



BERT & Generative Model

Sep 22th 2022/ 7기 김예진

1. After Transformer

- NLP BERT
- I/O
- Pre-trained
- Fine-tuning

2. Generative Model

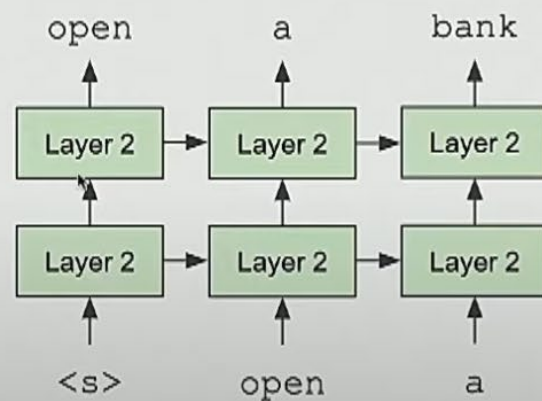
- Bayesian
- (Un)Supervised “Learning” - AE
- Generative & Discriminative Model
- Models - VAE

1. After Transformer

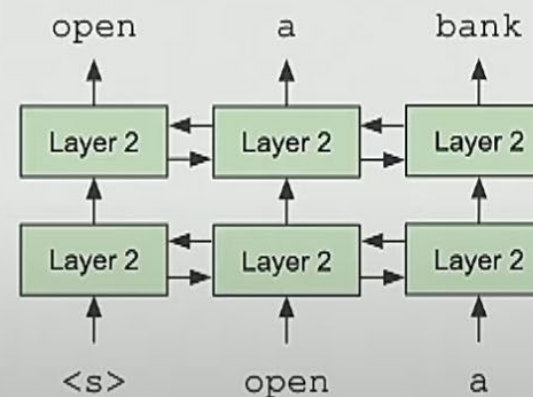
NLP – BERT

(sentence level) Bidirectional in RNN

Unidirectional context
Build representation incrementally



Bidirectional context
Words can “see themselves”



1. After Transformer

NLP – BERT

BERT = **Bidirectional Encoder Representation from Transformer!**

Key point)

- **Bidirectional**: 양방향의 정보를 어떻게 학습할까?
- **Encoder**: 기존 Transformer에서 Encoder

Run)

Task에 무관한 모델 구조를 가진다!

→ Representation도 무관해야 한다!

- **pre-trained**: Unsupervised
- **Fine-tuning**: (All param from pre-trained)

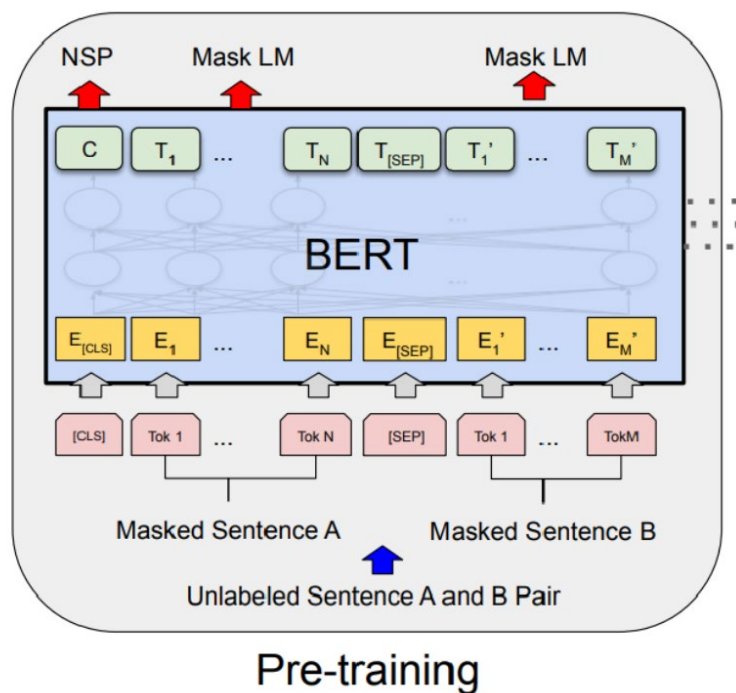
Supervised depending **on task**

유의사항)

- **Sentence**: 연속된 단어의 나열, sen
- **Sequence**: 하나 혹은 두 sentence, seq

1. After Transformer

NLP – BERT



Input = Masked Sen A & Masked Sen B

Output(Task): NSP와 Masked LM

NSP = Next Sentence Prediction

Sentence A와 Sentence B가 연속된 문장인지 분류

→ 맥락을 이해하는 학습 = **Encoder**?!

Masked LM: Mask된 내용물 예측

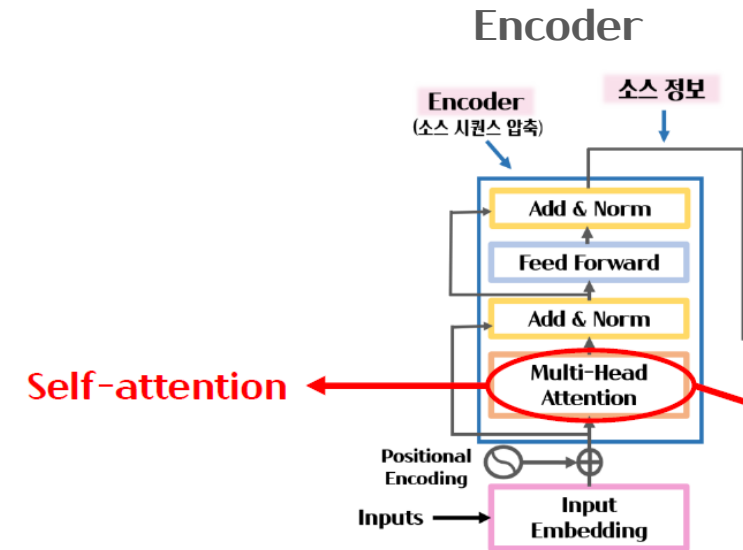
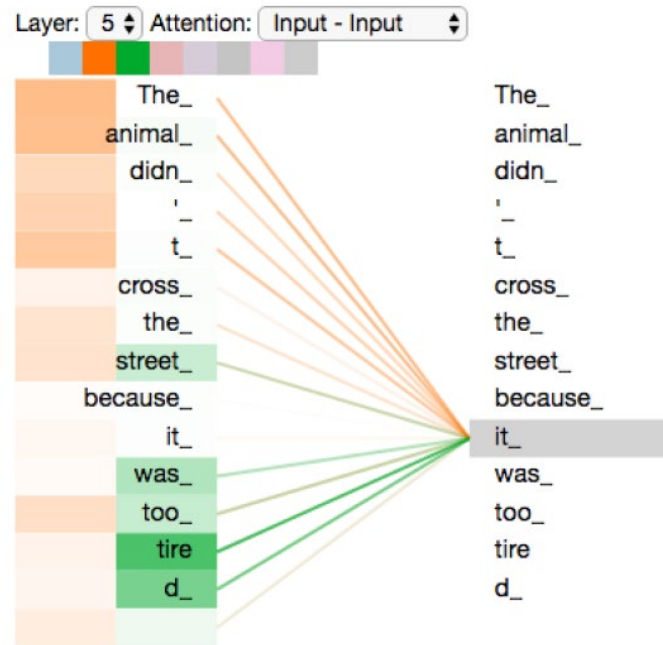
일부 토큰이 뭔지 모를 때 그 토큰이 뭔지 예측

→ 앞뒤의 토큰을 보고 어떤 토큰인지 보기에 **Bidirectional**!

1. After Transformer

NLP - BERT

Bidirectionality & Transformer.Encoder



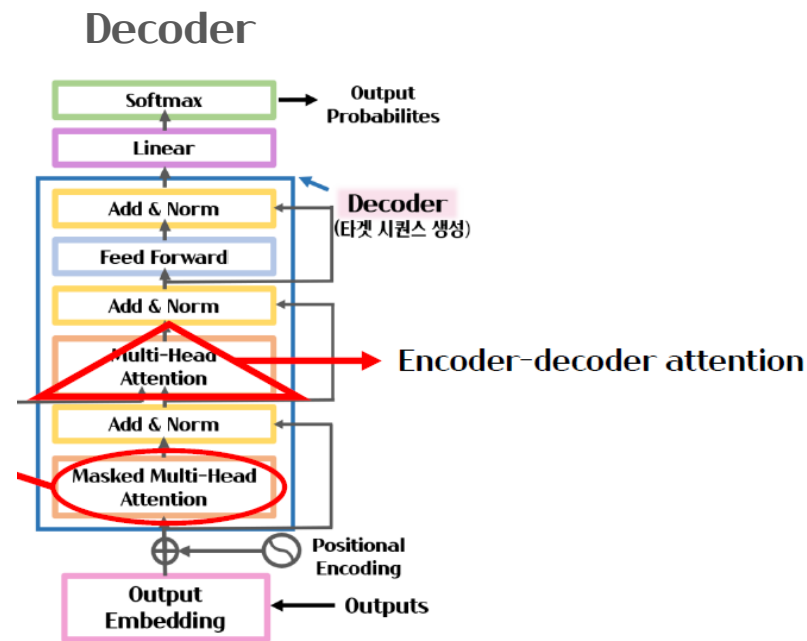
Self-attention 또한 일종의 Transformer에서 Bidirectionality를 녹여낸 부분이긴 함

1. After Transformer

NLP – BERT

Bidirectionality & Transformer.Encoder

	Query	Key
	<s>	어제
	I	카페
masking	went	갔었어
	to	거기
	the	사람
	café	많더라
	...	

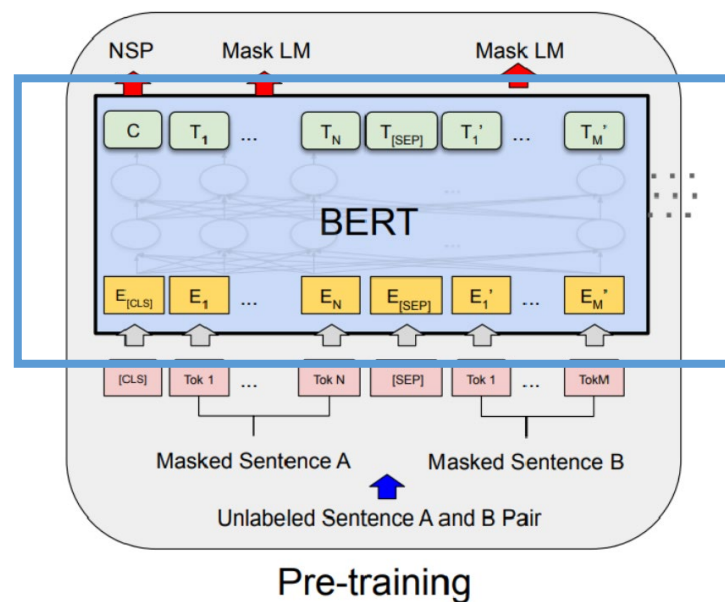
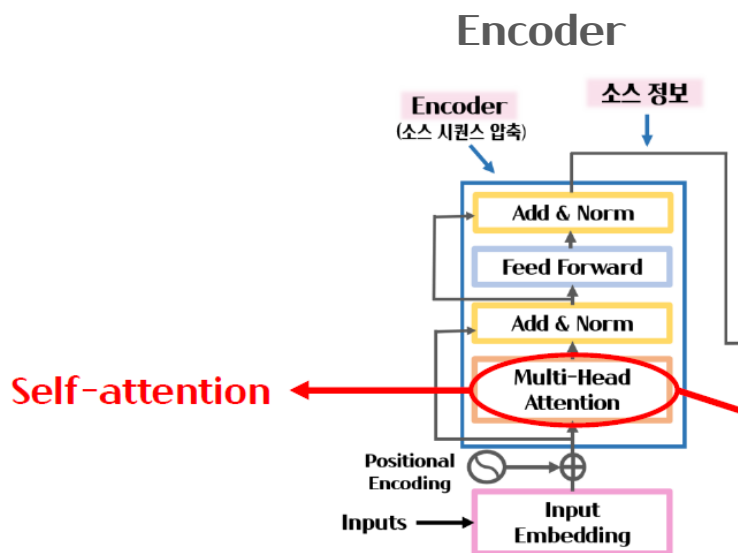


다만 Decoder에서 하는 Task는 결국 순서대로 들어오는 방식으로 진행...!

1. After Transformer

NLP - BERT

Bidirectionality & Transformer.Encoder



그래서 Transformer의 **Encoder** 구조만을 이용해서 Task(**NSP**, **Masked LM**) 수행하고 학습

1. After Transformer

I/O

Input	[CLS]	my	dog	is	cute	[SEP]	he	likes	play	##ing	[SEP]
Token Embeddings	$E_{[CLS]}$	E_{my}	E_{dog}	E_{is}	E_{cute}	$E_{[SEP]}$	E_{he}	E_{likes}	E_{play}	$E_{\# \# ing}$	$E_{[SEP]}$
	+	+	+	+	+	+	+	+	+	+	+
Segment Embeddings	E_A	E_A	E_A	E_A	E_A	E_A	E_B	E_B	E_B	E_B	E_B
	+	+	+	+	+	+	+	+	+	+	+
Position Embeddings	E_0	E_1	E_2	E_3	E_4	E_5	E_6	E_7	E_8	E_9	E_{10}

Input = **Token** + Segment + Position

Token

- Tokenizer : Word Piece

- CLS: seq의 시작에 위치

Aggregate sequence representation

For classification

- SEP: 두 sentence 구분

Cf) Word piece tokenizer:

Merge bigram which most increases the likelihood of the data

→ 자주 등장하는 연속된 token 쌍은 하나의 token으로 사용

1. After Transformer

I/O

Input	[CLS]	my	dog	is	cute	[SEP]	he	likes	play	##ing	[SEP]
Token Embeddings	$E_{[CLS]}$	E_{my}	E_{dog}	E_{is}	E_{cute}	$E_{[SEP]}$	E_{he}	E_{likes}	E_{play}	$E_{\#ing}$	$E_{[SEP]}$
	+	+	+	+	+	+	+	+	+	+	+
Segment Embeddings	E_A	E_A	E_A	E_A	E_A	E_A	E_B	E_B	E_B	E_B	E_B
	+	+	+	+	+	+	+	+	+	+	+
Position Embeddings	E_0	E_1	E_2	E_3	E_4	E_5	E_6	E_7	E_8	E_9	E_{10}

Input = Token + Segment + Position

Segment

- 각 token들이 sen A인지 sen B인지
- CLS와 SEP는 sen A에 포함

1. After Transformer

I/O

Input	[CLS]	my	dog	is	cute	[SEP]	he	likes	play	##ing	[SEP]
Token Embeddings	$E_{[CLS]}$	E_{my}	E_{dog}	E_{is}	E_{cute}	$E_{[SEP]}$	E_{he}	E_{likes}	E_{play}	$E_{\# \# ing}$	$E_{[SEP]}$
+	+	+	+	+	+	+	+	+	+	+	+
Segment Embeddings	E_A	E_A	E_A	E_A	E_A	E_A	E_B	E_B	E_B	E_B	E_B
+	+	+	+	+	+	+	+	+	+	+	+
Position Embeddings	E_0	E_1	E_2	E_3	E_4	E_5	E_6	E_7	E_8	E_9	E_{10}

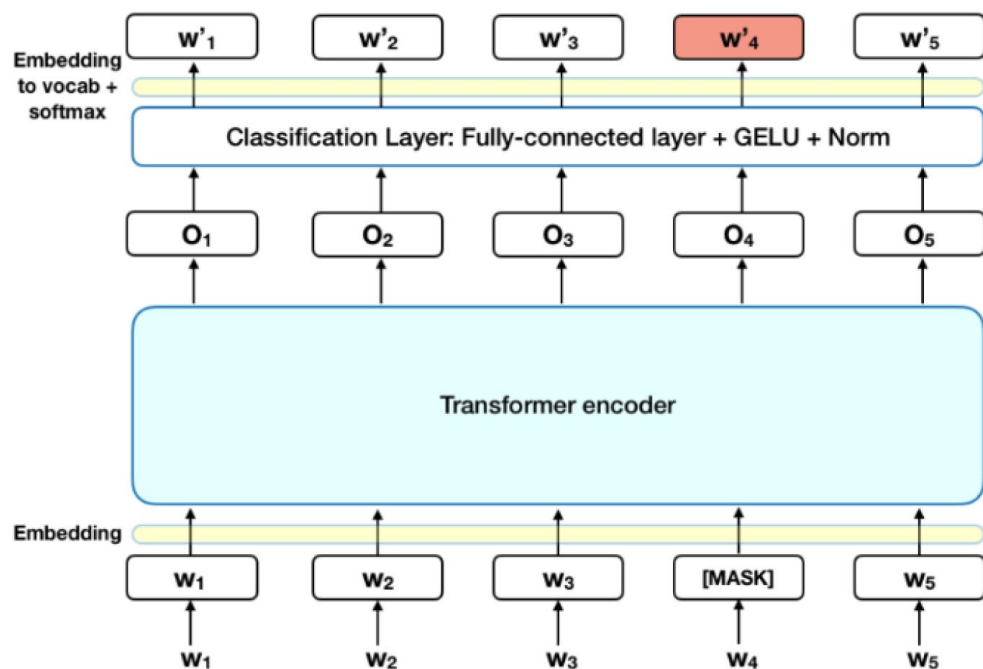
Input = Token + Segment + **Position**

Position

- Transformer에서 사용하던 positional embedding

1. After Transformer

Pre-trained



Output(Task): **Masked LM**와 NSP

Pretrain 과정에서 진행!

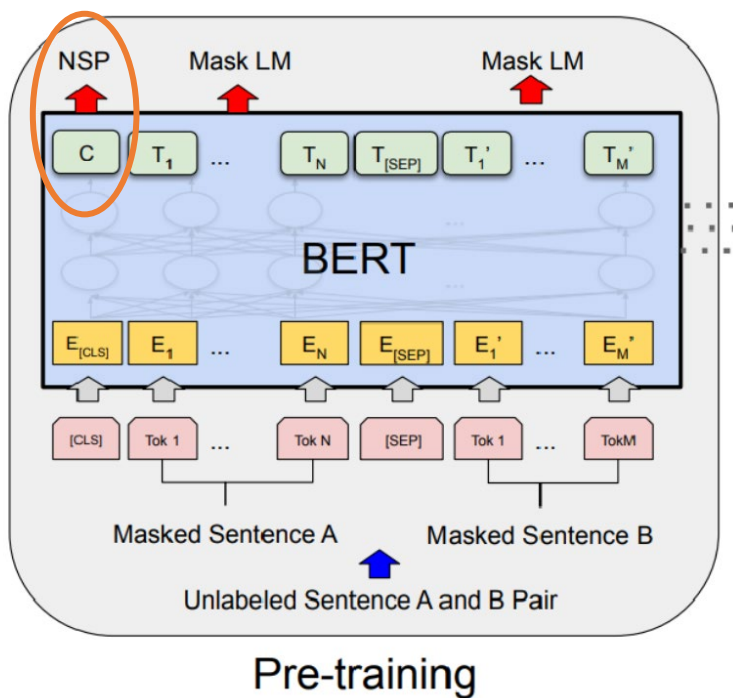
약 15% token "Masking" 그리고 예측

→ 앞뒤의 단어를 보고 학습하기에 **bidirectional!**

그 외에 몇가지 디테일은 Mask LM을 더 찾아보기로

1. After Transformer

Pre-trained



Output(Task): Masked LM와 NSP

Pretrain 과정에서 진행!

[SEP]로 구분된 두 sen이 연속된 유무 예측
50%는 연속된 것, 50%는 비연속 문장들

→ 맥락을 이해하도록 학습

→ 일부 task들에 대해선 성능 개선!

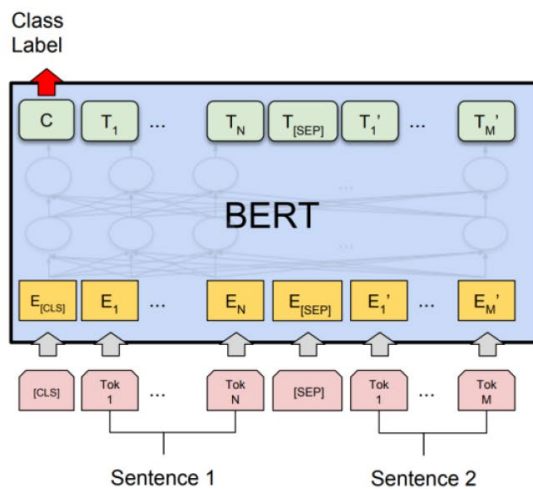
1. After Transformer

NLP - BERT

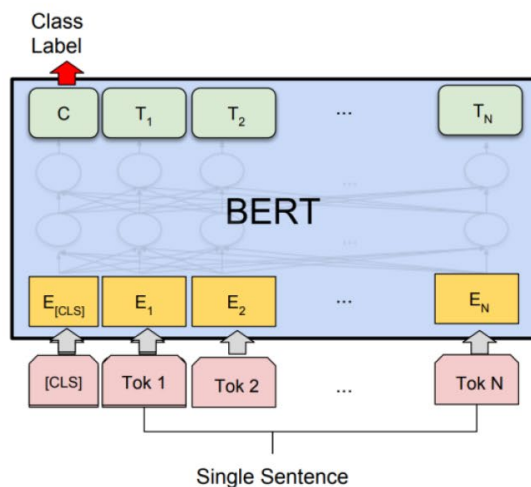


1. After Transformer

Fine-tuning



(a) Sentence Pair Classification Tasks:
MNLI, QQP, QNLI, STS-B, MRPC,
RTE, SWAG



(b) Single Sentence Classification Tasks:
SST-2, CoLA

Fine-tuning: 실제로 풀 수 있는 문제들!

→ **Classification**, QnA, Tagging

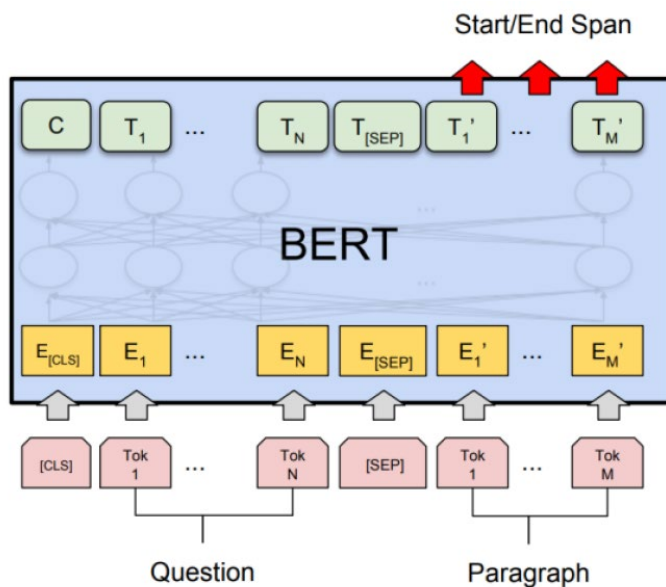
sequence-level

Input이 2문장인 경우: (예) 유사한지 분류

Input이 1문장인 경우: (예) 감성 분류!

1. After Transformer

Fine-tuning



(c) Question Answering Tasks:
SQuAD v1.1

Fine-tuning: 실제로 풀 수 있는 문제들!

→ Classification, QnA, Tagging

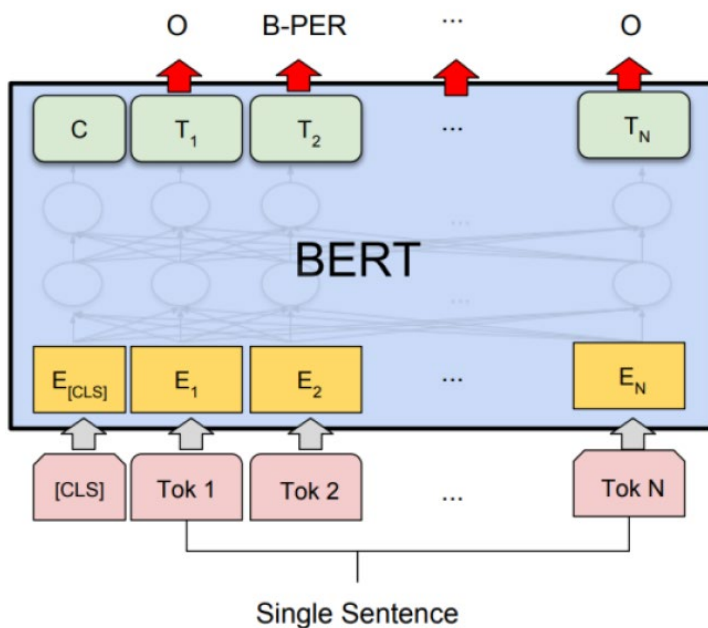
token-level

두 sen이 들어오는 경우:

Question과 Paragraph 쌍을 입력했을 때
답안 작성!

1. After Transformer

NLP – BERT



(d) Single Sentence Tagging Tasks:
CoNLL-2003 NER

Fine-tuning: 실제로 풀 수 있는 문제들!

→ Classification, QnA, **Tagging**

Token-level

한 sen이 들어오는 경우:

각각의 token마다 응답(ex. 형태소)를 예측

1. After Transformer

YONSEI Data Science Lab | DSL

NLP - BERT



동네 생활 데이터 통계 (학습 데이터에 사용)

Aa Name	≡ 데이터 수	📌 분류
동네 관련 이야기	7290	동네생활
중고거래 관련 신고	143	땀 이야기
중고거래 게시물	495	땀 이야기
기타(너무 짧은 글, 업체 홍보 글, 기능 문의 글, 제안, 등등)	395	땀 이야기

당근마켓에서 게시물들에 대한 분류 문제를 해결!

데이터 전처리부터 데이터셋 구축, 모델 구현 등
좋은 구현 예시

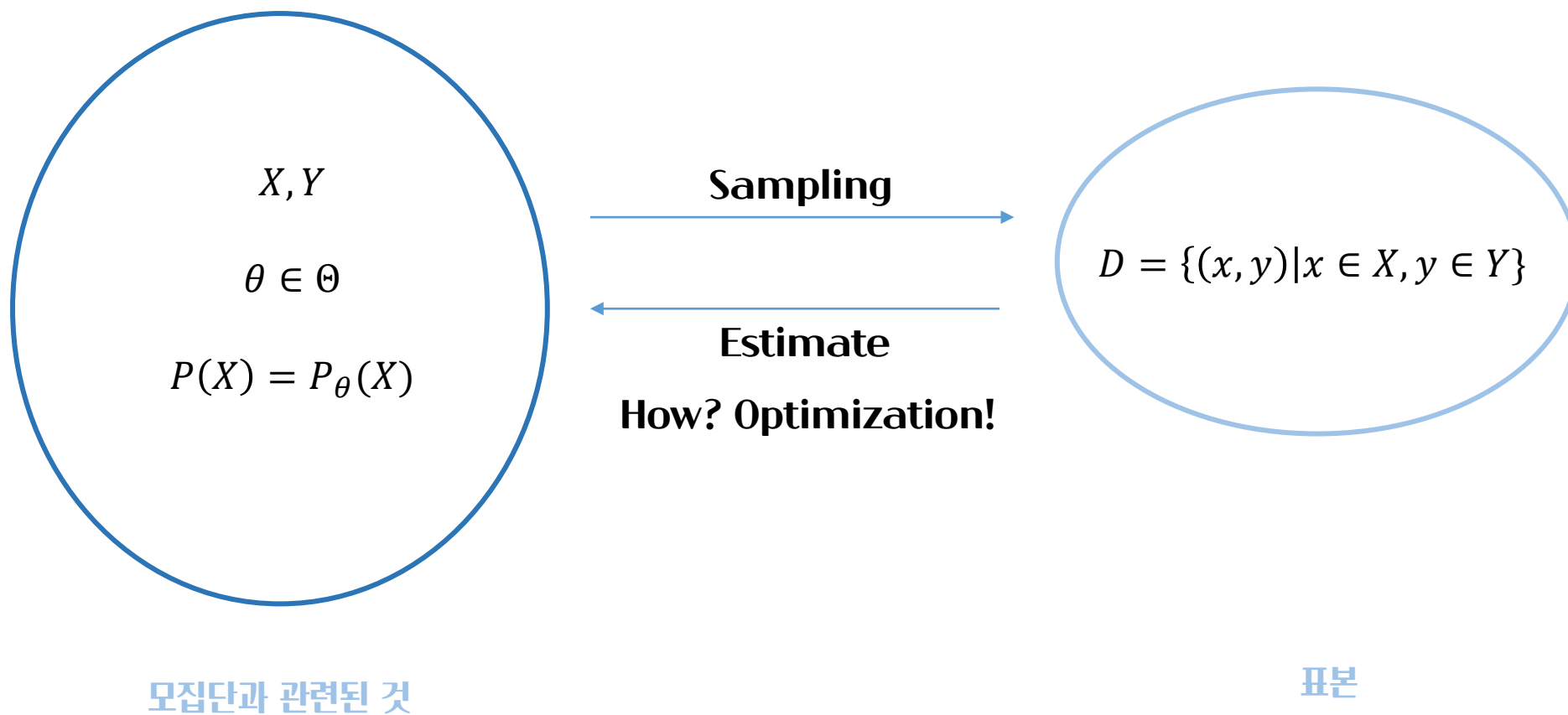
[동네이야기] 안녕하세요^^ 아버지 날을 맞아 부모님과 맛있는 식사를 하고 싶은데 마땅히 생각나는 맛... (노출범위2단계) - 제주 제주시 삼도1동 (14)	동네이야기 score: 1.0 버전: 20190418 🔄
[동네이야기] 이 꽃이를 공급해요 (노출범위2단계) - 제주 제주시 삼양동 (13)	동네이야기 score: 0.802 버전: 20190418 🔄
[동네이야기] 제주도청 운영시간 아세요? 여권만들러 가야하는데ㅠㅠ (노출범위2단계) - (9)	동네이야기 score: 1.0 버전: 20190418 🔄
[동네이야기] 운중동 청골송어집 문닫았나요? 전화가 없는번호라고 나오네요 혹시아시는분~ (노출범위인접동네) - 경기도 성남시 분당구 서현동 (5)	동네이야기 score: 1.0 버전: 20190418 🔄
[동네이야기] 중학생 졸업사진을 찍으려고 하는데요 혹시 캐릭터 의상 대여할수 있는곳좀 알려주세요~ (노출범위2단계) - 경기도 용인시 수지구 상현동 (15)	동네이야기 score: 0.979 버전: 20190418 🔄

<https://medium.com/daangn/딥러닝으로-동네생활-게시글-필터링하기-263cfe4bc58d>

2. Generative Model

YONSEI Data Science Lab | DSL

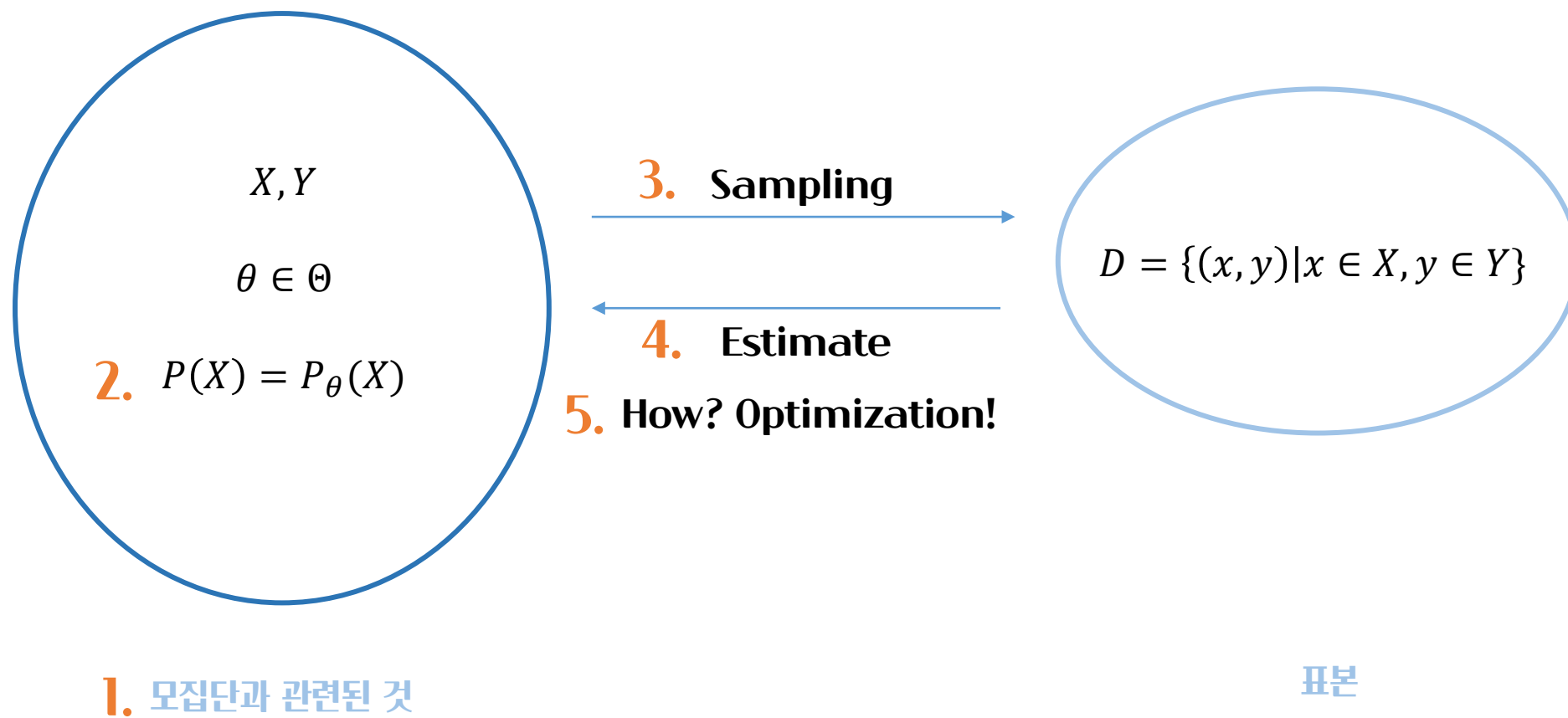
통계학이란



2. Generative Model

YONSEI Data Science Lab | DSL

통계학이란



2. Generative Model

YONSEI Data Science Lab | DSL

Bayesian

Parameter is not a constant value but a “Random Variable”!

$X|\theta \sim f(x|\theta)$: **conditional pdf** of X given $\Theta = \theta$

$\Theta \sim h(\theta)$: **Prior pdf** of $\Theta \leftarrow$ **표본이 없어서** 과거의 경험 혹은 주관에 따른 “사전 정보”의 역할

X' : Random sample from $f(x|\theta)$

x' : observation of X'

2. Generative Model

Bayesian

Parameter is not a constant value but a “Random Variable”!

$g_1(x, \theta) = f(x|\theta)h(\theta)$: Joint pdf of X and θ

$g_2(x) = \int g_1(x, \theta) d\theta$: Marginal pdf of X , “Evidence”

Conditional pdf of θ given $X' = x'$

“Posterior”: $k(\theta|x) = \frac{g_1(x,\theta)}{g_2(x)} = \frac{f(x|\theta)h(\theta)}{g_2(x)} \leftarrow$ 표본이 있으니 생각해본 파라미터의 분포?!

그럼 파라미터는 어떻게 “추정”? Bayes estimate, loss function, risk function...

2. Generative Model

YONSEI Data Science Lab | DSL

Bayesian

Prior & Posterior

결국 우리가 파라미터에 대해 모르거나, 경험적으로 알고 있는 정보가 있을 때(prior)
데이터를 보고 파라미터의 값 분포를 업데이트 해본다(posterior)!

Conjugate family of distribution

이론적으로, Prior 분포에 따라서 posterior의 분포가 어느정도 결정된다!

	Distribution		
$f(x \theta)$	Poisson	Binomial	Normal
$h(\theta)$	Gamma	Beta	Normal
$k(\theta x)$	Gamma	Beta	Normal

2. Generative Model

Supervised vs Unsupervised “Learning”

Supervised: Classification, Regression, Object Detection, ...

Data: $D = \{(x, y) | x \in X, y \in Y\}$

Model: $y = \mathcal{M}_\theta(x)$

→ Mapping!

올바르게 mapping하도록 학습



In

Whale!

Out

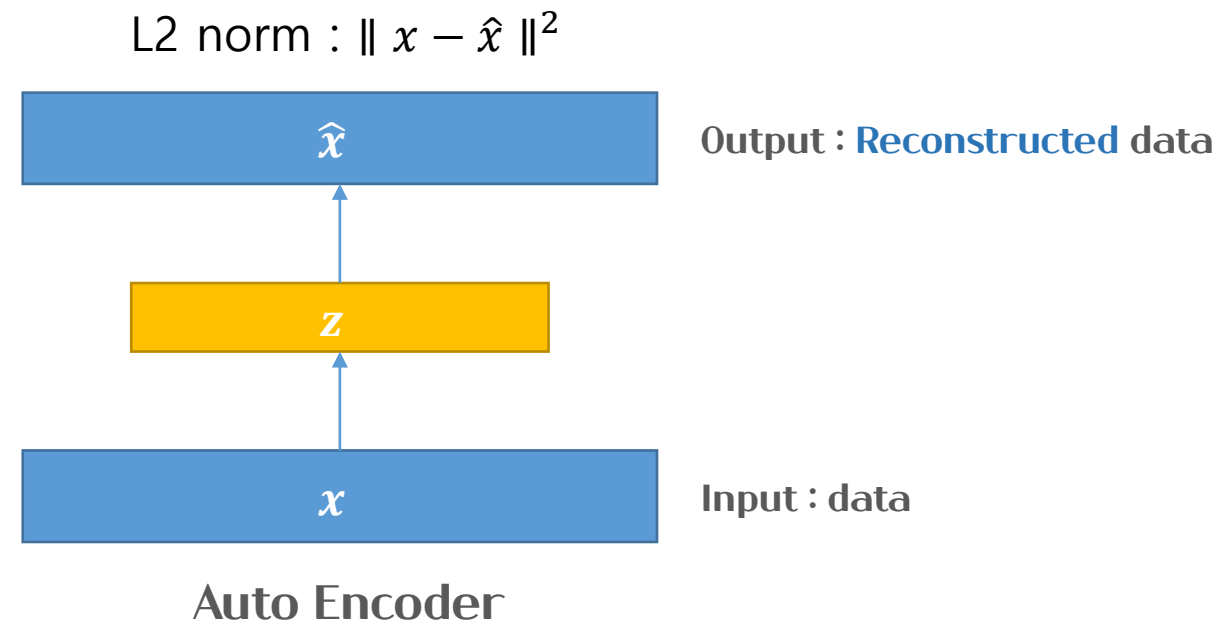
2. Generative Model

Supervised vs Unsupervised “Learning”

Unsupervised: Clustering, Feature learning, Density Estimation ...

Data: $D = \{x \mid x \in X\}$

Model: Data structure!



2. Generative Model

Cf) Auto Encoder(AE)

Encoder: 더 작은 차원의 벡터로 표현하는 방법을 학습!

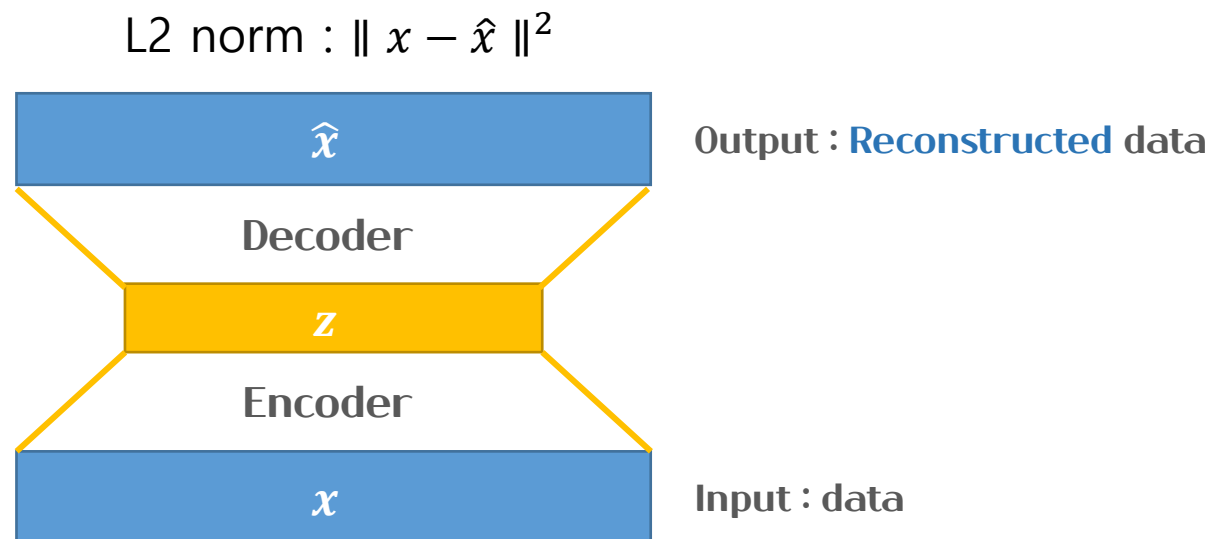
→ Feature Learning, Dimensionality Reduction

→ How? CNN, FC, ...

Decoder: 압축된 정보를 원본으로 복원 생성!

→ Generate...? NO!

→ How? TransposeCNN, FC, ...



관련 Task: Semantic segmentation → U-net

2. Generative Model

Encoder as a feature extractor!

Cf) Auto Encoder(AE)

Encoder can "represent" raw data containing data structure information!

Encoder: 더 작은 차원의 벡터로 표현하는 방법을 학습!

Ex) Feature Extractor or initial value!

→ Feature Learning, Dimensionality Reduction

→ How? CNN, FC, ...

CNN + Classifier(Softmax Classifier, SVD etc)

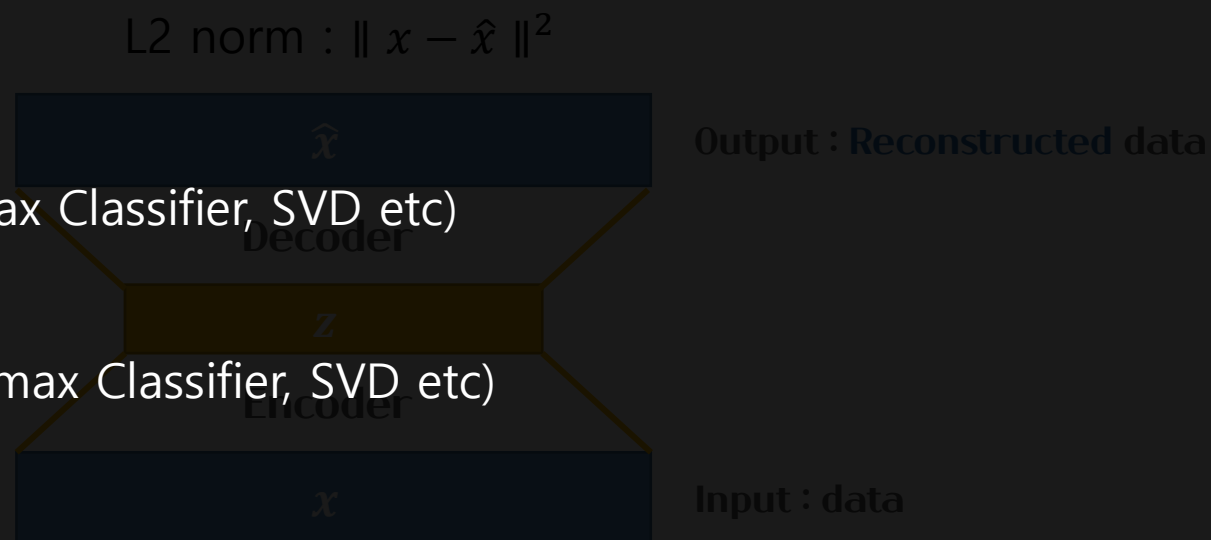
Decoder: 압축된 정보를 원본으로 복원 생성!

Encoder + Classifier(Softmax Classifier, SVD etc)

Feature Learning, Dimensionality Reduction

→ How? TransposeCNN, FC, ...

BUT, cannot **sample** "new" data.



WHY?

관련 Task: Semantic segmentation → U-net

2. Generative Model

Discriminative vs Generative “Model”

Discriminative: Learn $P(y|x)$

레이블들 간의 확률을 계산하는 모델이며
사진에 대한 확률이 아니다!

그렇기에 이상한 이미지가 들어온다면,
분류하지 못할 수 있다.



Given

Whale!
Dog..?
Bird?

prob

2. Generative Model

Discriminative vs Generative “Model”

Generative: Learn $P(x)$

등장하지 않을 사진은 낮은 확률을 가지는 분포를 계산하는 모델!
게다가, 분포를 안다면 x 의 **sampling**도 가능하지 않을까?

→ 확률분포로 표현해야 샘플링이 가능!

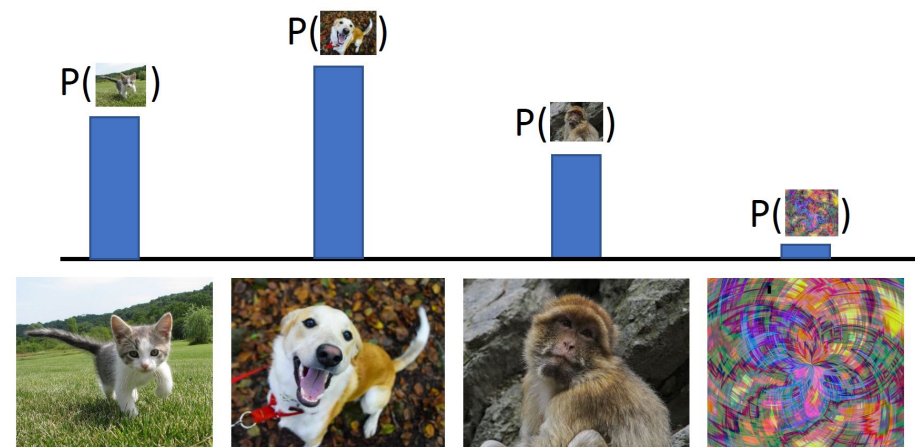
Discriminative Model

$$P(x|y) = \frac{P(y|x)}{P(y)} P(x)$$

(conditional)
Generative Model

Prior over label

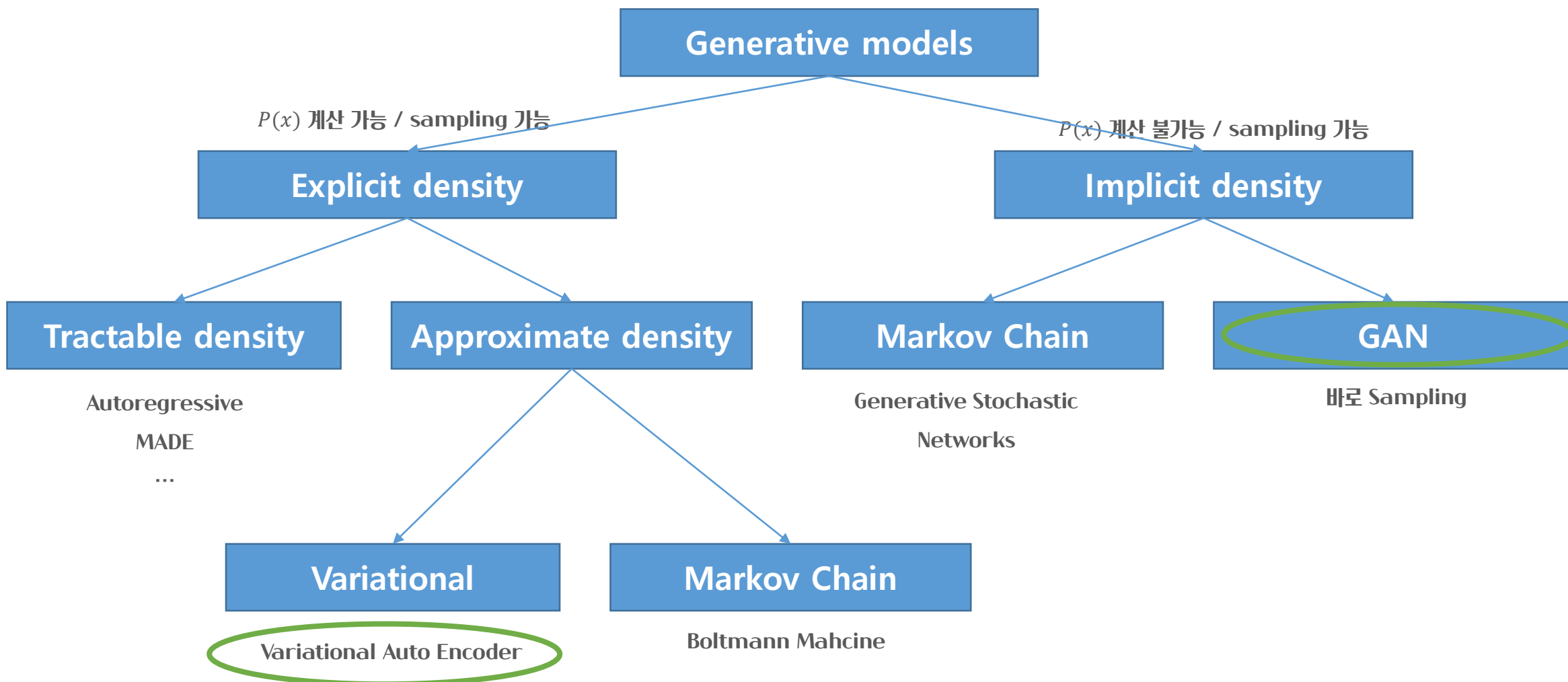
(Unconditional)
Generative Model



2. Generative Model

YONSEI Data Science Lab | DSL

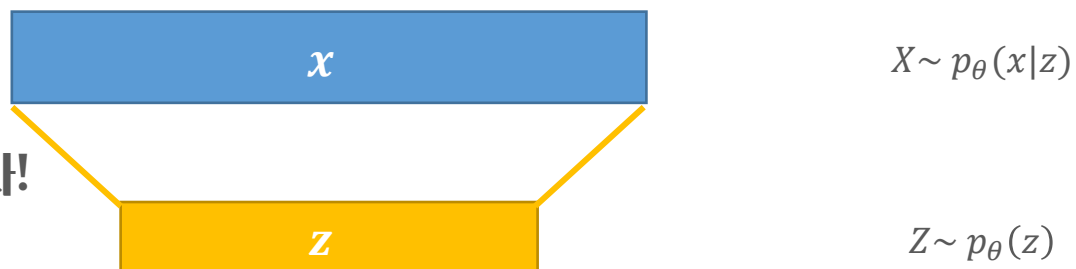
Models



2. Generative Model

Models-VAE

1. Latent feature를 잘 표현해본다
2. Given Latent feature, X가 등장할 확률을 구해보자!



Prior=사전정보 없는데…?

그렇다면 Gaussian 분포를 따른다고 해보자!

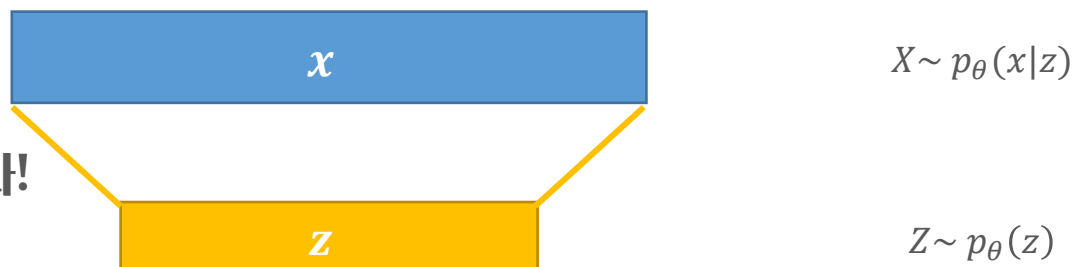
$$\underbrace{P(x|z)}_{\text{(conditional) Generative Model}} = \frac{P(z|x)}{\underbrace{P(z)}_{\text{Prior over latent feature}}} \underbrace{P(x)}_{\text{(Unconditional) Generative Model}}$$

2. Generative Model

YONSEI Data Science Lab | DSL

Models-VAE

1. Latent feature를 잘 표현해본다
2. Given Latent feature, X가 등장할 확률을 구해보자!



Conjugate family

Prior를 Gaussian가정 해봤으니,

Posterior도 Gaussian을 따른다고 해보자!

Mean : $\mu_{x|z}$

Covariance: $\Sigma_{x|z}$

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

(Unconditional) Generative Model

(conditional) Generative Model

Prior over latent feature

Discriminative Model?!

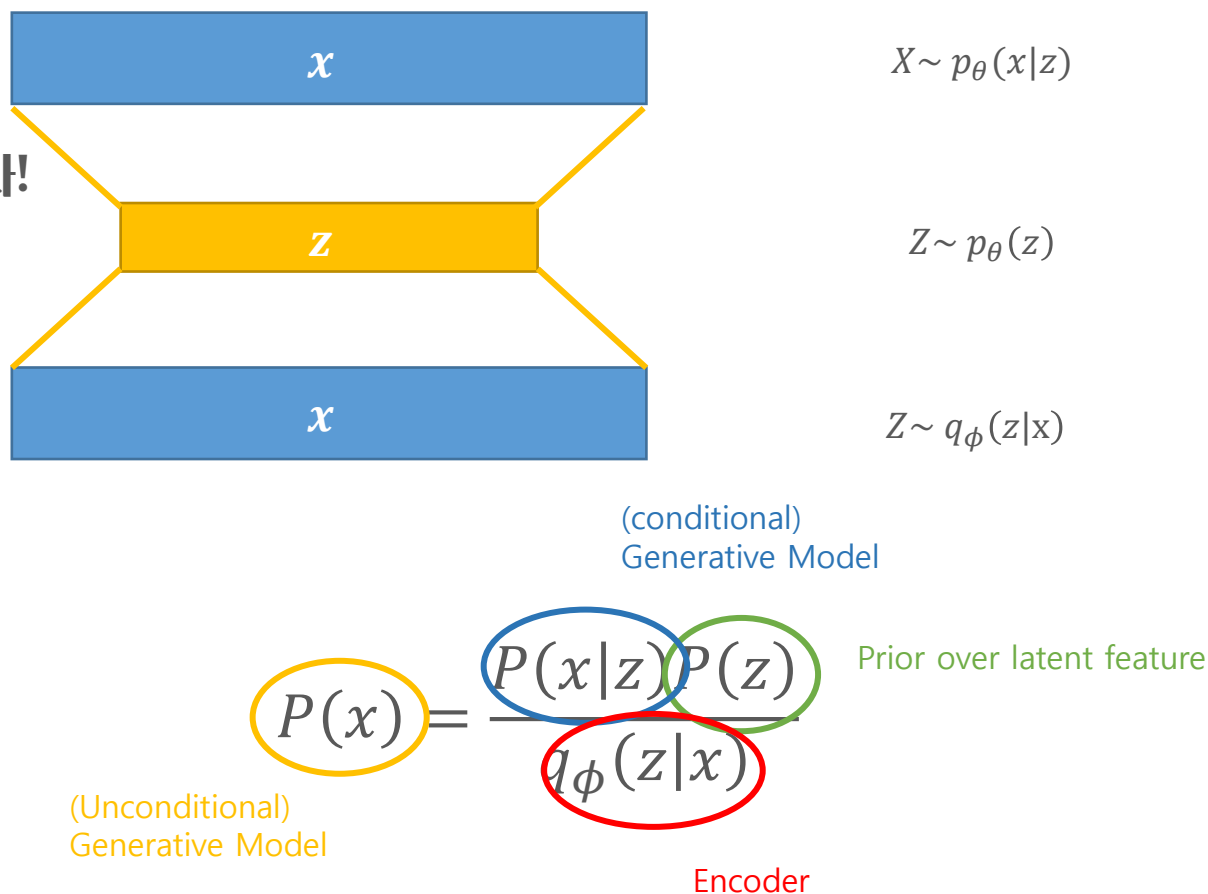
2. Generative Model

Models-VAE

1. Latent feature를 잘 표현해본다
2. Given Latent feature, X가 등장할 확률을 구해보자!

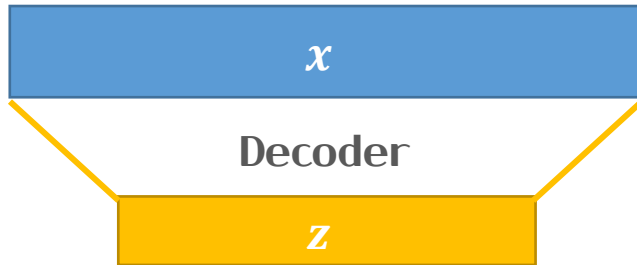
(conditional) Generative Model = Decoder

Discriminative model = Encoder(proxy)

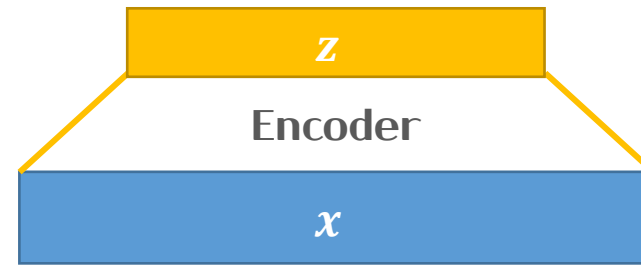


2. Generative Model

Models-VAE



$$p_{\theta}(x|z) = N(\mu_{x|z}, \Sigma_{x|z})$$



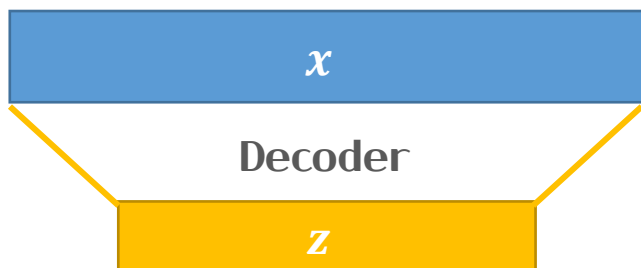
$$q_{\phi}(z|x) = N(\mu_{z|x}, \Sigma_{z|x})$$

$$\rightarrow \text{Approximate } p_{\theta}(x) \cong \frac{p_{\theta}(x|Z)p(z)}{q_{\phi}(Z|x)}$$

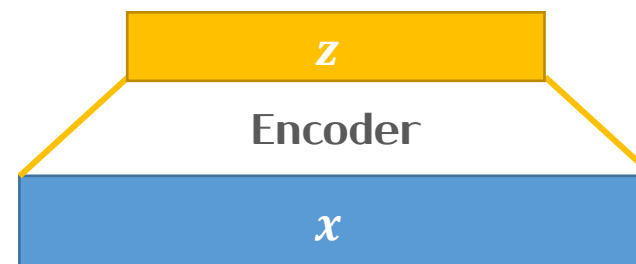
2. Generative Model

YONSEI Data Science Lab | DSL

Models-VAE



$$p_{\theta}(x|z) = N(\mu_{x|z}, \Sigma_{x|z})$$



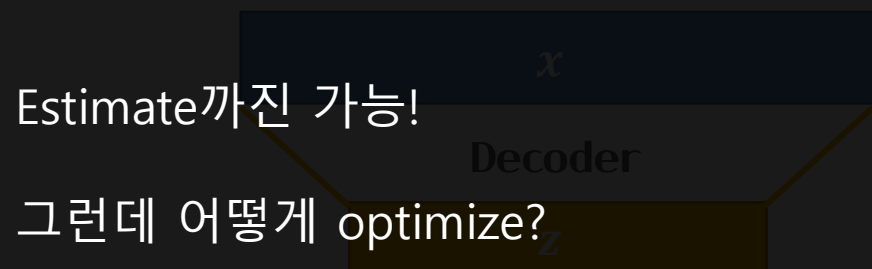
$$q_{\phi}(z|x) = N(\mu_{z|x}, \Sigma_{z|x})$$

이것만 계산해서 알아내면 된다!

2. Generative Model

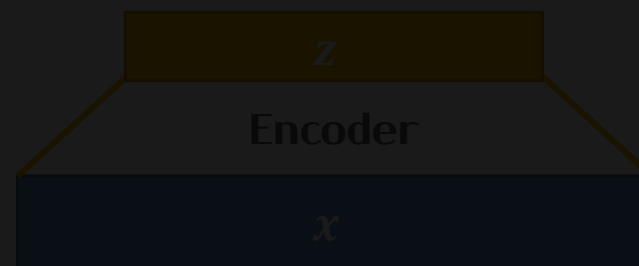
YONSEI Data Science Lab | DSL

Models-VAE



Estimate까진 가능!

그런데 어떻게 optimize?



이어서 자세한 식을 드리지만, 이번 시간에서는 그 개념과 통계적인 해석만을 가져봅시다.

$$p_{\theta}(x|z) = N(\mu_{x|z}, \Sigma_{x|z})$$

$$q_{\phi}(z|x) = N(\mu_{z|x}, \Sigma_{z|x})$$

$$\rightarrow \text{Approximate } p_{\theta}(x) \cong \frac{p_{\theta}(x|Z)p(z)}{p_{\theta}(x|Z)}$$

2. Generative Model

$$\begin{aligned}\log p_{\theta}(x) &= \log \frac{p_{\theta}(x|z)p(z)}{p_{\theta}(z|x)} = \log \frac{p_{\theta}(x|z)p(z)q_{\phi}(z|x)}{p_{\theta}(z|x)q_{\phi}(z|x)} \\ &= \log p_{\theta}(x|z) - \log \frac{q_{\phi}(z|x)}{p(z)} + \log \frac{q_{\phi}(z|x)}{p_{\theta}(z|x)} \\ &= E_z[\log p_{\theta}(x|z)] - E_z \left[\log \frac{q_{\phi}(z|x)}{p(z)} \right] + E_z \left[\log \frac{q_{\phi}(z|x)}{p_{\theta}(z|x)} \right]\end{aligned}$$

2. Generative Model

$$= E_{z \sim q_{\phi}(z|x)} [\log p_{\theta}(x|z)] - D_{KL} \left(q_{\phi}(z|x), p(z) \right) + D_{KL}(q_{\phi}(z|x), p_{\theta}(z|x))$$

Data Reconstruction

Divergence btw
samples from Encoder
& samples from prior

Divergence btw
Encoder & counterpart

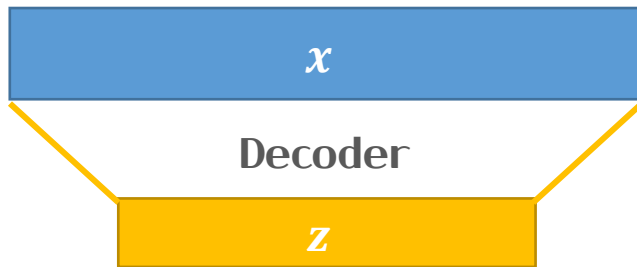
$$\log p_{\theta}(x) \geq E_{z \sim q_{\phi}(z|x)} [\log p_{\theta}(x|z)] - D_{KL} \left(q_{\phi}(z|x), p(z) \right)$$

Gap btw ELBo & Evidence

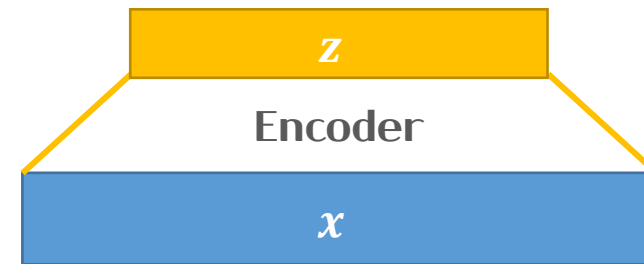
Evidence Lower Bound : ELBo

2. Generative Model

Models-VAE



$$p_{\theta}(x|z) = N(\mu_{x|z}, \Sigma_{x|z})$$



$$q_{\phi}(z|x) = N(\mu_{z|x}, \Sigma_{z|x})$$

Decoder는 z 를 주면 $p_{\theta}(x|z)$ 를 계산해줌!

2. Generative Model

$$= E_{z \sim q_{\phi}(z|x)} [\log p_{\theta}(x|z)] - D_{KL} \left(q_{\phi}(z|x), p(z) \right) + D_{KL}(q_{\phi}(z|x), p_{\theta}(z|x))$$

Data Reconstruction

Divergence btw
samples from Encoder
& samples from prior

Divergence btw
Encoder & counterpart

$$\log p_{\theta}(x) \geq \boxed{E_{z \sim q_{\phi}(z|x)} [\log p_{\theta}(x|z)]} - D_{KL} \left(q_{\phi}(z|x), p(z) \right)$$

Encoder로 sample z 를 생성 후 평균값 취하기!

2. Generative Model

$$= E_{z \sim q_{\phi}(z|x)} [\log p_{\theta}(x|z)] - D_{KL} \left(q_{\phi}(z|x), p(z) \right) + D_{KL}(q_{\phi}(z|x), p_{\theta}(z|x))$$

Data Reconstruction

Divergence btw
samples from Encoder
& samples from prior

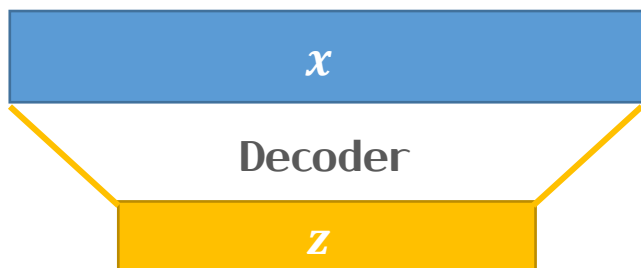
Divergence btw
Encoder & counterpart

$$\log p_{\theta}(x) \geq E_{z \sim q_{\phi}(z|x)} [\log p_{\theta}(x|z)] - D_{KL} \left(q_{\phi}(z|x), p(z) \right)$$

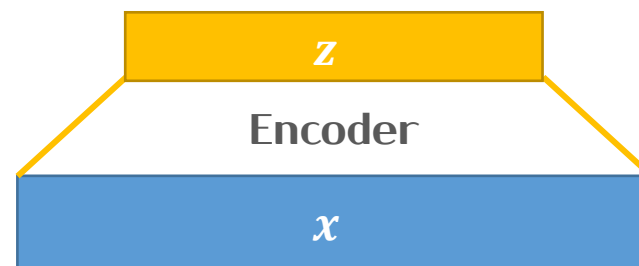
$q_{\phi}(z|x)$ 가 diagonal covariance
 $p(z)$ 가 std Normal 분포라면 closed form

2. Generative Model

Models-VAE



$$p_{\theta}(x|z) = N(\mu_{x|z}, \Sigma_{x|z})$$



$$q_{\phi}(z|x) = N(\mu_{z|x}, \Sigma_{z|x})$$

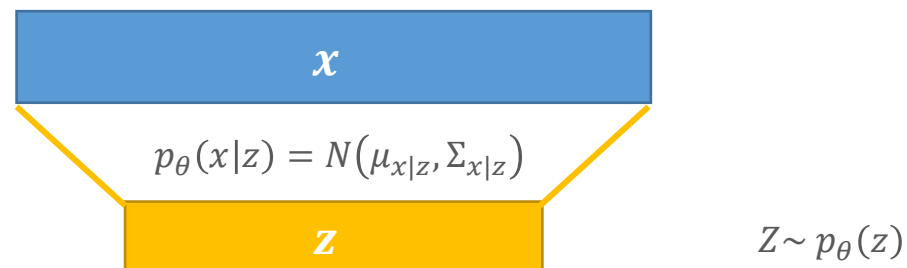
$$\rightarrow \text{Approximate } p_{\theta}(x) \cong \frac{p_{\theta}(x|Z)p(z)}{q_{\phi}(Z|x)}$$

2. Generative Model

Models-VAE

Sampling(생성) 법

1. Prior of z (Gaussian)으로 표본 생성 : $p_{\theta}(z)$
2. Decoder로 posterior의 Parameter 계산: $\mu_{x|z}, \Sigma_{x|z}$
3. Posterior 분포로 샘플링: $p_{\theta}(x|z)$



그렇지만,,, 좋은 퀄리티의 sample이 아니었다!(blurred...)

게다가 Evidence 자체가 아닌, “lower bound”를 maximize를 한 것

2. Generative Model

YONSEI Data Science Lab | DSL

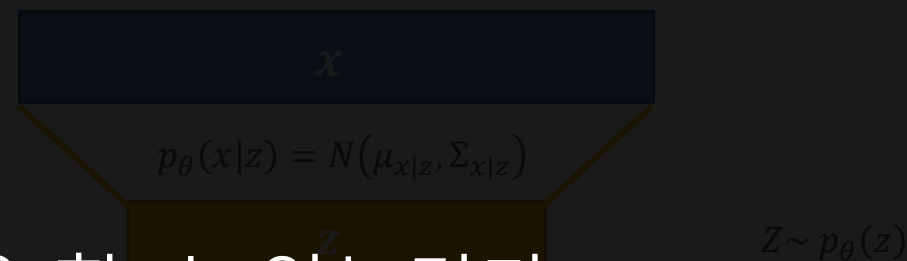
Models-VAE

Sampling(생성) 법

1. Prior of z (Gaussian)으로 표본 생성 : $p_{\theta}(z)$ 그런데...

2. Decoder로 posterior의 Parameter 계산: $\mu_{x|z}, \Sigma_{x|z}$

3. Posterior 분포로 샘플링 : $p_{\theta}(x|z)$ 꼭 확률을 계산해야 sampling을 할 수 있는건가?

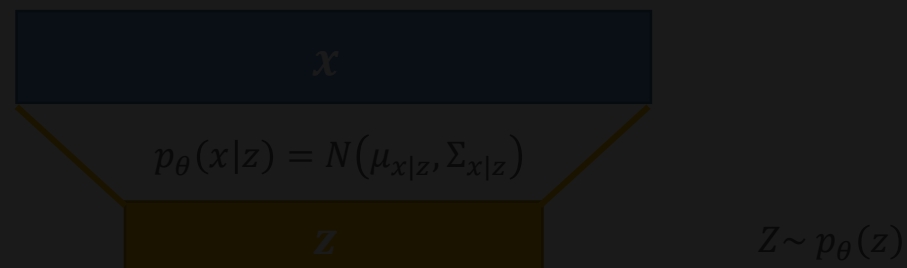


2. Generative Model

Models-VAE

Sampling(생성) 법

1. Prior of z (Gaussian)으로 표본 생성 : $p_{\theta}(z)$
2. Decoder로 posterior의 Parameter 계산 **다음주에!**
3. Posterior 분포로 샘플링: $p_{\theta}(x|z)$



고려대학교 산업경영공학부 DSBA 연구실 <https://www.youtube.com/watch?v=lwtexRHoWG0>

<https://github.com/pilsung-kang/Text-Analytics>

https://www.youtube.com/watch?v=MN_lSncZBs

<https://medium.com/daangn/딥러닝으로-동네생활-게시글-필터링하기-263cfe4bc58d>

Stanford CS224N: <https://www.youtube.com/watch?v=knTc-NQsJkA>

이기복 교수님 “통계적 머신러닝” 강의안

강승호 교수님 “이론통계학(1)” 강의안

임종호 교수님 “수리통계학(2) 강의

5기 한영웅님 AutoEncoder 세션 자료

왜 $q_\phi(z|x)$ 가 diagonal covariance이고 $p(z)$ 가 std Normal 분포라면 closed form인지

$$\begin{aligned} -D_{KL} \left(q_\phi(z|x), p(z) \right) &= \int_Z q_\phi(z|x) \log \frac{p(z)}{q_\phi(z|x)} dz \\ &= \int_Z N(z; \mu_{z|x}, \Sigma_{z|x}) \log \frac{N(z; 0, I)}{N(z; \mu_{z|x}, \Sigma_{z|x})} dz \\ &= \frac{1}{2} \sum_{j=1}^J \left(1 + \log \left((\Sigma_{z|x})_j^2 \right) - (\mu_{z|x})_j^2 - (\Sigma_{z|x})_j^2 \right) \end{aligned}$$