
Generative Models

Introduction

Example 1

자율주행 자동차 / 로봇
시뮬레이션:

트랙 주행은 할 수 있지만,
만, 실제 공도 환경을 재현하기는 어려움.

Example 2

파생상품 시장:

파생상품의 가격 변화를 시뮬레이션해
금융사가 판매할 수 있도록 검증(계리)할
수 있음.

Example 3

이미지 복원/생성:

손상되거나 화질이 낮은 이미지를 복원하거나, 원하는 이미지를 기준 방법보다 훨씬 경제적으로 합성 처리할 수 있음.
(딥페이크)

<https://qiskit.org/documentation/finance/tutorials/index.html>

Grounding Knowledge

1

Likelihood, Cross Entropy, KL Divergence

2

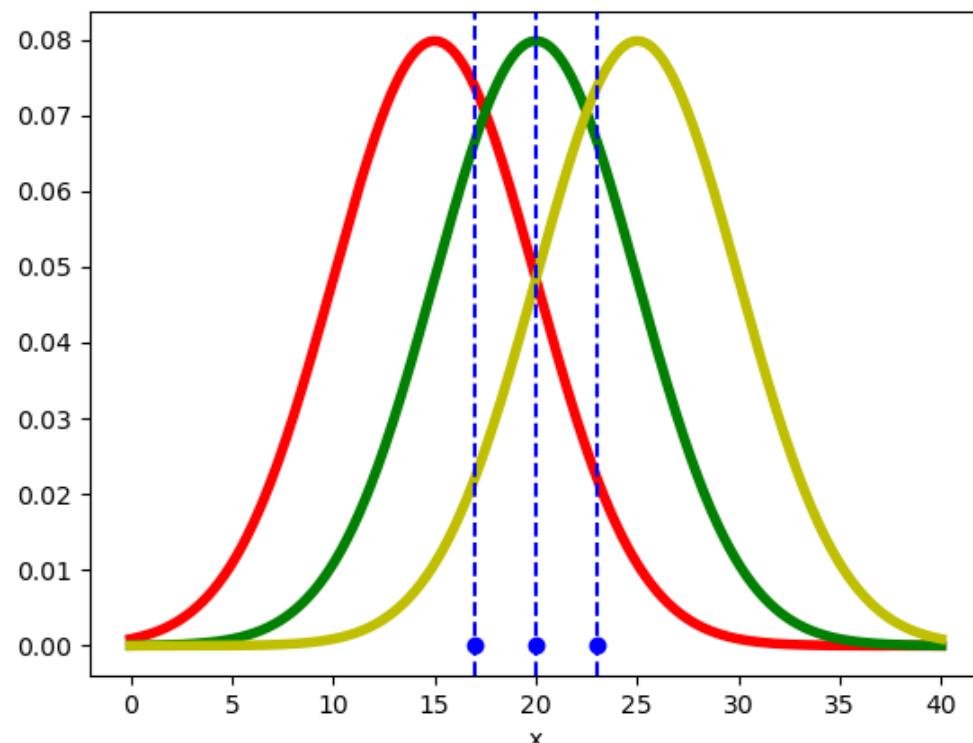
Density Estimation

3

Adversarial Game (for GAN)

Likelihood

가능도



A

연속 분포의 그래프에서:
pdf의 값 특정 사건이 일어날 가능성 ≠ 확률

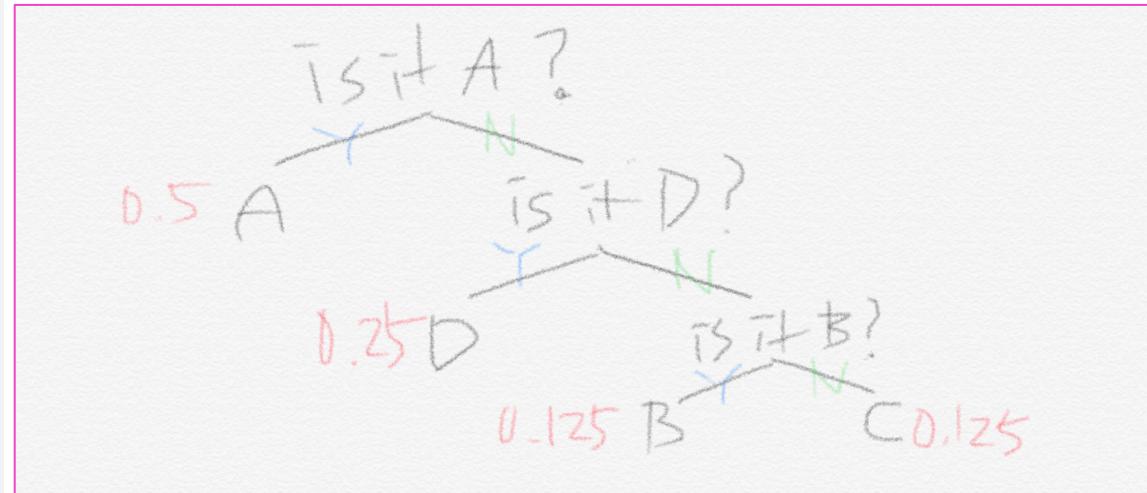
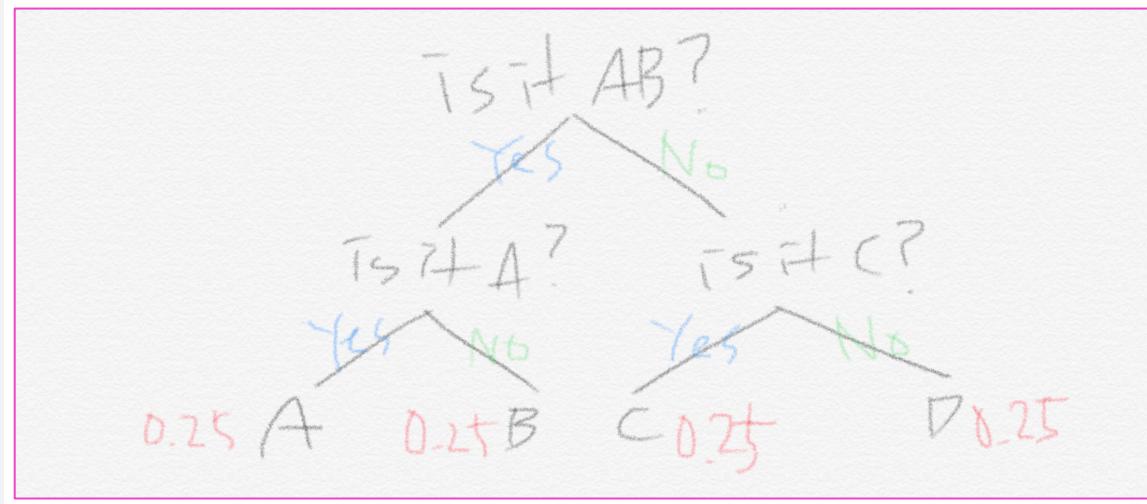
B

$\prod_i f(x_i; \theta)$:
확률변수의 어떤 파라미터에 대해, 표본 $\{x_i\}$ 가 도출
될 가능성.

C

MLE:
likelihood를 최대화하는 $\{x_i\}$ 를 이용한 통계량. 이때
iid로 추출된 표본일 경우, 계산의 편의성을 위해 log
likelihood를 대신 사용해도 무방.

Cross entropy



A

Entropy: 기계 1과 2의 출력 값을 예측할 때, 최적의 전략 하에서 필요한 질문 개수의 기댓값.

$$\sum_{i=1}^n p_i \log_2 \left(\frac{1}{p_i} \right) = - \sum_{i=1}^n p_i \log_2 p_i$$

B

Cross entropy: 기계 2의 출력 값을 예측할 때, 최적의 전략이 아닌 확률 분포 $\{q_i\}$ 로 표현된 특정 전략을 사용했을 때의 엔트로피 값.

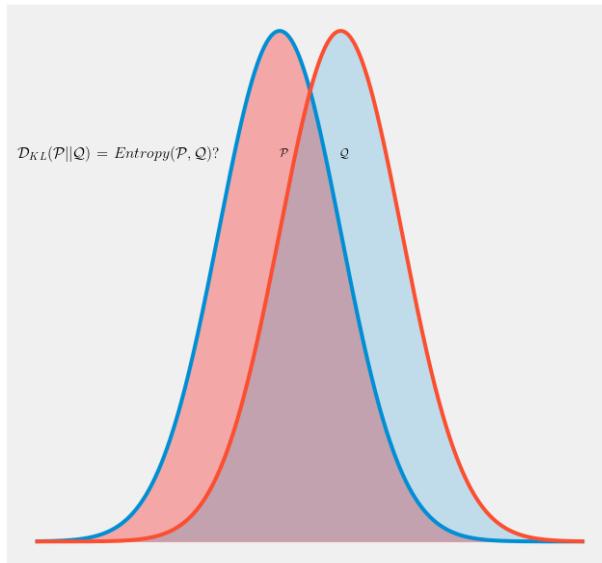
$$H(p, q) = - \sum_{i=1}^n p_i \log_2 q_i$$

C

Log loss function은 정확히 Cross entropy와 같고, 이것을 최소화하는 것은 negative log likelihood를 최소화하는 문제와 같음.

$$-y \log \hat{y} - (1 - y) \log(1 - \hat{y})$$

KL Divergence



$$KL(p \parallel q) = \begin{cases} \sum_i p_i \log \frac{p_i}{q_i} \text{ 또는 } -\sum_i p_i \log \frac{q_i}{p_i} & (\text{이산형}) \\ \int p(x) \log \frac{p(x)}{q(x)} dx \text{ 또는 } -\int p(x) \log \frac{q(x)}{p(x)} dx & (\text{연속형}) \end{cases}$$

A

두 확률 분포의 차이를 계산하는 데에 사용하는 함수. 어떤 이상적인 분포에 대해, 그 분포를 근사하는 다른 분포를 사용해 샘플링을 한다면 발생할 수 있는 정보 엔트로피 차이.

$$KL(p \parallel q) = H(p, q) - H(p)$$

B

1. $KL(p \parallel q) \geq 0$ by Jensen's inequality
 2. $KL(p \parallel q) \neq KL(q \parallel p)$
- * KL divergence는 거리 개념이 아니다.

$$JSD(p \parallel q) = \frac{1}{2}KL(p \parallel M) + \frac{1}{2}KL(q \parallel M)$$

C

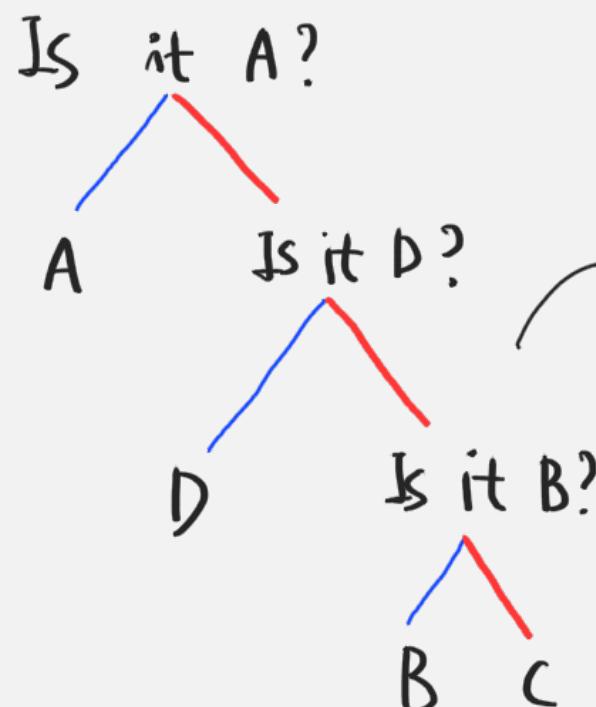
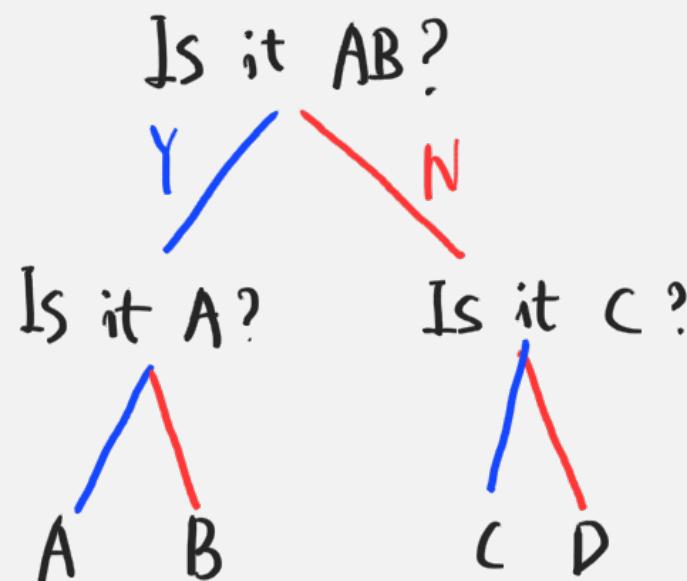
$$KL(p \parallel q) \cong \frac{1}{n} \sum_{i=1}^n \{-\ln q(x_i; \theta) + \ln p(x_i)\}$$

어떤 분포 $p(x)$ 를 $q(x)$ 를 이용해 근사시킬 때, KL divergence를 최소화시키는 것은 Likelihood를 최대화하는 것과 같음.

도화지

Machine 1			
A	B	C	D
$\frac{1}{4}$	$\frac{1}{4}$	$\frac{1}{4}$	$\frac{1}{4}$

Machine 2			
A	B	C	D
$\frac{1}{2}$	$\frac{1}{8}$	$\frac{1}{8}$	$\frac{1}{4}$



필요한 질문수

$$: \log_2 (\text{가능한 경우의 수})$$

$$= \log_2 \left(\frac{1}{p_i} \right)$$

$$\log_2 \left(\frac{1}{1/2} \right), \log_2 \left(\frac{1}{1/4} \right), \log_2 \left(\frac{1}{1/8} \right) = 1, 2, 3$$

$$\text{Entropy}_{H(P)} : \sum P_i \log \left(\frac{1}{P_i} \right)$$

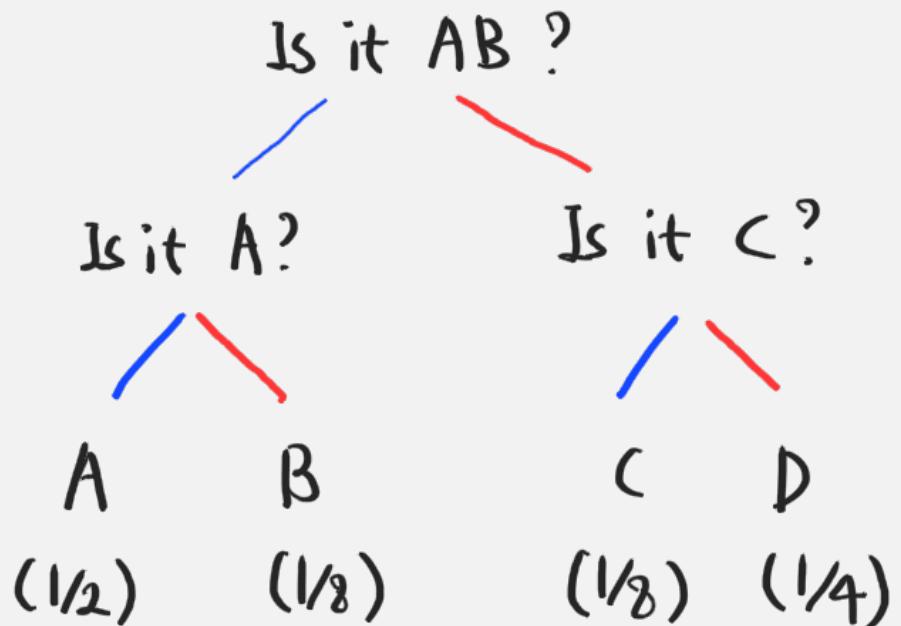
$$= 4 \cdot 0.25 \cdot 2 = 2$$

$$0.5 \cdot 1 + 0.25 \cdot 2 + 0.125 \cdot 3$$

$$= 1.75$$

Machine 2

A	B	C	D
$\frac{1}{2}$	$\frac{1}{8}$	$\frac{1}{8}$	$\frac{1}{4}$



$$P_i = \left[\frac{1}{2} \quad \frac{1}{8} \quad \frac{1}{8} \quad \frac{1}{4} \right]$$

$$q_i = \left[\frac{1}{4} \quad \frac{1}{4} \quad \frac{1}{4} \quad \frac{1}{4} \right]$$

전략

$$\sum P_i \log \left(\frac{1}{q_i} \right) = - \sum P_i \log q_i = 2 = H(P, q)$$

전략 q_i 가 최적 전략 P_i 와 비슷해질수록 작아짐.

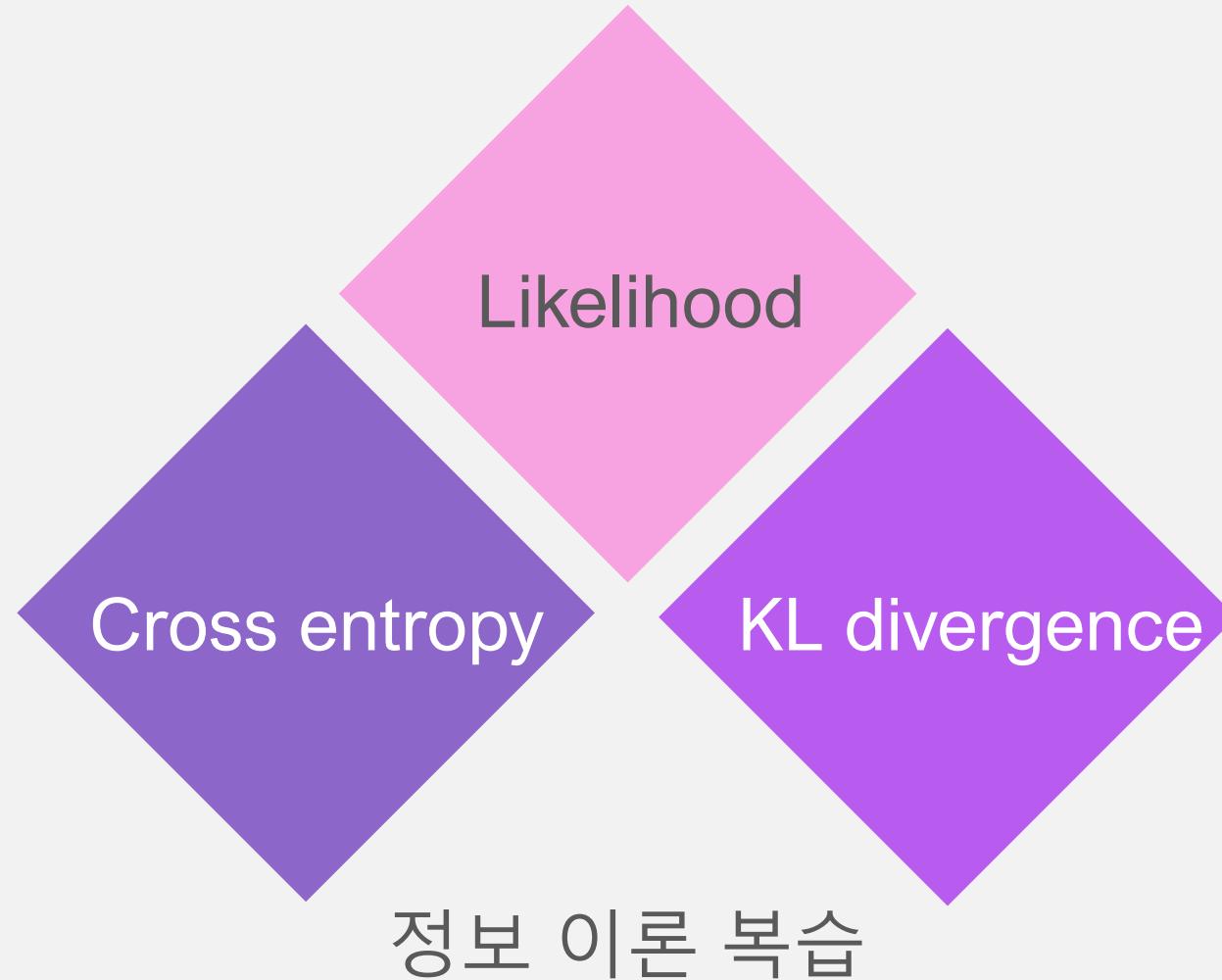
\equiv Likelihood 최대화 \equiv negative log likelihood 최소화

도화지

$$\begin{aligned} H(P, Q) &= -\sum P_i \log Q_i \\ &= -\sum P_i \log Q_i + \sum P_i \log P_i - \underbrace{\sum P_i \log P_i}_{H(P)} \\ &= \underbrace{\sum P_i \log \frac{P_i}{Q_i}}_{\text{'정보량 차이'}} + \underbrace{H(P)}_{P의 엔트로피} \end{aligned}$$

$$KL(P||Q) = H(P, Q) - H(P)$$

Summary



1. Cross entropy는 negative log likelihood와 같다. 그래서 cross entropy를 최소화하는 것이 log likelihood를 최대화하는 것과 같다.
2. 확률분포 p, q 에 대한 cross entropy는 $H(p) + KL(p||q)$ 와 같으므로 KL divergence를 최소화하는 것 또한 결국 log likelihood를 최대화하는 것과 같다.

Density estimation

분포 추정

A

Parametric case

추정하고자 하는 분포 p_{data} 와 가장 유사한 p_{model} 을 도출할 수 있는 파라미터 θ 는?

$$\theta^* = \operatorname{argmax}_{\theta} E_{x \sim p_{data}} [\ln p_{model}(x|\theta)]$$

B

Non-parametric case

간단한 분포 (uniform, gaussian) 등에 비선형 변환 G 를 적용해 p_{data} 에 근사한다.

$p_{code}(z)$: unif[0,1] or $N(0,1)$ 등. 확률 변수 Z
 $Z = \text{latent space}$: Z 의 확률 공간

$p_{data}(x)$: 근사하고자 하는 분포. 확률 변수 X
 $X = \text{data space}$: X 의 확률 공간

$G: \mathcal{Z} \rightarrow \mathcal{X}$ s.t. $G(z) = x$: 비선형 변환

Density estimation

분포 추정

A

Single variate case
:example

Uniform[0,1]을 코드 분포로 사용한다면?

B

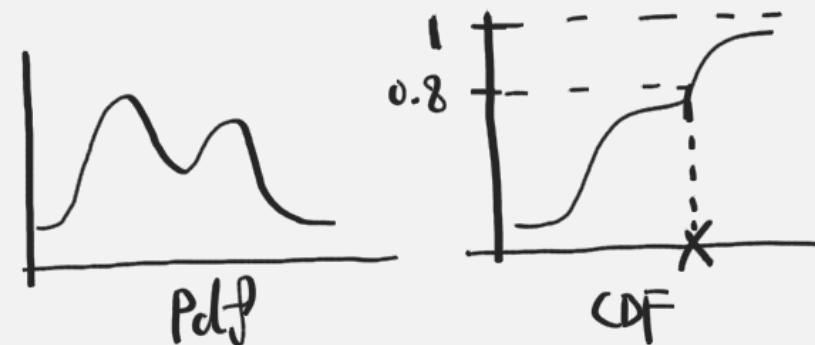
Multivariate case
:example

다변량 정규분포 $\mathcal{N}(\mu, \Sigma)$ 를 근사해 보자.

도화지

$P_{code} : \text{Unif}[0,1]$ $Z : [0,1]$

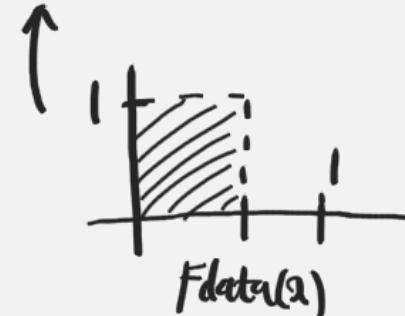
P_{model} , X : 일의의 일변수 확률분포



$G : F^{-1}_{\text{data}} : [0,1] \rightarrow \mathbb{R}$

Assuming $G(Z) = X$, we gain

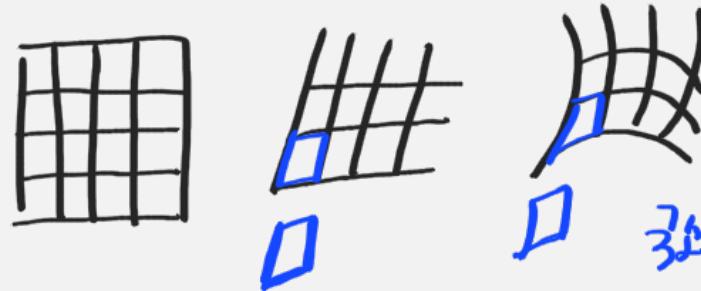
$$F_{G(Z)}(x) = P(G(Z) \leq x) = P(Z \leq F_{\text{data}}(x)) = F_{\text{data}}(x)$$



도화지

$$P_{\text{model}}(x) = P_x(x) = \frac{P_{\text{code}}(G^{-1}(x))}{|\det J_G(x)|}$$

비선형 변환을 선형 변환으로 고사.



ex) $P_{\text{model}} = \mathcal{N}(\mu, \Sigma) : x$

$$P_{\text{code}} = \mathcal{N}(0, I_n) : z$$

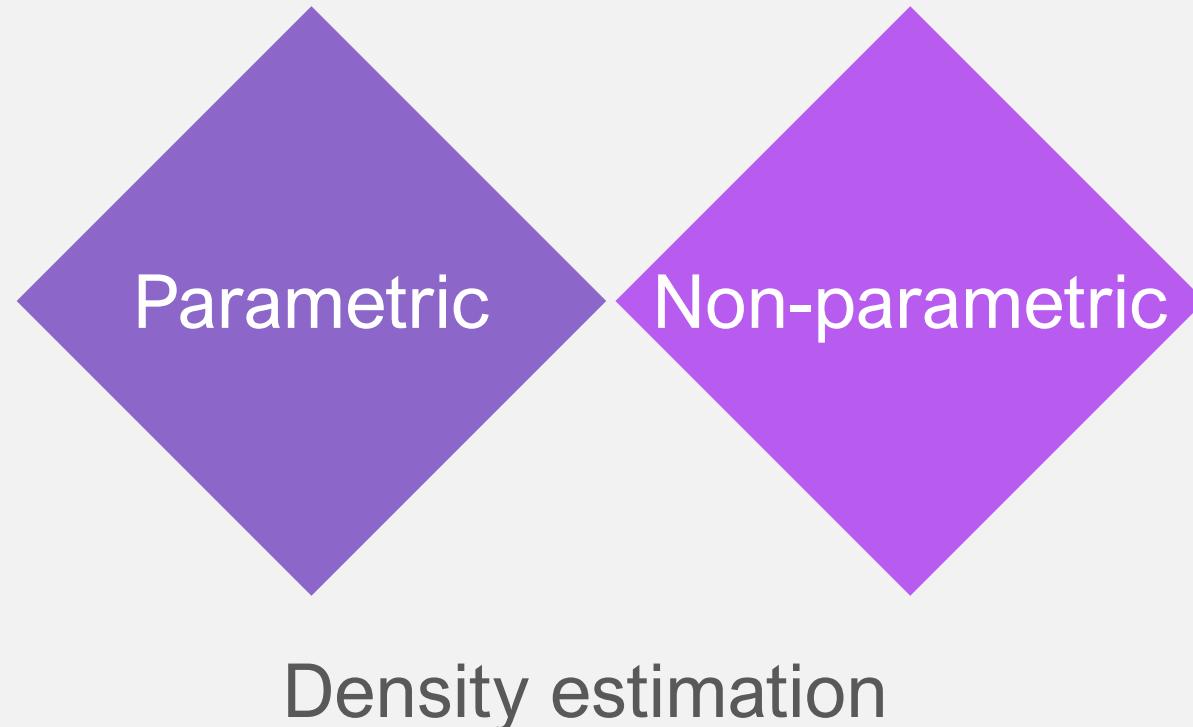
$$G: \mathbb{R}^n \rightarrow \mathbb{R}^n = \mu + Ax, \text{ where } \Sigma = AA^T \quad \text{Then } G^{-1}(x) = A^{-1}(x - \mu), \quad J = \left| \frac{\partial G}{\partial x} \right| = |\det(A)|$$

$$P_{\text{code}}(z) = \frac{1}{(2\pi)^{2n}} \exp\left(-\frac{1}{2} \|z\|^2\right)$$

$$P_{\text{model}}(x) = \frac{1}{|\det(A)|} \frac{1}{(2\pi)^{2n}} \exp\left(-\frac{1}{2} \|A^{-1}(x - \mu)\|^2\right) \Rightarrow \frac{1}{\sqrt{(2\pi)^n \det \Sigma}} \exp\left(-\frac{1}{2} (x - \mu)^T \Sigma^{-1} (x - \mu)\right)$$

교재 595p.

Summary + additional information



1. Parametric case는 MLE 찾기로 해석할 수 있다.
2. Non-parametric case에서 실제로 중요한 문제는 적절한 선형 변환 G 를 찾는 문제이다. Latent space의 차원을 최대한 작게 해야 필요한 파라메터의 개수가 줄어들기 때문이다.
3. P model이 parametric/non-parametric 여부로 문제를 분류할 수도 있지만, 분포의 형태를 어느 정도 알고 있느냐/알 수 없느냐 여부에 따라서 분류할 수도 있다.

Adversarial games

적대적인 게임

001

두 명의 경쟁자들이,
한 명은 효용 함수를
최소화하고자, 다른
한 명은 최대화하고
자 하는 목적을 갖고
각각 parameter 1,
parameter 2를 조정.

$$(x^*, y^*) = \arg \max_y \min_x V(x, y)$$

002

Gradient
descent로 균형을
찾는다고 할 때,
시간 변수 t에 대해
다음 식으로 나타
낼 수 있음.

$$x(t + \Delta t) = x(t) - \Delta t \frac{\partial V}{\partial x}$$

$$y(t + \Delta t) = y(t) + \Delta t \frac{\partial V}{\partial y}$$

003

$\frac{dx}{dt} = -\frac{\partial V}{\partial x}$,
 $\frac{dy}{dt} = \frac{\partial V}{\partial y}$ 이고, 균형에
서 x, y 가 변할 유인이
없으므로 도함수 값이
0이어야 하고, 이는
그래디언트 값이 0인
것과 동치.

$$(x^*, y^*) = \lim_{t \rightarrow \infty} (x(t), y(t))$$

in the optimal point,

$$\frac{\partial V}{\partial x} = 0, \quad \frac{\partial V}{\partial y} = 0$$

도화지

$$V(x, y) = \frac{1}{2}(x^2 - y^2) \quad \max_y \min_x V(x, y)$$

$$\begin{aligned} \frac{dx}{dt} = -\frac{\partial V}{\partial x} &= -x & \rightsquigarrow \frac{c}{e^t} \begin{cases} x(t) = x_0 e^{-t} \\ y(t) = y_0 e^{-t} \end{cases} & (x^*, y^*) = \lim_{t \rightarrow \infty} (x(t), y(t)) \\ \frac{dy}{dt} = \frac{\partial V}{\partial y} &= -y & & = (0, 0) \end{aligned}$$

discrete case

$$\begin{aligned} x_{n+1} &= x_n - \eta x_n \\ y_{n+1} &= y_n - \eta y_n \end{aligned} \quad \left\{ \begin{array}{l} x_n = (1-\eta)^n x_0 \\ y_n = (1-\eta)^n y_0 \end{array} \right. \rightsquigarrow (0, 0)$$



Let's move on!

Variational AutoEncoder

Jeong Yujin

ESC, YONSEI UNIVERSITY

April 6, 2022

Table of Contents

1 Introduction

- Generative model
- Latent variable

2 VAE

- Variational Autoencoder
- Lower bound
- Reparameterization trick
- Loss function
- code example

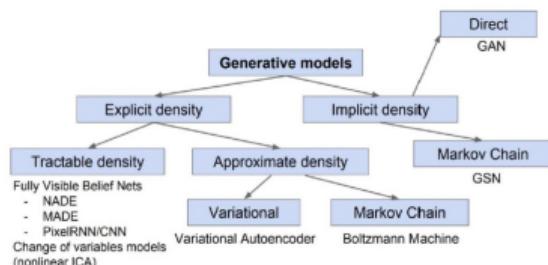
Generative model

Generative model은 주어진 학습데이터를 학습하여 학습 데이터의 분포를 따르는 유사한 데이터를 생성해내는 model이다.

대표적인 model:

- GAN
- VAE

VAE explicitly define and solve for $p_{model}(x)$ and define intractable density function with latent variable z .

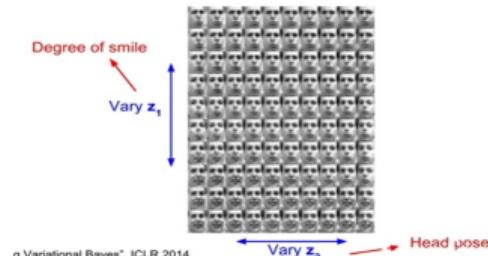


Latent variable

- Latent variable is a random variable that can not be observed. However, the presence of latent variables can be detected by their effects on variables that are observable.



- In VAE, x is a data and z is latent factors used to generate x .
e.g. x is a face image and z can be 'how much of smile is on the face', 'position of the eyebrow', 'orientation of head', ...

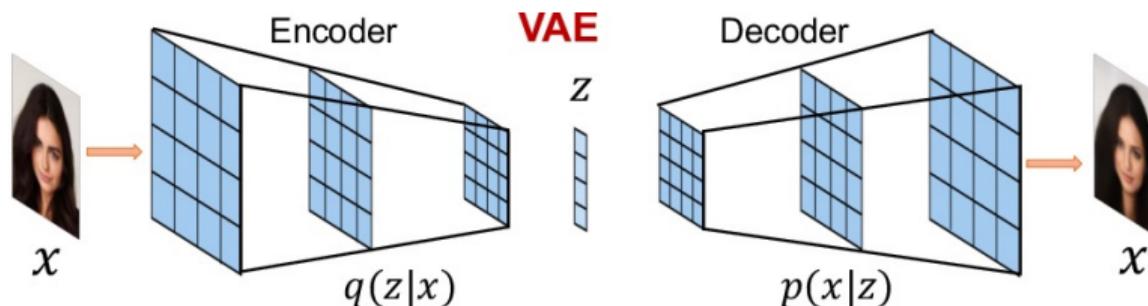


g "Variational Bayes", ICLR 2014

II. VAE

Variational Auto Encoder

- Variational AutoEncoder는 generative model의 한 종류로, variational bayes와 graphical model에서 영감을 받았다.
- Using **Stochastic Gradient descent method**, VAE learns the parameters ϕ, θ
- VAE approximates **encoder network** $q_\phi(z|x)$ and learns **decoder network** $p_\theta(x|z)$ by using **Neural Network**.
-



VAE

For the case of i.i.d dataset and continuous latent variables per datapoint, the **Auto Encoding VB(AEVB)** algorithm is proposed.

In the AEVB algorithm,

- We make inference and learning efficient by using the SGVB estimator.
- The learned approximate posterior inference model can also be used for recognition, denoising, representation and visualization purposes.
- When a NN is used for the recognition model, we arrive at **variational auto encoder**.

Approximated posterior

- **Variational Inference:**

Approximate the target distribution $p(x)$ by using a tractable density $q(x)$

$$q^* = \underset{q}{\operatorname{argmin}} KL(q||p)$$

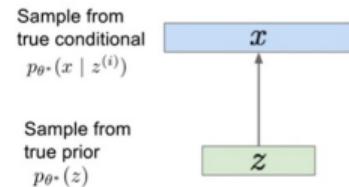
Problem scenario

$\mathbf{X} = \{\mathbf{x}^{(i)}\}_{i=1}^N$, i.i.d samples of some continuous and discrete variable \mathbf{x} .

We assume that the data are generated by some random process, involving unobserved continuous random variable \mathbf{z} .

The random process consists of two steps:

- ① a value $\mathbf{z}^{(i)}$ is generated from some prior distribution $p_{\theta^*}(\mathbf{z})$
 $\mathbf{z}^{(i)} \sim p_{\theta^*}(\mathbf{z})$
- ② a value $\mathbf{x}^{(i)}$ is generated from some conditional distribution $p_{\theta^*}(\mathbf{x}|\mathbf{z})$
 $\mathbf{x}^{(i)} \sim p_{\theta^*}(\mathbf{x}|\mathbf{z})$



But a lot of this process is hidden from our view: the true parameter θ^* as well as the values of the latent variable $\mathbf{z}^{(i)}$ are unknown to us.

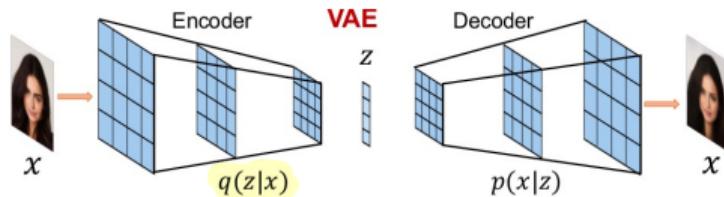
VAE

- In generative model, we want to know $p_\theta(x)$: marginal likelihood.
 - However, we can not compute the data likelihood
- $$p_\theta(x) = \int p_\theta(z)p_\theta(x|z)dz$$
- Posterior density $p_\theta(z|x)$ also intractable because of $p_\theta(x)$

$\therefore P_\theta(z|x) = \frac{P_\theta(x|z)p(z)}{P_\theta(x)}$ (by Bayes rule)

(true posterior)

- Then, define additional **encoder network** $q_\phi(z|x)$ that approximates $p_\theta(z|x)$



$q_\phi(z|x)$: **encoder network**
 $p_\theta(x|z)$: **decoder network**

$q_\phi(z|x) \sim \mathcal{Z}$: Gaussian

$p_\theta(x|z) \sim \mathcal{X}$

VAE

- 왜 $q_\phi(z|x)$ 를 정의했나?

Because we can derive a **lower bound** on the **data likelihood**.

- $\log p_\theta(x) \geq L(\theta, x)$

위와 같은 **lower bound(tractable)**를 maximize함으로써 log likelihood를 최대화할 수 있기 때문이다.(maximum likelihood)

VAE

- The marginal likelihood:

$$\log p_{\theta}(\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(N)}) = \sum_{i=1}^N \log p_{\theta}(\mathbf{x}^{(i)})$$

- The marginal likelihoods of individual datapoints can be written as:

$$\begin{aligned} \log p_{\theta}(\mathbf{x}^{(i)}) &= E_{z \sim q_{\phi}(z|\mathbf{x}^{(i)})} \left[\log p_{\theta}(\mathbf{x}^{(i)}) \right] \\ &= E_{q_{\phi}} \left[\log \frac{p_{\theta}(\mathbf{x}^{(i)}|z)p_{\theta}(z)}{p_{\theta}(z|\mathbf{x}^{(i)})} \right] \quad \left(\because \frac{p_{\theta}(\mathbf{x}^{(i)}|z)p_{\theta}(z)}{p_{\theta}(z|\mathbf{x}^{(i)})} = \frac{\frac{p_{\theta}(\mathbf{x}^{(i)}, z)}{p_{\theta}(z)}}{\frac{p_{\theta}(z, \mathbf{x}^{(i)})}{p_{\theta}(\mathbf{x}^{(i)})}} = p_{\theta}(\mathbf{x}^{(i)}) \right) \\ &= E_{q_{\phi}} \left[\log \frac{p_{\theta}(\mathbf{x}^{(i)}|z)p_{\theta}(z)}{p_{\theta}(z|\mathbf{x}^{(i)})} * \frac{q_{\phi}(z|\mathbf{x}^{(i)})}{q_{\phi}(z|\mathbf{x}^{(i)})} \right] \end{aligned}$$

VAE

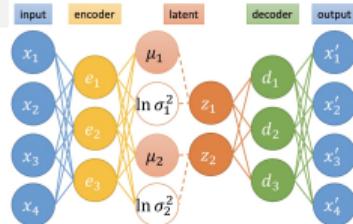
$$\begin{aligned}
 &= E_{q_{\phi}} \left[\log \frac{p_{\theta}(\mathbf{x}^{(i)}|z)p_{\theta}(z)}{p_{\theta}(z|\mathbf{x}^{(i)})} * \frac{q_{\phi}(z|\mathbf{x}^{(i)})}{q_{\phi}(z)} \right] \\
 &= E_{q_{\phi}(z)} \left[\log p_{\theta}(\mathbf{x}^{(i)}|z) \right] - E_{q_{\phi}(z)} \left[\log \frac{q_{\phi}(z|\mathbf{x}^{(i)})}{p_{\theta}(z)} \right] + E_{q_{\phi}(z)} \left[\log \frac{q_{\phi}(z|\mathbf{x}^{(i)})}{p_{\theta}(z|\mathbf{x}^{(i)})} \right] \\
 &= E_{q_{\phi}(z)} \left[\log p_{\theta}(\mathbf{x}^{(i)}|z) \right] - \underset{\textcircled{1}}{D_{KL}} \left(q_{\phi}(z|\mathbf{x}^{(i)}) || p_{\theta}(z) \right) + \underset{\textcircled{2}}{D_{KL}} \left(q_{\phi}(z|\mathbf{x}^{(i)}) || p_{\theta}(z|\mathbf{x}^{(i)}) \right) \\
 &\quad \underset{\textcircled{3}}{\underbrace{\geq 0}} \quad (p_{\theta}(z|\mathbf{x}^{(i)}) \text{ is intractable})
 \end{aligned}$$

- $D_{KL} (q_{\phi}(z|\mathbf{x}^{(i)}) || p_{\theta}(z|\mathbf{x}^{(i)})) \geq 0$ 이므로,
 $L(\theta, \phi; \mathbf{x}^{(i)}) = E_{q_{\phi}(z)} \left[\log p_{\theta}(\mathbf{x}^{(i)}|z) \right] - D_{KL} (q_{\phi}(z|\mathbf{x}^{(i)}) || p_{\theta}(z))$: maximize (tractable)
- We can compute estimate of the first term through sampling. ($z \sim q_{\phi}(z|\mathbf{x}^{(i)})$ 에서 뽑기, $E_{q_{\phi}(z)} [\log p_{\theta}(\mathbf{x}^{(i)}|z)]$ 에 넣어서 계산)
- The second term has nice closed form solution. ($q_{\phi}(z|\mathbf{x}^{(i)}) \sim \mathcal{N}(\mu_{\phi}, \text{diag}_{\phi}(D))$)
- In third term, $p_{\theta}(z|x)$ is intractable. but ≥ 0

VAE

Hence,

- Variational bound:



$$L(\theta, \phi; \mathbf{x}^{(i)}) = -D_{KL}(q_\phi(\mathbf{z}|\mathbf{x}^{(i)})||p_\theta(\mathbf{z})) + E_{q_\phi(\mathbf{z}|\mathbf{x}^{(i)})} \left[\log p_\theta(\mathbf{x}^{(i)}|\mathbf{z}) \right]$$

$$\therefore \log p_\theta(\mathbf{x}^{(i)}) \geq L(\theta, \phi; \mathbf{x}^{(i)})$$

Our goal is to maximize the lower bound $L(\theta, \phi; \mathbf{x}^{(i)})$

따라서 lower bound를 최대화하도록 θ 와 ϕ 를 optimize해야 한다.

$$(\theta^*, \phi^*) = \operatorname{argmax}_{\theta, \phi} L(\theta, \phi; \mathbf{x}^{(i)})$$

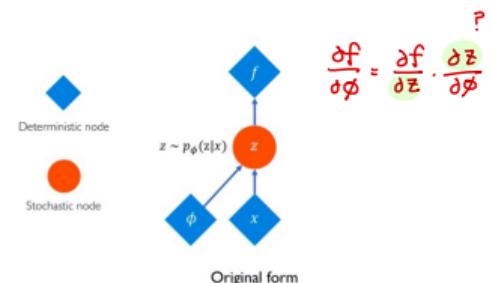
VAE

- How can we optimize ϕ and θ ?

Differentiate and optimize the lower bound w.r.t both the variational parameters ϕ , and generative parameters θ .

- In forward pass, We can compute this term by sampling $z \sim q_\phi(z|x)$. $E_{z \sim q_\phi(z|x^{(i)})} [\log p_\theta(x^{(i)}|z)]$.

- 그러나 back propagation은 불가능! sampling은 미분 가능하지 않다. Gradient descent method 사용 불가능.
Gradient descent method를 사용하려면 parameter에 대해 미분이 가능해야 하고, model은 deterministic해야 한다.(Stochasticity는 input에만 있어야 한다.)

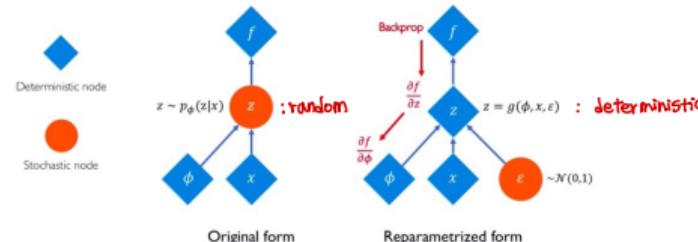


Reparameterization trick

e.g.

 \mathbf{t} : random $\mathbf{t} = \mathbf{a}\mathbf{x} + \mathbf{b}$: deterministic
random

: 같은 input에 대해서 같은 결과가 나옴.



z: 애 자체가 stochastic.

- 원래 $z \sim q_\phi(z|x)$: 애는 미분 불가능.
- 따라서 parameter에 대해 미분 가능하도록 reparameterize해줌.
 $\tilde{z} \sim g_\phi(\epsilon, x)$ with $\epsilon \sim p(\epsilon)$
: deterministic 해짐 random.
- Reparameterization trick을 사용하여 z를 deterministic하게 바꾸어주었지만, z는 $q_\phi(z|x)$ 에서 sampling 된 것과 마찬가지임.

Take, for example, the univariate Gaussian case: let $z \sim p(z|x) = \mathcal{N}(\mu, \sigma^2)$. In this case, a valid reparameterization is $z = \mu + \sigma\epsilon$, where ϵ is an auxiliary noise variable $\epsilon \sim \mathcal{N}(0, 1)$. Therefore, $\mathbb{E}_{\mathcal{N}(z;\mu,\sigma^2)} [f(z)] = \mathbb{E}_{\mathcal{N}(\epsilon;0,1)} [f(\mu + \sigma\epsilon)] \simeq \frac{1}{L} \sum_{l=1}^L f(\mu + \sigma\epsilon^{(l)})$ where $\epsilon^{(l)} \sim \mathcal{N}(0, 1)$.

Algorithm 1 Minibatch version of the Auto-Encoding VB (AEVB) algorithm. Either of the two SGVB estimators in section 2.3 can be used. We use settings $M = 100$ and $L = 1$ in experiments.

```
 $\theta, \phi \leftarrow$  Initialize parameters
repeat
   $\mathbf{X}^M \leftarrow$  Random minibatch of  $M$  datapoints (drawn from full dataset)
   $\epsilon \leftarrow$  Random samples from noise distribution  $p(\epsilon)$ 
   $\mathbf{g} \leftarrow \nabla_{\theta, \phi} \tilde{\mathcal{L}}^M(\theta, \phi; \mathbf{X}^M, \epsilon)$  (Gradients of minibatch estimator (8))
   $\theta, \phi \leftarrow$  Update parameters using gradients  $\mathbf{g}$  (e.g. SGD or Adagrad [DHS10])
until convergence of parameters  $(\theta, \phi)$ 
return  $\theta, \phi$ 
```

VAE

- Now, we can form Monte Carlo estimates of expectations of some function w.r.t $q_\phi(z|x)$ as follows:

$$E_{\substack{z \sim q_\phi(z|x^{(i)})}} [f(z)] = E_{\substack{\epsilon \sim p(\epsilon)}} \left[f \left(\underbrace{g_\phi(\epsilon, x^{(i)})}_{\substack{\text{Monte Carlo estimates} \\ z}} \right) \right] \simeq \frac{1}{L} \sum_{l=1}^L f \left(g_\phi(\epsilon^{(l)}, x^{(i)}) \right)$$

where $\epsilon^{(l)} \sim p(\epsilon)$

- Apply this technique to the lower bd. $L(\theta, \phi; \mathbf{x}^{(i)}) = E_{q_\phi(z|x^{(i)})} [-\log q_\phi(z|x^{(i)}) + \log p_\theta(\mathbf{x}^{(i)}, z)]$
- Stochastic Gradient Variational Bayes** estimator $\tilde{L}^A(\theta, \phi; \mathbf{x}^{(i)}) \simeq L(\theta, \phi, \mathbf{x}^{(i)})$:

$$\tilde{L}^A(\theta, \phi; \mathbf{x}^{(i)}) = \frac{1}{L} \sum_{l=1}^L [\log p_\theta(\mathbf{x}^{(i)}, \mathbf{z}^{(i,l)}) - \log q_\phi(\mathbf{z}^{(i,l)}|\mathbf{x}^{(i)})]$$

where $\mathbf{z}^{(i,l)} = g_\phi(\epsilon^{(i,l)}, \mathbf{x}^{(i)})$ and $\epsilon^{(l)} \sim p(\epsilon)$

VAE

$$-D_{KL}(q_{\phi}(\mathbf{z}|\mathbf{x}^{(i)})||p_{\theta}(\mathbf{z})) + E_{q_{\phi}(\mathbf{z}|\mathbf{x}^{(i)})} [\log p_{\theta}(\mathbf{x}^{(i)}|\mathbf{z})]$$

- Often, the KL term $-D_{KL}(q_{\phi}(\mathbf{z}|\mathbf{x}^{(i)})||p_{\theta}(\mathbf{z}))$ can be integrated **analytically**. (when $q_{\phi}(\mathbf{z}|\mathbf{x}^{(i)})$ is Gaussian)
- only the expected reconstruction error $E_{q_{\phi}(\mathbf{z}|\mathbf{x}^{(i)})} [\log p_{\theta}(\mathbf{x}^{(i)}|\mathbf{z})]$ requires estimation by **sampling**
- $\tilde{L}^B(\theta, \phi; \mathbf{x}^{(i)}) = -D_{KL}(q_{\phi}(\mathbf{z}|\mathbf{x}^{(i)})||p_{\theta}(\mathbf{z})) + \frac{1}{L} \sum_{l=1}^L \log p_{\theta}(\mathbf{x}^{(i)}|\mathbf{z}^{(i,l)})$ if minibatch size L is large (e.g. $L=1000$), $\Leftarrow L=1$ in VAE.
- where $\mathbf{z}^{(i,l)} = g_{\phi}(\epsilon^{(i,l)}, \mathbf{x}^{(i)})$ and $\epsilon^{(l)} \sim p(\epsilon)$
- $-D_{KL}((q_{\phi}(\mathbf{z})||p_{\theta}(\mathbf{z})) = \int q_{\theta}(\mathbf{z}) (\log p_{\theta}(\mathbf{z}) - \log q_{\theta}(\mathbf{z})) d\mathbf{z}$

$$= \frac{1}{2} \sum_{j=1}^J (1 + \log((\sigma_j)^2) - (\mu_j)^2 - (\sigma_j)^2)$$

Pf)

$$-D_{KL}(\text{#}) = \int q_{\theta}(\mathbf{z}) (\log p_{\theta}(\mathbf{z}) - \log q_{\theta}(\mathbf{z})) d\mathbf{z} = \underbrace{\int q_{\theta}(\mathbf{z}) \log p_{\theta}(\mathbf{z}) d\mathbf{z}}_{\text{#1}} - \underbrace{\int q_{\theta}(\mathbf{z}) \log q_{\theta}(\mathbf{z}) d\mathbf{z}}_{\text{#2}}$$

①: $\int q_{\theta}(\mathbf{z}) \log p_{\theta}(\mathbf{z}) d\mathbf{z}$, where $q_{\theta}(\mathbf{z}_j) \sim N(\mu_j, \sigma_j^2)$ and $p_{\theta}(\mathbf{z}_j) \sim N(0, 1)$ (\mathbf{z}_j is independent)

$$\sim \log p_{\theta}(\mathbf{z}) = -\frac{1}{2} \log(2\pi) - \frac{1}{2} \sum_{j=1}^J z_j^2$$

$$\therefore \text{#1} = -\frac{1}{2} \log(2\pi) \cdot \int q_{\theta}(\mathbf{z}) d\mathbf{z} - \frac{1}{2} \int q_{\theta}(\mathbf{z}) \cdot \sum_{j=1}^J z_j^2 d\mathbf{z} = -\frac{1}{2} \log(2\pi) - \frac{1}{2} \int q_{\theta}(\mathbf{z}) \cdot (\sum_{j=1}^J z_j^2) d\mathbf{z} = -\frac{1}{2} \log(2\pi) - \frac{1}{2} \int q_{\theta}(\mathbf{z}) \cdot \mathbf{z}^2 d\mathbf{z}$$

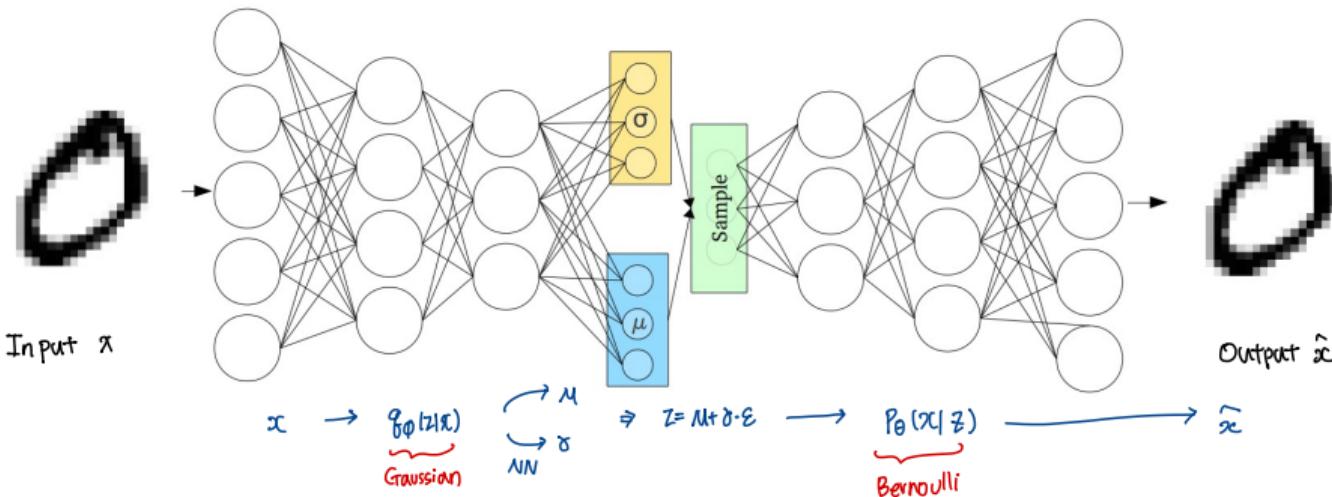
$$\boxed{\int q_{\theta}(\mathbf{z}) \cdot \mathbf{z}^2 d\mathbf{z}} = \int q_{\theta}(\mathbf{z}) \cdot z_1^2 d\mathbf{z} + \dots + \int q_{\theta}(\mathbf{z}) \cdot z_J^2 d\mathbf{z} = E_{q_{\theta}(\mathbf{z})}[z_1^2] + \dots + E_{q_{\theta}(\mathbf{z})}[z_J^2] = [V_{q_{\theta}}(z_1) + (E_{q_{\theta}}[z_1])^2] + \dots + [V_{q_{\theta}}(z_J) + (E_{q_{\theta}}[z_J])^2] + \frac{J}{2} (V_{q_{\theta}}(\mathbf{z}))^2$$

$$\text{#2} \text{ #2} = -\frac{1}{2} \log(2\pi) \cdot \frac{1}{2} \sum_{j=1}^J (\log \sigma_j^2)$$

$$\therefore \text{#1} - \text{#2} = -\frac{1}{2} \log(2\pi) \cdot \frac{1}{2} \sum_{j=1}^J (1 + \log \sigma_j^2)$$

$$= \frac{1}{2} \sum_{j=1}^J (1 + \log \sigma_j^2 - \sigma_j^2 - \mu_j^2)$$

VAE (MNIST)



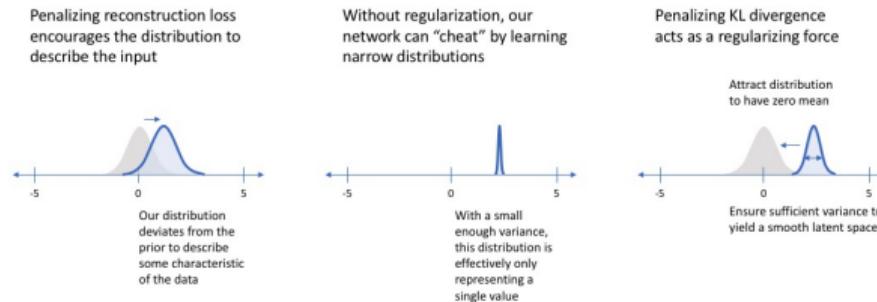
- $p(z) \sim N(0, I)$
- $q_{\phi}(z|x) \sim N(\mu_{\phi}(x), \text{diag}_{\phi}(x))$
- $p_{\theta}(x|z) \sim N(\mu_{\theta}(z), \delta_{\theta}(z))$ or $\text{Bernoulli}(p_{\theta}(z))$

Loss function

- $-D_{KL}(q_{\phi}(\mathbf{z}|\mathbf{x}^{(i)})||p_{\theta}(\mathbf{z}))$: $N(0, I)$

regularizer 역할. $q_{\phi}(z|x)$ 가 prior $p_{\theta}(z)$ 와 가까워지게 함.

$q(z|x)$ 가 gaussian distribution을 따르면 analytic하게 계산됨.



- $\log p_{\theta}(\mathbf{x}^{(i)}|\mathbf{z}^{(i)})$:

reconstruction term.

In our model, $p_{\theta}(x)$ is multivariate Bernoulli distribution. The $\log p_{\theta}(x|z)$ is calculated as follows:

In this case let $p_\theta(\mathbf{x}|\mathbf{z})$ be a multivariate Bernoulli whose probabilities are computed from \mathbf{z} with a fully-connected neural network with a single hidden layer:

$$\log p(\mathbf{x}|\mathbf{z}) = \sum_{i=1}^D x_i \log y_i + (1 - x_i) \cdot \log(1 - y_i)$$

where $\mathbf{y} = f_\sigma(\mathbf{W}_2 \tanh(\mathbf{W}_1 \mathbf{z} + \mathbf{b}_1) + \mathbf{b}_2)$ (11)

where $f_\sigma(\cdot)$ is the elementwise sigmoid activation function, and where $\theta = \{\mathbf{W}_1, \mathbf{W}_2, \mathbf{b}_1, \mathbf{b}_2\}$ are the weights and biases of the MLP.

- The total loss=reconstruction error(BCE)+KLD

$$\tilde{L}^B(\theta, \phi; \mathbf{x}^{(i)}) = \underbrace{-D_{KL}(q_\phi(\mathbf{z}|\mathbf{x}^{(i)}) || p_\theta(\mathbf{z}))}_{\text{minimize}} + \underbrace{\log p_\theta(\mathbf{x}^{(i)}|\mathbf{z}^{(i)})}_{\substack{\text{maximize} \\ \Updownarrow \\ \text{minimize binary cross entropy}}}$$

Code Example

```
# VAE model
class VAE(nn.Module):
    def __init__(self, image_size, hidden_size_1, hidden_size_2, latent_size):
        super(VAE, self).__init__()

        self.fc1 = nn.Linear(image_size, hidden_size_1)
        self.fc2 = nn.Linear(hidden_size_1, hidden_size_2)
        self.fc31 = nn.Linear(hidden_size_2, latent_size)
        self.fc32 = nn.Linear(hidden_size_2, latent_size)

        self.fc4 = nn.Linear(latent_size, hidden_size_2)
        self.fc5 = nn.Linear(hidden_size_2, hidden_size_1)
        self.fc6 = nn.Linear(hidden_size_1, image_size)

    def encode(self, x):
        h1 = F.relu(self.fc1(x))
        h2 = F.relu(self.fc2(h1))
        return self.fc31(h2), self.fc32(h2)

    def reparameterize(self, mu, logvar):
        std = torch.exp(0.5 * logvar)
        eps = torch.randn_like(std)
        return mu + std * eps

    def decode(self, z):
        h3 = F.relu(self.fc4(z))
        h4 = F.relu(self.fc5(h3))
        return torch.sigmoid(self.fc6(h4))

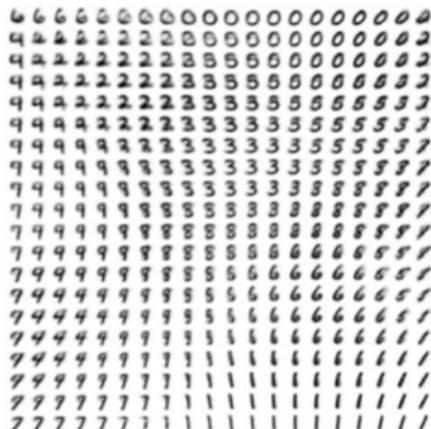
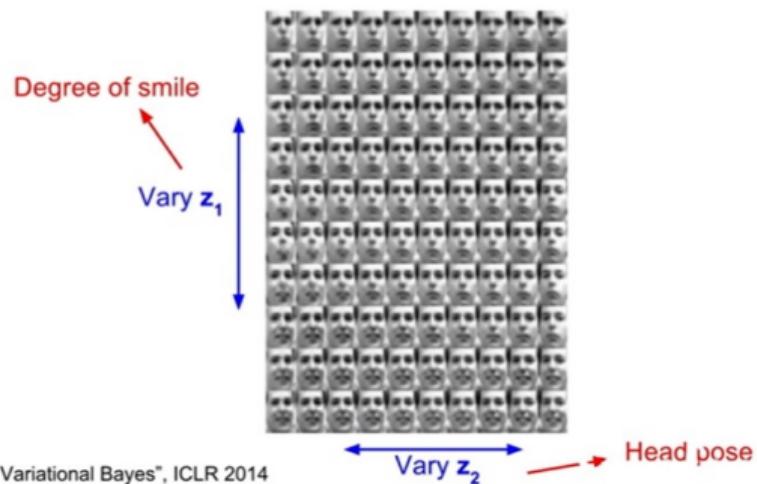
    def forward(self, x):
        mu, logvar = self.encode(x.view(-1, 784))
        z = self.reparameterize(mu, logvar)
        return self.decode(z), mu, logvar
```

```
def loss_function(recon_x, x, mu, logvar):  
    BCE = F.binary_cross_entropy(recon_x, x.view(-1, 784), reduction = 'sum')  
    KLD = -0.5 * torch.sum(1 + logvar - mu.pow(2) - logvar.exp())  
    return BCE, KLD
```

```
def train(epoch, model, train_loader, optimizer):  
    model.train()  
    train_loss = 0  
    for batch_idx, (data, _) in enumerate(train_loader):  
        data = data.to(DEVICE)  
        optimizer.zero_grad()  
  
        recon_batch, mu, logvar = model(data)  
  
        BCE, KLD = loss_function(recon_batch, data, mu, logvar)  
  
        loss = BCE + KLD  
  
        writer.add_scalar("Train/Reconstruction Error", BCE.item(), batch_idx + epoch * (len(train_loader.dataset)/BATCH_SIZE) )  
        writer.add_scalar("Train/KL-Divergence", KLD.item(), batch_idx + epoch * (len(train_loader.dataset)/BATCH_SIZE) )  
        writer.add_scalar("Train/Total Loss" , loss.item(), batch_idx + epoch * (len(train_loader.dataset)/BATCH_SIZE) )  
  
        loss.backward()  
  
        train_loss += loss.item()  
  
        optimizer.step()  
  
    if batch_idx % 100 == 0:  
        print('Train Epoch: {} [{}/{} ({:.0f}%)]\t Loss: {:.6f}'.format(  
            epoch, batch_idx * len(data), len(train_loader.dataset),  
            100. * batch_idx / len(train_loader),  
            loss.item() / len(data)))  
  
    print("===== > Epoch: {} Average loss: {:.4f}".format(  
        epoch, train_loss / len(train_loader.dataset)  
    ))
```

VAE

- 따라서 VAE의 장점은 posterior인 $q_\phi(z|x)$ 와 likelihood인 $p_\theta(x|z)$ 를 각각 encoder/decoder로 modeling하여 이를 한번에 gradient descent method를 사용해 update한다는 것.
- 단점은 생성된 image들이 다소 blurry하다는 것!



Reference

Auto-Encoding Variational Bayes, Kingma and Maw Welling 2014

<https://www.youtube.com/watch?v=5WoltGTWV54>

<https://cumulu-s.tistory.com/25>

GAN

Generative Adversarial Networks

ESC 27기 윤이경

1

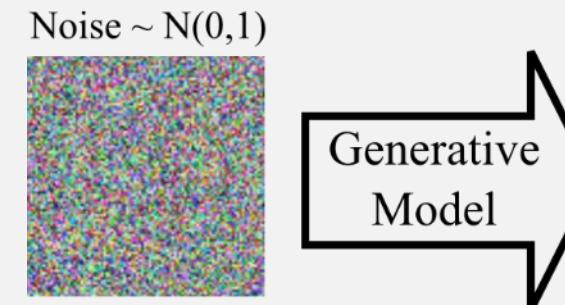
Generative Adversarial Networks

- Review : Density Estimation, Adversarial Games
- What is Generative Adversarial Networks?
- GAN _ Goal of Formulation
- Training and Performance Evaluation, Challenges
- Generative Moment Matching Networks

Density Estimation

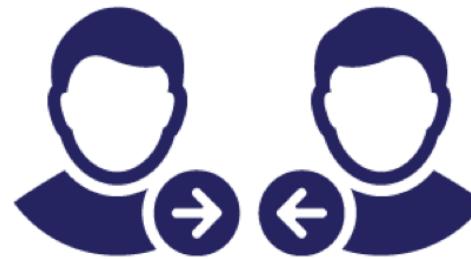
GAN – GAN just want ability to sample, not explicitly modeling density

- 1) Sample from a simple distribution (random noise)
- 2) Learn a transformation to the training distribution.



- Q. How to represent this complex transformation?
A. By using a neural net.

Adversarial Games (Game-theory)



BUYER / SELLER

GAN also take a game-theoretic approach.
Learn to generate from training distribution through 2-player game.

$$(x^*, y^*) = \arg \max_y \min_x V(x, y).$$

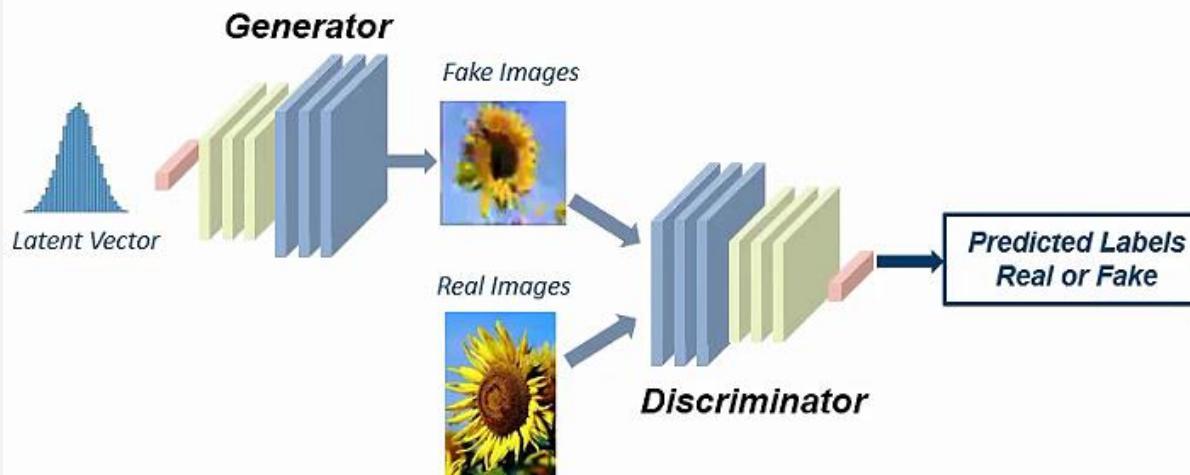
Generative Adversarial Networks



Generative Adversarial Networks

What is GAN(Generative Adversarial Network)?

Train to trick the Discriminator



Train to judge real / fake correctly



Generate an image similar to real images



Generative Adversarial Networks

생성자와 판별자 두개의 네트워크를 활용한 생성 모델.

1) Generator network :

Goal : try to fool discriminator by generating real-looking images

Input : some random noise

Output : image x

$$x = G(z; \theta^{(g)})$$

2) Discriminator network :

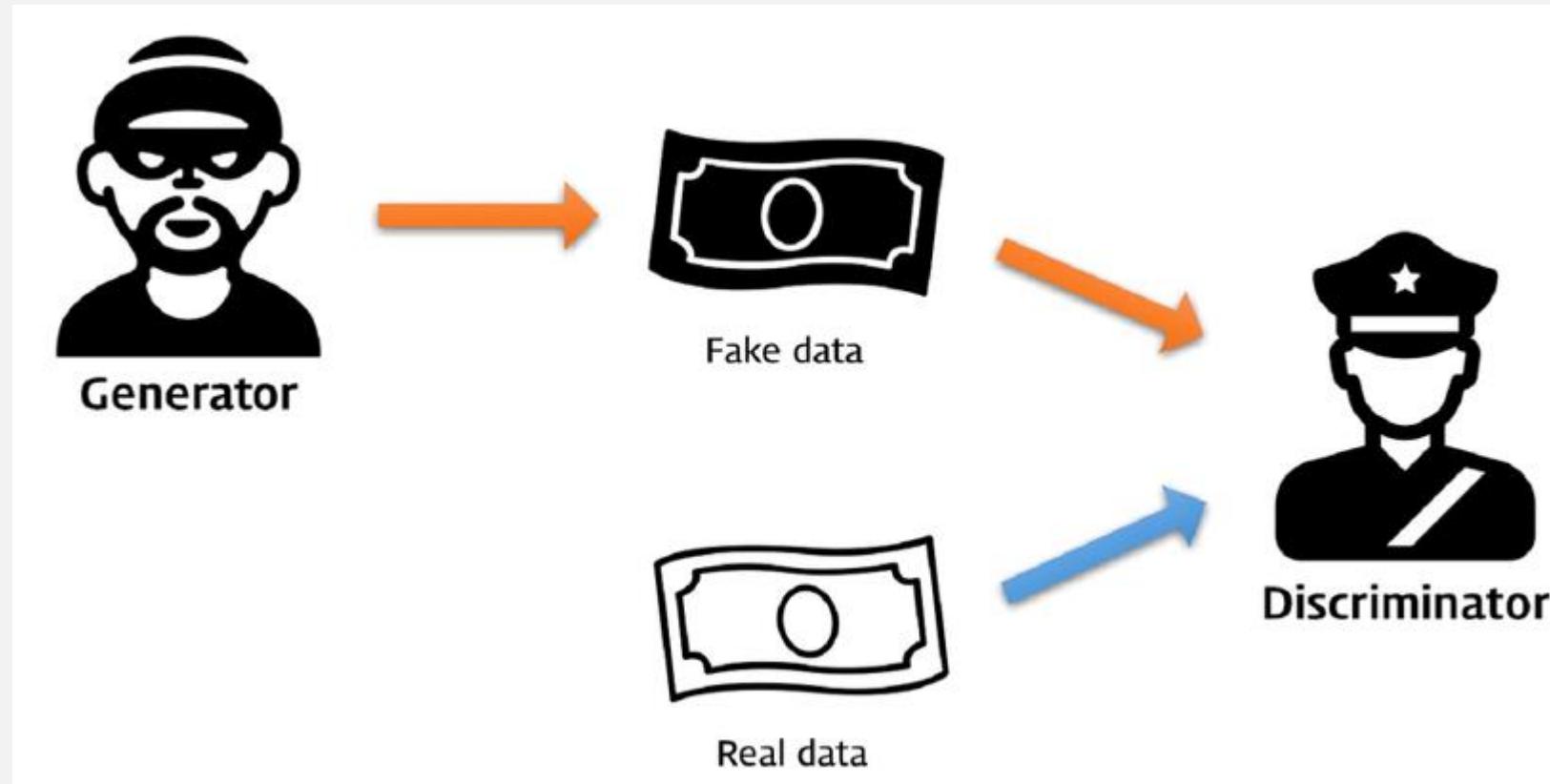
Goal : try to distinguish between real and fake images

Input : image x

Output : probability between 0 and 1
(real = 1 ~ fake = 0)

$$D(x; \theta^{(d)})$$

Generative Adversarial Networks



Generative Adversarial Networks



VS



판별자 D (Discriminator)

“진짜 화폐랑 G가 만들어내는 가짜 화폐를
잘 구별해야지!!”

생성자 G (Generator)

“진짜 화폐랑 내가 만들어낸 가짜 화폐를
D가 구분할 수 없도록
완벽하게 진짜같은 가짜 화폐를 만들어낼거야!!”

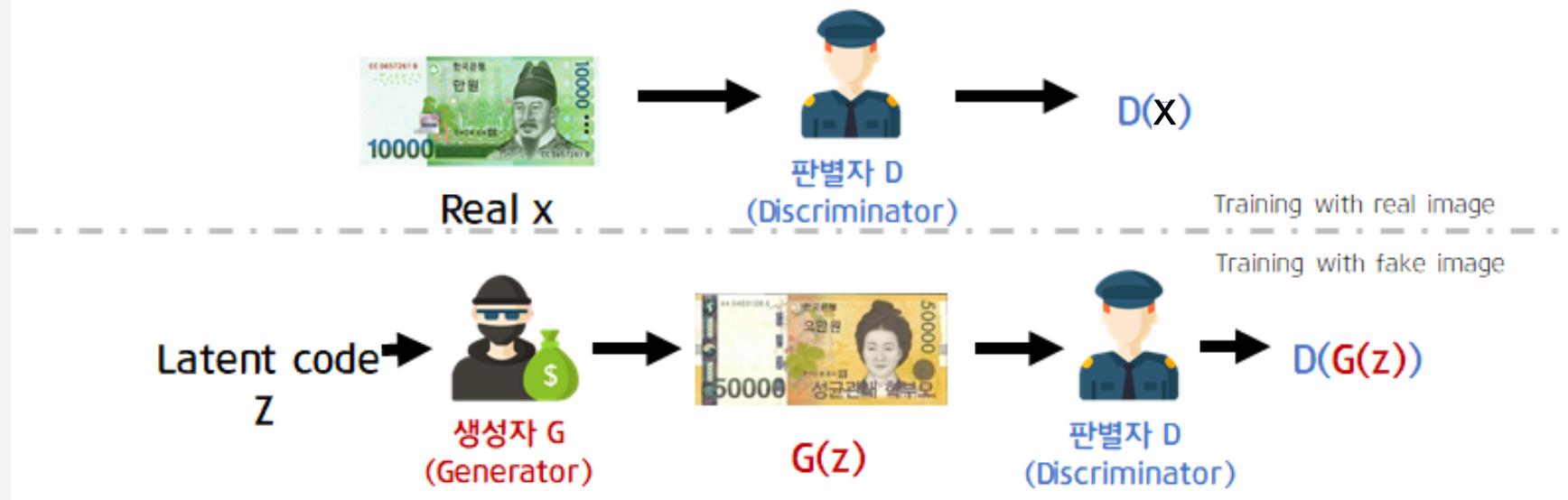
Generative Adversarial Networks

$$\min_G \max_D V(D, G) = E_{x \sim p_{data}(x)}[\log D(x)] + E_{z \sim P_z(z)}[\log(1 - D(G(z)))]$$

$D(x) = 1$ 일 때 최대

$D(G(z)) = 1$ 일 때 최소

$D(G(z)) = 0$ 일 때 최대



Generative Adversarial Networks

Training objectives

$$P_{model} \sim G(z)$$

Loss function The discriminator tries to maximize the payoff

$$V(G, D) = \mathbb{E}_{x \sim p_{data}} [\ln D(x)] + \mathbb{E}_{x \sim p_{model}} [\ln(1 - D(x))], \quad (19.4.4)$$

Maximizing the first term assures for a correct classification of true data, while maximizing the second term ensures for the correct classification of the model generated data. This follows from the fact that the discriminator's output $D(x)$ tends to be close to 1 for true data x , while $1 - D(x)$ tends to be close to 1 for generated data x .

$$\min_G \max_D V(D, G) = \mathbb{E}_{\mathbf{x} \sim p_{data}(\mathbf{x})} [\log D(\mathbf{x})] + \mathbb{E}_{\mathbf{z} \sim p_{\mathbf{z}}(\mathbf{z})} [\log(1 - D(G(\mathbf{z})))].$$

Generative Adversarial Networks

Minimax problem

$$\min_G \max_D V(D, G) = \mathbb{E}_{\mathbf{x} \sim p_{\text{data}}(\mathbf{x})} [\log D(\mathbf{x})] + \mathbb{E}_{\mathbf{z} \sim p_{\mathbf{z}}(\mathbf{z})} [\log(1 - D(G(\mathbf{z})))].$$

Discriminator : θ_d wants to **maximize** objective

- $D(\mathbf{x})$: how likely \mathbf{x} is real \rightarrow close to 1 (real)
- $D(G(\mathbf{z}))$: how likely $G(\mathbf{z})$ is fake \rightarrow close to 0 (fake)

이 두 관계 적대적

Generator : θ_g wants to **minimize** objective

- $D(G(\mathbf{z}))$: how likely $G(\mathbf{z})$ is real \rightarrow close to 1 (look-real)
- discriminator is fooled into thinking that generated $G(\mathbf{z})$ is real

Generative Adversarial Networks

Each gradient update

1. Discriminator update (repeat k times per iteration)
 - 1) prepare 'm' synthetic samples (using generator) and 'm' real samples
 - 2) feed these '2m' samples to discriminator
 - 3) Gradient ascent : backprop on discriminator network to update θ_d
2. Generator update (once per iteration)
 - 1) prepare 'm' synthetic samples using generator
 - 2) feed these 'm' samples to discriminator
 - 3) Gradient decent : backprop on whole network
 - backward gradient from discriminator output → generator
 - but update θ_g only

Generative Adversarial Networks

Algorithm 1 Minibatch stochastic gradient descent training of generative adversarial nets. The number of steps to apply to the discriminator, k , is a hyperparameter. We used $k = 1$, the least expensive option, in our experiments.

for number of training iterations **do**

for k steps **do**

- Sample minibatch of m noise samples $\{\mathbf{z}^{(1)}, \dots, \mathbf{z}^{(m)}\}$ from noise prior $p_g(\mathbf{z})$.
- Sample minibatch of m examples $\{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(m)}\}$ from data generating distribution $p_{\text{data}}(\mathbf{x})$.
- Update the discriminator by ascending its stochastic gradient:

$$\nabla_{\theta_d} \frac{1}{m} \sum_{i=1}^m \left[\log D(\mathbf{x}^{(i)}) + \log (1 - D(G(\mathbf{z}^{(i)}))) \right].$$

end for

- Sample minibatch of m noise samples $\{\mathbf{z}^{(1)}, \dots, \mathbf{z}^{(m)}\}$ from noise prior $p_g(\mathbf{z})$.

- Update the generator by descending its stochastic gradient:

$$\nabla_{\theta_g} \frac{1}{m} \sum_{i=1}^m \log (1 - D(G(\mathbf{z}^{(i)}))).$$

end for

The gradient-based updates can use any standard gradient-based learning rule. We used momentum in our experiments.

```
data_loader = torch.utils.data.DataLoader(dataset=mnist,
                                         batch_size=batch_size,
                                         shuffle=True)

# Discriminator
D = nn.Sequential(
    nn.Linear(image_size, hidden_size),
    nn.LeakyReLU(0.2),
    nn.Linear(hidden_size, hidden_size),
    nn.LeakyReLU(0.2),
    nn.Linear(hidden_size, 1),
    nn.Sigmoid())

# Generator
G = nn.Sequential(
    nn.Linear(latent_size, hidden_size),
    nn.ReLU(),
    nn.Linear(hidden_size, hidden_size),
    nn.ReLU(),
    nn.Linear(hidden_size, image_size),
    nn.Tanh())

# Device setting
D = D.to(device)
G = G.to(device)

# Binary cross entropy loss and optimizer
criterion = nn.BCELoss()
d_optimizer = torch.optim.Adam(D.parameters(), lr=0.0002)
g_optimizer = torch.optim.Adam(G.parameters(), lr=0.0002)
```

```
def denorm(x):
    out = (x + 1) / 2
    return out.clamp(0, 1)

def reset_grad():
    d_optimizer.zero_grad()
    g_optimizer.zero_grad()
```

<https://github.com/yunjey/pytorch-tutorial>

```

# Start training
total_step = len(data_loader)
for epoch in range(num_epochs):
    for i, (images, _) in enumerate(data_loader):
        images = images.reshape(batch_size, -1).to(device)

        # Create the labels which are later used as input for the BCE loss
        real_labels = torch.ones(batch_size, 1).to(device)
        fake_labels = torch.zeros(batch_size, 1).to(device)

        # =====#
        # Train the discriminator
        # =====#
        # Compute BCE_Loss using real images where BCE_Loss(x, y): - y * log(D(x)) - (1-y) * log(1-D(x))
        # Second term of the loss is always zero since real_labels == 1
        outputs = D(images)
        d_loss_real = criterion(outputs, real_labels)
        real_score = outputs

        # Compute BCELoss using fake images
        # First term of the loss is always zero since fake_labels == 0
        z = torch.randn(batch_size, latent_size).to(device)
        fake_images = G(z)
        outputs = D(fake_images)
        d_loss_fake = criterion(outputs, fake_labels)
        fake_score = outputs

        # Backprop and optimize
        d_loss = d_loss_real + d_loss_fake
        reset_grad()
        d_loss.backward()
        d_optimizer.step()

# =====#
# Train the generator
# =====#
# Compute loss with fake images
z = torch.randn(batch_size, latent_size).to(device)
fake_images = G(z)
outputs = D(fake_images)

# We train G to maximize log(D(G(z)) instead of minimizing log(1-D(G(z)))
# For the reason, see the last paragraph of section 3. https://arxiv.org/pdf/1406.2661.pdf
g_loss = criterion(outputs, real_labels)

# Backprop and optimize
reset_grad()
g_loss.backward()
g_optimizer.step()

if (i+1) % 200 == 0:
    print('Epoch [{}/{}], Step [{}/{}], d_loss: {:.4f}, g_loss: {:.4f}, D(x): {:.2f}, D(G(z)): {:.2f}'
          .format(epoch, num_epochs, i+1, total_step, d_loss.item(), g_loss.item(),
                  real_score.mean().item(), fake_score.mean().item()))

# Save real images
if (epoch+1) == 1:
    images = images.reshape(images.size(0), 1, 28, 28)
    save_image(denorm(images), os.path.join(sample_dir, 'real_images.png'))

# Save sampled images
fake_images = fake_images.reshape(fake_images.size(0), 1, 28, 28)
save_image(denorm(fake_images), os.path.join(sample_dir, 'fake_images-{}.png'.format(epoch+1)))

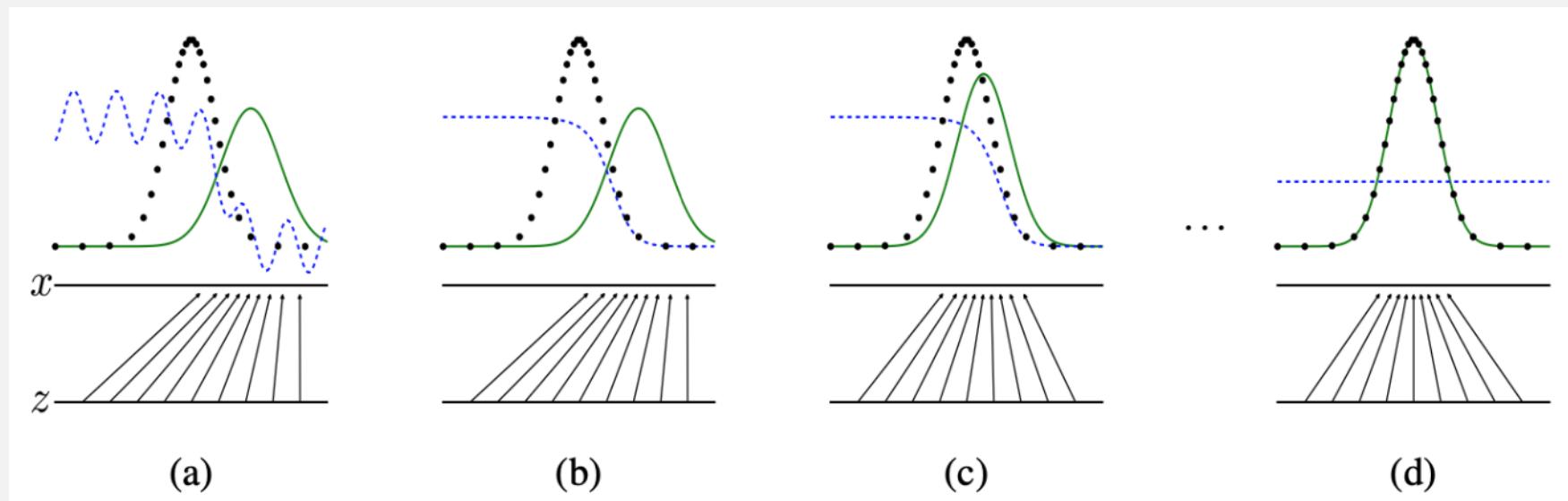
# Save the model checkpoints
torch.save(G.state_dict(), 'G.ckpt')
torch.save(D.state_dict(), 'D.ckpt')

```

Generative Adversarial Networks

Goal of Formulation

$$P_{g(z)}(P_{model}) \rightarrow P_{data}, \quad D(G(z)) = \frac{1}{2} \quad (G(z) \text{ is not distinguishable by } D)$$



How can the formulation lead $P_{g(z)}$ converge to P_{data} ?

Generative Adversarial Networks

Global optimal

Proposition 19.4.3 *For a fixed generator, G , the optimal discriminator (19.4.5) is given by*

$$D_G^*(x) = \frac{p_{data}(x)}{p_{data}(x) + p_{model}(x)}. \quad (19.4.6)$$

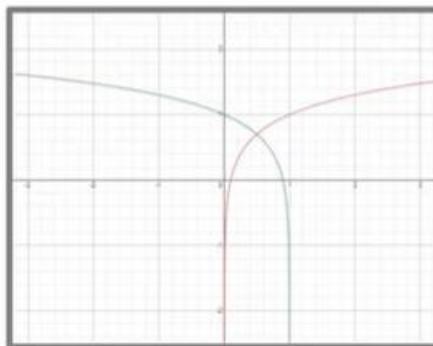
Generative Adversarial Networks

$$\text{Proposition: } D_G^*(x) = \frac{p_{data}(x)}{p_{data}(x) + p_g(x)}$$

Proof: For G fixed,

$$\begin{aligned} V(G, D) &= E_{x \sim p_{data}(x)}[\log D(x)] + E_{z \sim p_z(z)}[\log(1 - D(G(z)))] \\ &= \int_x p_{data}(x) \log(D(x)) dx + \int_z p_z(z) \log(1 - D(g(z))) dz \\ &= \int_x p_{data}(x) \log(D(x)) + p_g(x) \log(1 - D(x)) dx \end{aligned}$$

$$E[X] = \int_{-\infty}^{\infty} xf(x)dx$$



function $y \rightarrow a \log(y) + b \log(1 - y)$ achieves its maximum in $[0, 1]$ at $\frac{a}{a+b}$

same as *optimal control*: $\frac{\delta V(G, D)}{\delta D}[D^*(x)] = 0$

Generative Adversarial Networks

Proof: Considering distributions p_{data} and p_{model} defined on the space \mathcal{X} , the payoff function can be written as the following integral:

$$V(G, D) = \int_{\mathcal{X}} \left[p_{data}(x) \ln D(x) + p_{model}(x) \ln(1 - D(x)) \right] dx.$$

Since G is fixed, we consider the payoff as a functional of $D(x)$ as

$$F(D) = \int_{\mathcal{X}} L(D(x)) dx,$$

with the Lagrangian function

$$L(G) = p_{data}(x) \ln D(x) + p_{model}(x) \ln(1 - D(x)).$$

Solving for D we obtain the solution

$$D_G^*(x) = \frac{p_{data}(x)}{p_{data}(x) + p_{model}(x)}.$$

In order to show that this critical point corresponds to a maximum, we consider the second variation and show that it has a negative value

$$\frac{d^2 L(D)}{dD^2} = -\frac{p_{data}(x)}{D(x)^2} - \frac{p_{model}(x)}{(1 - D(x))^2}.$$

Hence, $D_G^*(x)$ corresponds to a maximum of $V(G, D)$ for a given G . ■

Critical points satisfy the variational equation

$$\frac{dL(D)}{dD} = 0.$$

This becomes

$$\frac{p_{data}(x)}{D(x)} - \frac{p_{model}(x)}{1 - D(x)} = 0.$$

Generative Adversarial Networks

Global optimum point

Proposition 19.4.6 *The global minimum*

$$G^* = \arg \min_G V(G, D_G^*)$$

is achieved if and only if $p_{model} = p_{data}$. At this minimum point the value is $-\ln 4$.

Generative Adversarial Networks

Proposition: Global optimum point is $p_g = p_{data}$

Proof:

$$C(G) = \max_D V(G, D) = E_{x \sim p_{data}(x)}[\log D^*(x)] + E_{z \sim p_z(z)}[\log(1 - D^*(G(z)))]$$

$$= E_{x \sim p_{data}(x)} \left[\log \frac{p_{data}(x)}{p_{data}(x) + p_g(x)} \right] + E_{x \sim p_g(x)} \left[\log \frac{p_g(x)}{p_{data}(x) + p_g(x)} \right]$$

removed
when $p_g = p_{data}$

$$D_G^*(x) = \frac{p_{data}(x)}{p_{data}(x) + p_g(x)}$$

$$= E_{x \sim p_{data}(x)} \left[\log \frac{2 * p_{data}(x)}{p_{data}(x) + p_g(x)} \right] + E_{x \sim p_g(x)} \left[\log \frac{2 * p_g(x)}{p_{data}(x) + p_g(x)} \right] - \log(4)$$

$$KL(p_{data} || p_g) = \int_{-\infty}^{\infty} p_{data}(x) \log \left(\frac{p_{data}(x)}{p_g(x)} \right) dx$$

$$= KL(p_{data} || \frac{p_{data}(x) + p_g(x)}{2}) + KL(p_g || \frac{p_{data}(x) + p_g(x)}{2}) - \log(4)$$

$$= 2 * JSD(p_{data} || p_g) - \log(4)$$

$$JSD(p || q) = \frac{1}{2} KL(p || \frac{p + q}{2}) + \frac{1}{2} KL(q || \frac{p + q}{2})$$

- The *Kullback-Leibler* (KL) divergence

$$KL(\mathbb{P}_r \parallel \mathbb{P}_g) = \int \log \left(\frac{P_r(x)}{P_g(x)} \right) P_r(x) d\mu(x) ,$$

- The *Jensen-Shannon* (JS) divergence

$$JS(\mathbb{P}_r, \mathbb{P}_g) = KL(\mathbb{P}_r \parallel \mathbb{P}_m) + KL(\mathbb{P}_g \parallel \mathbb{P}_m) ,$$

Generative Adversarial Networks

Performance evaluation

Desideratum #1 : Quality

- low entropy (highly predictable) conditional label distribution $p(y|x)$
= given an image x , should know its label y easily
- use neural net to classify generated images → gives distribution $p(y|x)$

Desideratum #2 : Diversity

- high entropy (less predictable) $p(y)$
$$p(y) = \int p(y | x = G(z)) dz$$

= if generated images are diverse, $p(y)$ should be uniform

There is often trade-off between quality and diversity

Generative Adversarial Networks

Lorem Ipsum is simply dummy text of the printing and typesetting industry.

Performance evaluation

1. IS: Inception Score

$$\text{IS}(G) = \exp \left(\mathbb{E}_{\mathbf{x} \sim p_g} D_{KL} (p(y|\mathbf{x}) \parallel p(y)) \right),$$

2. FID: Frechet Inception Distance

$$FID = \|\mu_r - \mu_g\|^2 + T_r(\Sigma_r + \Sigma_g - 2(\Sigma_r \Sigma_g)^{1/2})$$

= actually compares statistics of generated samples to real sample

Generative Adversarial Networks

GAN Training challenges

1. **non-convergence** : Model parameters oscillate/ destabilize/ never convergence.

Problem in GAN training :

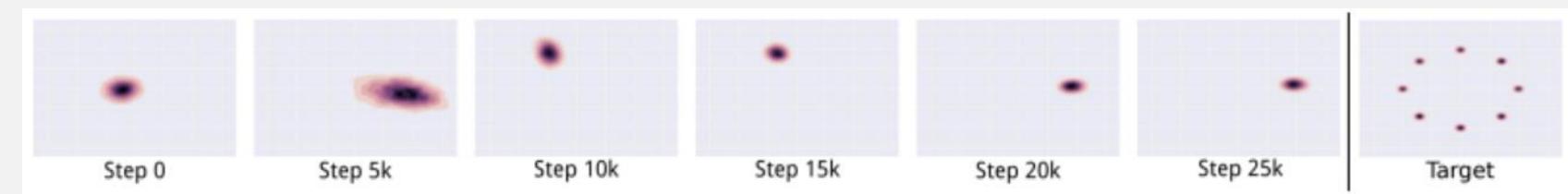
- downhill move of one player -> may push the other player uphill
- a general problem with games not unique to GAN

In theory : updates in function space : guaranteed convergence (a convex problem)

In practice : updates in parameter space (non-convex) : no such guarantee

Consequences :

- **oscillation** :
- **mode collapse** :

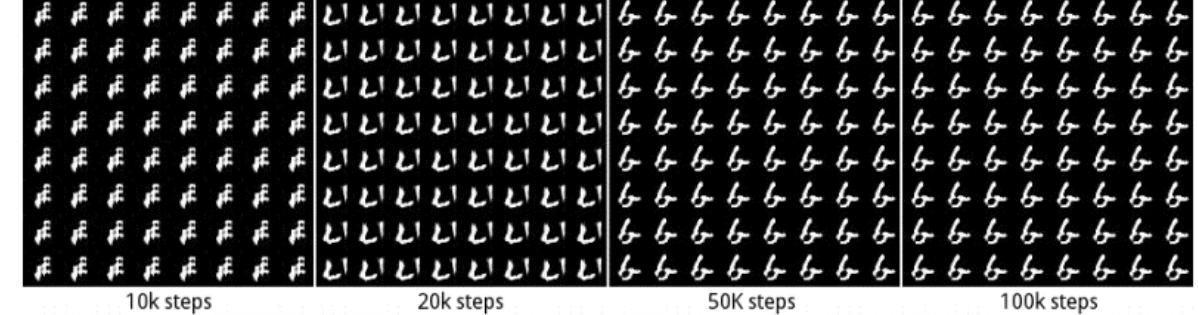


Generative Adversarial Networks

GAN Training challenges

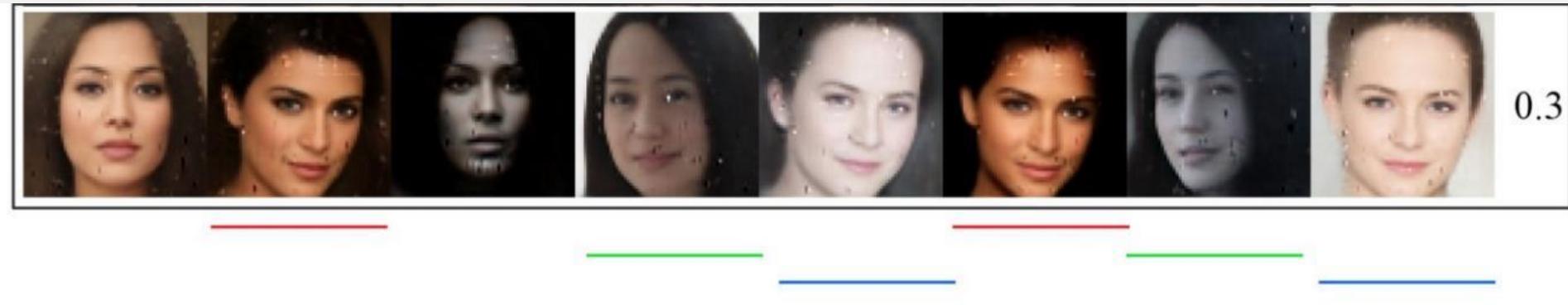
2. **Mode collapse** : most common form of harmful non-convergence
: several different input z values to the same output point

- : Discriminator와 generator가 학습이 서로 상호작용을 하면서 같이 학습이 진행 되어야 하는데, 한쪽이 너무 학습이 잘 되어버려서(학습의 불균형) mode collapse가 발생
- : Discriminator와 generator 사이의 능력에 적절한 균형을 이루면서, 두 네트워크가 안정적으로 전역해로 수렴하도록 만드는 것이 GAN이 해결해야 할 숙제



Generative Adversarial Networks

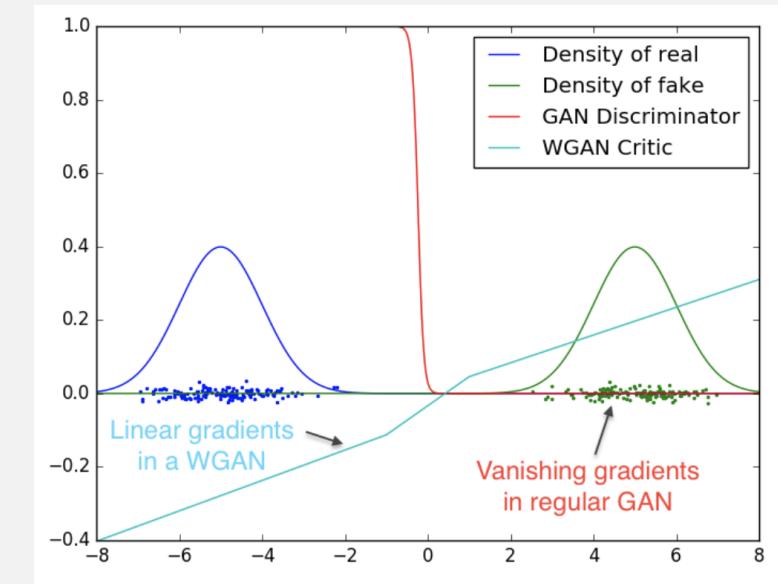
Partial mode collapse



Generative Adversarial Networks

[딥러닝 GAN 튜토리얼 - 시작부터 최신 트렌드까지 GAN 논문 순서 | mocha's machine learning \(ysbsb.github.io\)](https://mocha-ml.tistory.com/1)

- : Wasserstein GAN
- : Deep Convolutional GAN (DCGAN)
- : Least Squares GAN (LSGAN)
- : Semi-Supervised GAN (SGAN)
- : Auxiliary Classifier GAN (ACGAN)
- : CycleGAN



name change : discriminator \rightarrow critic

- ✓ GAN discriminator : predict the probability of generated images being real
- ✓ WGAN critic : score the realness/fakeness a given image

Generative Moment Matching Networks

Generative moment matching networks have been introduced by Li et al in 2015.

The generator is trained by ‘moment matching’, a technique by which the generator’s outcome has moments as close as possible to the corresponding moments of training data.

⇒ consider the parameter θ that minimizes the maximum mean discrepancy defined by

$$\theta^* = \arg \min_{\theta} d_{MMD}(p_{\theta}, q) = \arg \min_{\theta} \|\mu_{\phi}(G(Z; \theta)) - \mu_{\phi}(Y)\|_{Eu}.$$

$$\theta^* = \arg \min_{\theta} \left\{ \frac{1}{n^2} \sum_{i,j} K(x_i, x_j) + \frac{1}{m^2} \sum_{i,j} K(y_i, y_j) - \frac{2}{mn} \sum_{i,j} K(x_i, y_j) \right\},$$

* maximum mean discrepancy(MMD) : distance between two distributions