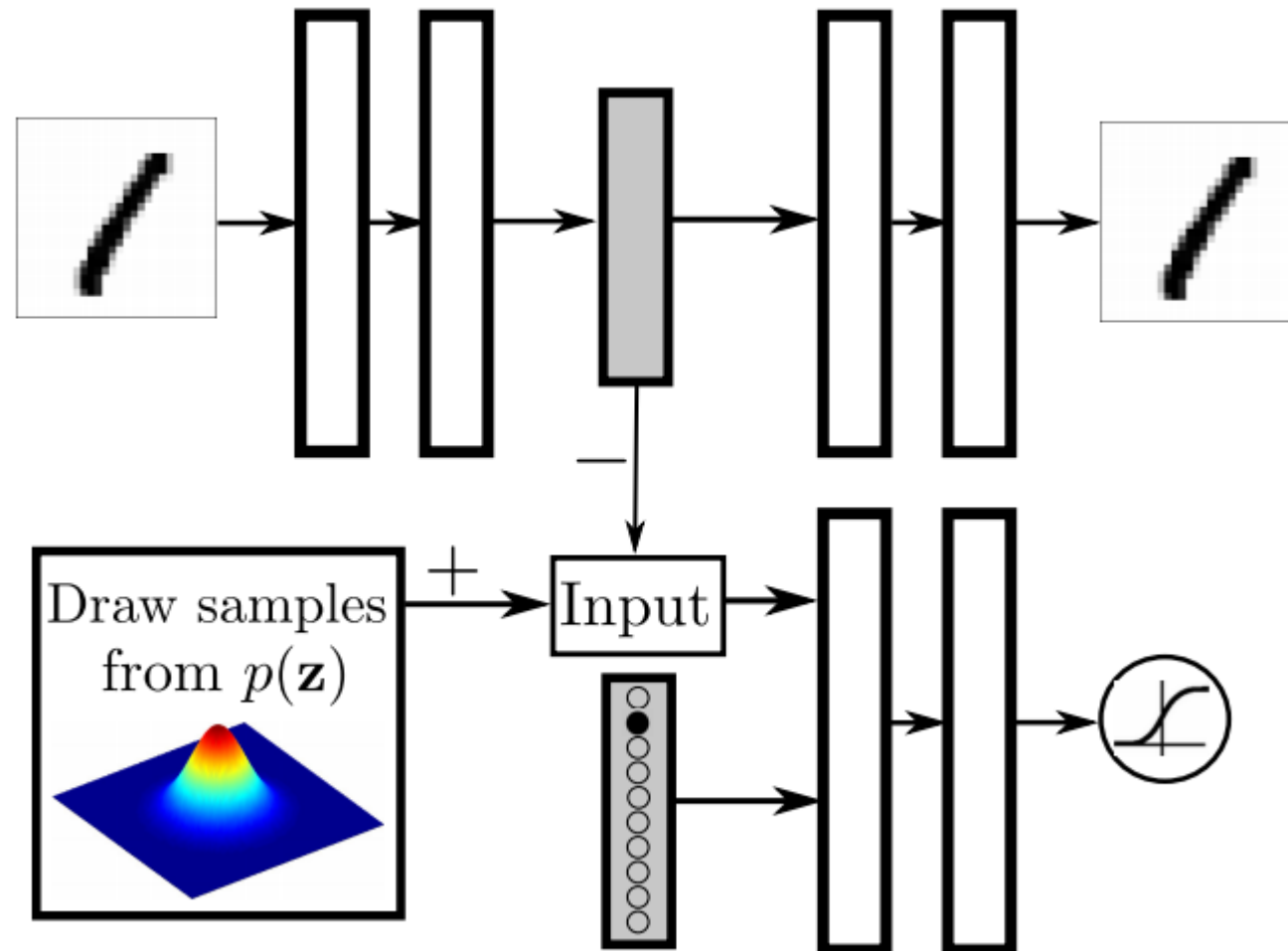
- **AAE(adversarial autoencoder)**

- ✓ 데이터의 분포를 원하는 확률 분포로 정규화 하여 해당 분포로부터 input과 비슷한 다른 output을 출력(generative model)



- **AAE(adversarial autoencoder) + condition**

- ✓ 데이터의 분포를 원하는 확률 분포로 정규화 하여 해당 분포로부터 input과 비슷한 다른 output을 출력(generative model)

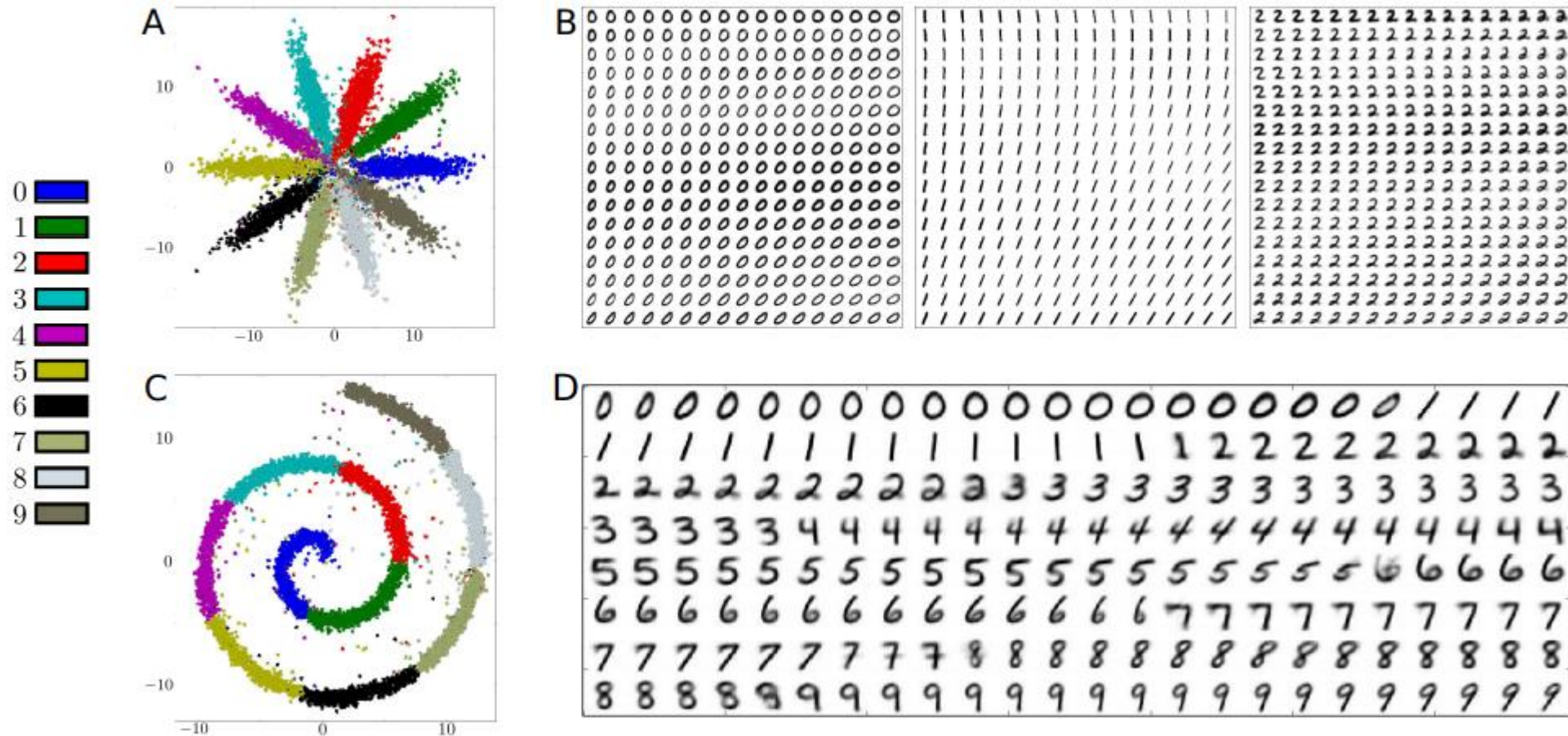


AE & GAN

VAE(variational autoencoder)

- **AAE(adversarial autoencoder)**

- ✓ 데이터의 분포를 원하는 확률 분포로 정규화 하여 해당 분포로부터 input과 비슷한 다른 output을 출력(generative model)



Reference

- Makhzani, Alireza, et al. "Adversarial autoencoders." *arXiv preprint arXiv:1511.05644* (2015).
- <https://www.youtube.com/watch?v=rNh2CrTFpm4>(오토인코더의 모든것)
- https://www.slideshare.net/NaverEngineering/ss-96581209?from_action=save(오토인코더의 모든것)

End

Appendix

