

Lecture5: Latent Variable Models

■ Ki Pyo	완료
■ 선택	VAE

1. Latent Variable Models: Motivation → What & Why

Latent Variable Model

what : 보이지 않지만 X 에 영향을 주는 요소를 확률 변수 z (관측되지 않은 원인)로 나타내어 x 와 z 의 joint probability를 학습하는 모델

why : 보이지 않는 요소 또한 고려하기 위해!

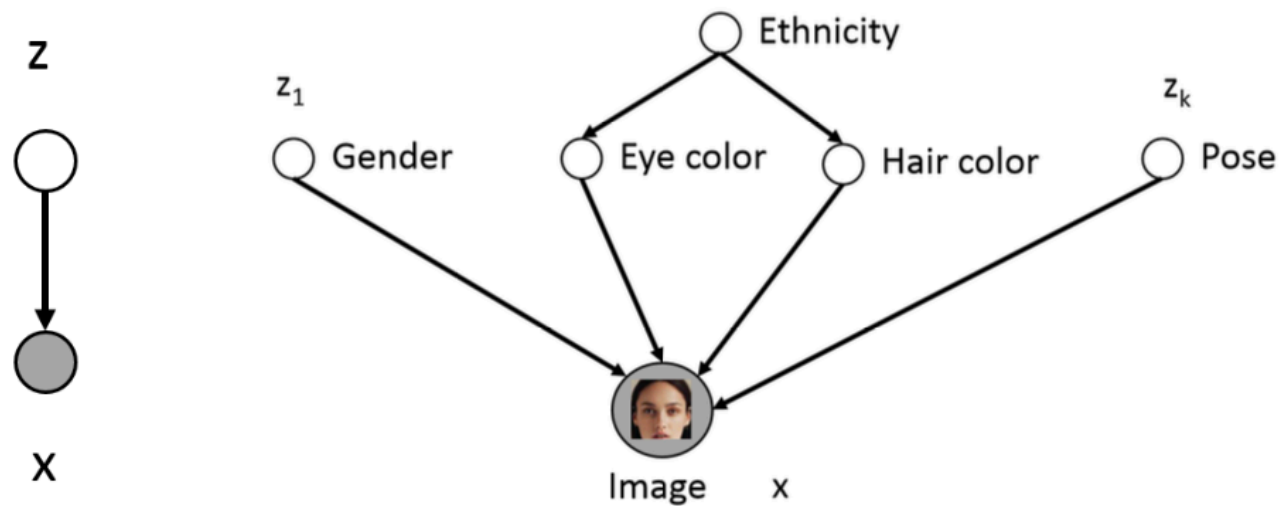
EX)



- 우리가 관측하고 있는 것들, 모집단 $P(x)$ 에 대한 샘플 데이터셋은 pixel의 값들임 -> X
- Latent Factors For Variation, z : 관측되지는 않지만 픽셀 값 X 에 영향을 미치는 요소들이 많을 것임.
→ 예를 들어 사람들의 사진이라면, 사람의 성별, 눈 색깔 등과 같은 보이지 않는 요소들이 현재 관측된 픽셀 값들에 영향을 줄 수 있음
- Latent Variable Model은 이런 관측되지 않은, X 의 variability에 영향을 주는 요소들을 Z 라는 잠재변수로 취급하여 X 와 Z 의 분포를 모델링

2. Latent Variable Models: Motivation: Cont'd

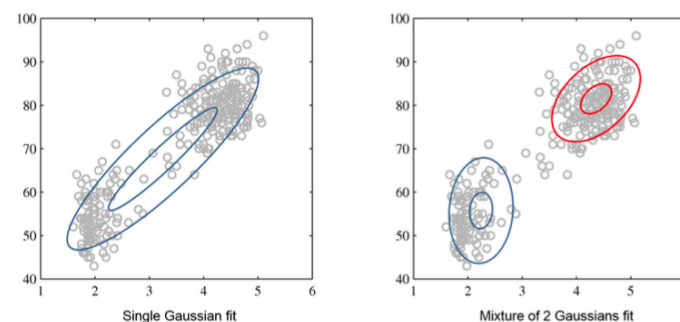
- Latent Variable Model에서 X , Z 둘의 관계를 Bayesian Network로 표현하면 다음과 같다
- $z \rightarrow x$ 의 ordering을 따르기에 joint probability $P(x, z)$ 를 $P(z)P(x|z)$ 로 모델링하고 각각의 확률 분포를 Learning할 것임
- $P(z)$ 는 Learnable한 대상이 될 수도 있고 Fix 할 수도 있음(eg. VAE)



- 잠재적인 요소가 여러개라면 $z_1 \sim z_k$ 는 하나의 Z 로 표현될 것임 : Bunch of Latent Factors for variation that will be helpful for describing different types of images that we have access to.

3. Advantages of Latent Variable model

1. Less Variability in each cluster, clustered by Z

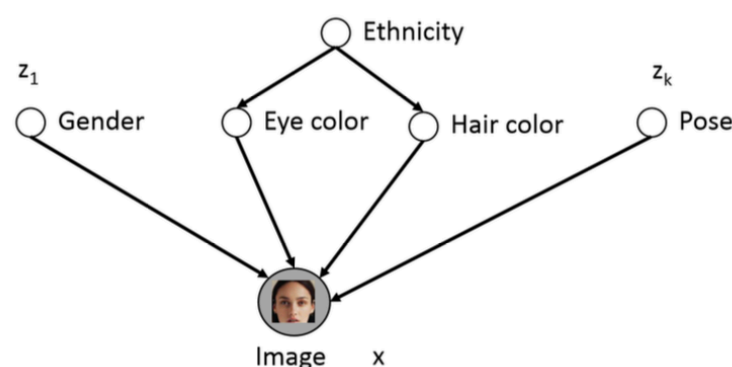


- 이미지들 간의 Cluster를 잠재적으로 잘 구별해내는 잠재적인 feature z 를 뽑아낸다면 그때의 각각 Clustering된 X distribution은 capture해야하는 variation이 적기에 더 모델링하기 편함 $P(x|z)$
 \Rightarrow "쉽게 말해서 Z 로 인해 X 들이 Clustering이 될 수 있고 clustering된 X 에 대해서 분포를 모델링하는 것이 Latent Variable z 없이, Clustering없이 썬으로 X 에 대한 분포 $P(x)$ 를 모델링하는 것보다 편함"

2. Unsupervised Representation Learning

- 관측 변수 X 에 대해 Label이 없는 상태로 $P(x, z)$ 를 학습
- Latent Variable Model을 통해 Z, X 의 분포를 잘 모델링 및 학습했다면 z 가 X 를 구별할 좋은 feature들을 가지고 있을 수도 있음 representation extraction에도 유용할 것임 $\rightarrow P_{\theta}(x, z)$ 를 모델링하니 $P_{\theta}(z|x)$ 를 구할 수 있음
- $P(z|x)$ 로 x 에 대한 feature와 같은 역할을 하는 z 를 compute할 수 있다

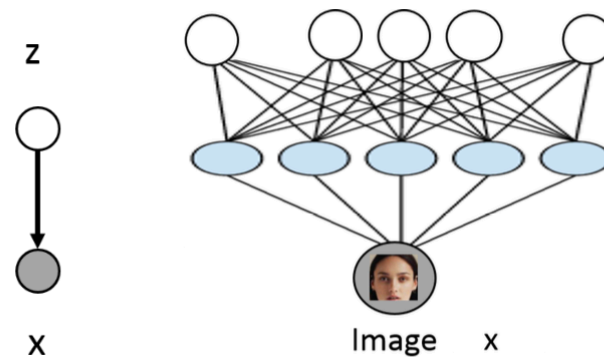
4. Deep Latent Variable Models



- 만약 위치럼
 - (1) 어떤 Latent Factor들이 X 의 Variability에 영향을 줄 수 있는지 명확하게 알 수 있고,
 - (2) Latent Factor간의 관계도 알 수 있다면,

Bayesian Network로 표현하여 X에 대한 Distribution을 명시적으로 나타낼 수 있음(위의 경우 각 Latent Factor끼리는 서로 어떤 영향을 주지 않음)

→ $P(x|z_1, z_2, \dots, z_k)$ 와 같이 !



- 1 $\mathbf{z} \sim \mathcal{N}(0, I)$
- 2 $p(\mathbf{x} | \mathbf{z}) = \mathcal{N}(\mu_{\theta}(\mathbf{z}), \Sigma_{\theta}(\mathbf{z}))$ where $\mu_{\theta}, \Sigma_{\theta}$ are neural networks

- 그러나 (1), (2)가 어떤지 모르는 경우가 대부분이니

Z factor개수도 무작위로 설정하고(보통 어떤 분포에서 나왔다고 가정)

이렇게 샘플링된 Z를 입력으로 받아 NN으로 Functional하게 Z factor간의 관계를 고려하여 X Distribution에 대한 파라미터를 출력하는 방식으로 X에 대한 Distribution을 모델링 및 학습을 진행

- 예를 들어 given z에 대해 X에 대한 분포가 Gaussian Distribution이라면 NN는 mean vector와 covariance matrix를 뽑아내면됨

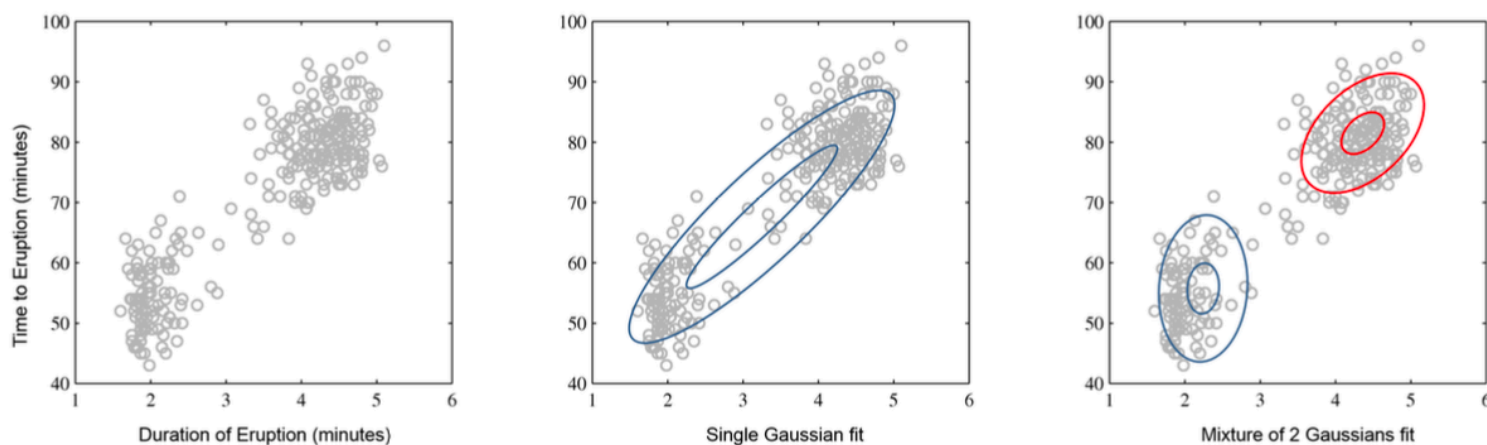
5. Summary

Goal of Latent Variable Model

- X에 대한 분포를 모델링 및 학습하고 싶은 것은 Generative Model의 목적과 같다.
- 다만 보이지 않은 X에 대한 Latent Factor of Variation Z 또한 고려하여 $P(x, z)$ 를 학습한다.
- $Z \rightarrow X$ 의 ordering을 고려하여 $P(x, z) = P(z)P(x|z)$ 로 분해하여 각각을 학습. X Sampling(generation)은 Z sampling → $X|Z$ sampling 순으로 진행

2 Advantages

- 구조 단순화:** Z를 통해 복잡한 $p(x)$ 를 여러 개의 단순한 $p(x|z)$ 로 나누어 모델링 가능(각각의 나누어진 영역은 X의 Variability가 적음)



- 표현 학습:** Z는 X를 잘 구분·요약하는 의미 있는 representation이 될 수 있음

6. Mixture Of Gaussians

- z는 이산적인 값을 가진다: Mixture의 type를 의미함

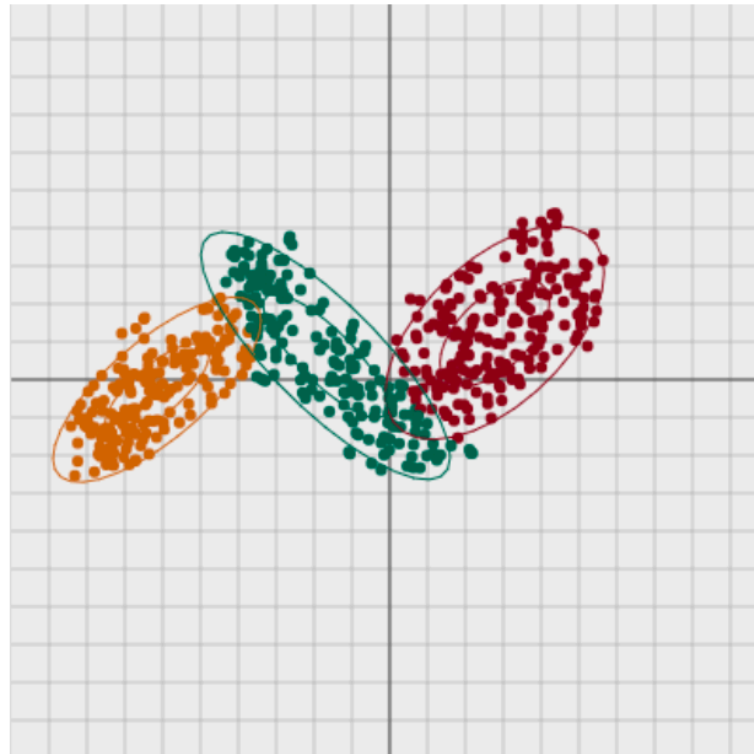
- $P(x|z=k)$ 는 각 z 에 대한 X 의 Distribution이고 Gaussian Distribution이다.
- $z \rightarrow x$ ordering을 따라 먼저 $P(z)$ 에서 z 를 샘플링 한 후에 $P(x|z)$ 에서 x 를 샘플링

Mixture of Gaussians:

- ① $\mathbf{z} \sim \text{Categorical}(1, \dots, K)$
- ② $p(\mathbf{x} | \mathbf{z} = k) = \mathcal{N}(\mu_k, \Sigma_k)$

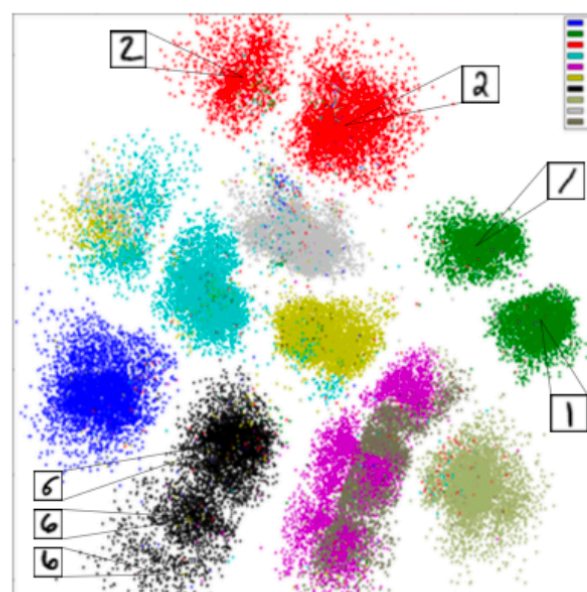
$P(z, x)$ 를 위와 같이 모델링

- z 를 입력으로 받아 functional 하게 x distribution의 파라미터를 뽑아내지는 않음 -> Deep Neural Network를 사용하지 않아 **Shallow Latent Variable Model**이라고 불림
- 각 Distribution에 대한 mean vector와 Covariance를 Functional하게 뽑아내지 않으므로 따로 저장해야함 (look-up table)
- Clustering : $P(z|x)$ 를 통해 각 x 가 어떤 z 에서 가장 나올 법한지를 알 수 있음



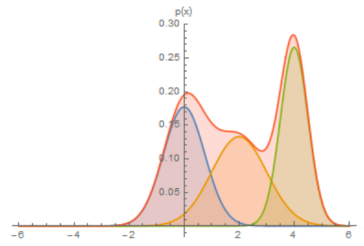
Shown is the posterior probability that a data point was generated by the i -th mixture component, $P(z = i|x)$

- Unsupervised Learning : 잘 훈련된 GMM의 각 z 는 X 를 잘 구별할 수 있는 feature의 역할을 할 수 있음 -> Z 도 cluster를 이룰 것임



7. Marginalization

Alternative motivation: Combine simple models into a more complex and expressive one

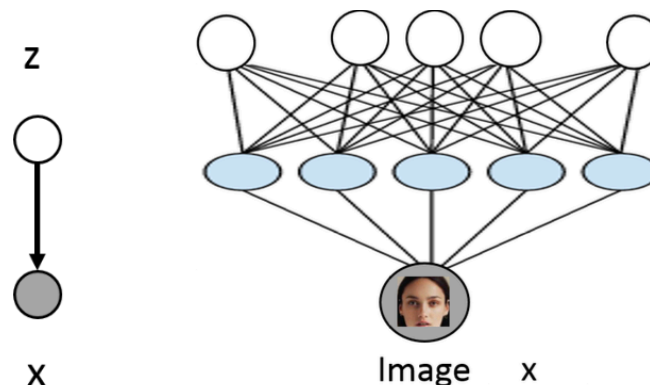


$$p(\mathbf{x}) = \sum_{\mathbf{z}} p(\mathbf{x}, \mathbf{z}) = \sum_{\mathbf{z}} p(\mathbf{z}) p(\mathbf{x} | \mathbf{z}) = \sum_{k=1}^K p(\mathbf{z} = k) \underbrace{\mathcal{N}(\mathbf{x}; \mu_k, \Sigma_k)}_{\text{component}}$$

- $P(\mathbf{x}, \mathbf{z})$ 를 학습하는 GMM에서 $P(\mathbf{x})$ density는 \mathbf{z} marginalization으로 구할 수 있음
- 각 \mathbf{z} 에 대한 $P(\mathbf{x} | \mathbf{z})$ 는 간단한 가우시안 분포이지만, $P(\mathbf{x}, \mathbf{z})$ 에서 전체적으로 \mathbf{x} 에 대한 Distribution $P(\mathbf{x})$ 은 Complex하고 Flexible한 모양을 띠며 → Latent Variable Model의 또다른 강점

8. VAE

- 1 $\mathbf{z} \sim \mathcal{N}(0, I)$
- 2 $p(\mathbf{x} | \mathbf{z}) = \mathcal{N}(\mu_{\theta}(\mathbf{z}), \Sigma_{\theta}(\mathbf{z}))$ where $\mu_{\theta}, \Sigma_{\theta}$ are neural networks
 - $\mu_{\theta}(\mathbf{z}) = \sigma(A\mathbf{z} + c) = (\sigma(a_1\mathbf{z} + c_1), \sigma(a_2\mathbf{z} + c_2)) = (\mu_1(\mathbf{z}), \mu_2(\mathbf{z}))$
 - $\Sigma_{\theta}(\mathbf{z}) = \text{diag}(\exp(\sigma(B\mathbf{z} + d))) = \begin{pmatrix} \exp(\sigma(b_1\mathbf{z} + d_1)) & 0 \\ 0 & \exp(\sigma(b_2\mathbf{z} + d_2)) \end{pmatrix}$
 - $\theta = (A, B, c, d)$
- 3 Even though $p(\mathbf{x} | \mathbf{z})$ is simple, the marginal $p(\mathbf{x})$ is very complex/flexible



- GMM은 \mathbf{z} 가 이산적인 값을 가진 것에 비해 VAE에서는 \mathbf{z} 에 대한 Prior가 연속확률분포를 따름 $\mathcal{N}(0, I)$. \mathbf{z} 가 discrete한 scalar가 아니라 \mathbf{z} 가 여러 개의 factor로 이루어져 있는 continual vector임.
- $P(\mathbf{z})$ 를 따르는 따라서 무수히 많은 \mathbf{z} vector에 대해 Decoder가 $P(\mathbf{x} | \mathbf{z})$ Gaussian Distribution의 파라미터를 뱉어내는 구조임. GMM과는 다르게 NN을 이용하여 sampling된 \mathbf{z} 의 각 factor 간의 관계를 고려하여 \mathbf{x} 의 distribution parameter를 출력함 → Mixture of an infinite number of Gaussians, **Deep Latent Variable Model**
- Decoder Output : mean vector and diagonal covariance matrix
- 이 역시 각각의 $P(\mathbf{x} | \mathbf{z})$ 는 간단한 Gaussian Distribution으로 모델링하지만 $P(\mathbf{x}, \mathbf{z})$ 에서 marginalization으로 구한 $P(\mathbf{x})$ distribution은 complex & flexible함.

9. What's the limit of Latent Variable Model

그러나 Training 과정 중에 \mathbf{x} 에 대한 Likelihood Estimation 계산 과정이 매우 복잡하고 intractable할 수도 있다. $P(\mathbf{x}, \mathbf{z})$ 에서 \mathbf{z} marginalization 때문에(\mathbf{z} 는 관측되지 않았기에, \mathbf{x} 생성에 관여한 \mathbf{z} 의 모든 가능한 값을 check해야함) fully observed된 \mathbf{x} 만을 다루는 autoregressive model보다 likelihood estimation이 복잡함.

10. Marginal Likelihood



- 이미지의 절반은 관측된 값 X 이고 절반은 관측되지 않은 알 수 없는 Z 값일 때 이에 대한 Latent generative model을 모델링한다고 하면 $P(x, z)$ 를 학습하게 될 것이다.
- 학습된 $P(x, z)$ 에서 $P(x)$ 를 구하려면 z 에 대한 marginalization이 필요하다.
- 위에서는 예시로 이미지의 가려진 알 수 없는 부분을 z 라고 표현하였는데, 이때 marginalization을 한다는 것은 가능한 모든 z 에 대해 $P(x, z)$ 를 계산 후 더하는 것과 같다.
 z 는 all white인 경우도 있을 테고 z 는 all black인 경우도 있을 테고 등등 (Need to consider all possible ways to complete the image)
- 심지어 z 가 highdimensional한 feature라면 + continuous한 값이라면.. marginalization이 어려워 질 수 있음

11. Marginal Likelihood of X in VAE

z



x

A mixture of an infinite number of Gaussians:

- ① $z \sim \mathcal{N}(0, I)$
- ② $p(x | z) = \mathcal{N}(\mu_\theta(z), \Sigma_\theta(z))$ where $\mu_\theta, \Sigma_\theta$ are neural networks
- ③ Z are unobserved at train time (also called hidden or latent)
- ④ Suppose we have a model for the joint distribution. What is the probability $p(X = \bar{x}; \theta)$ of observing a training data point \bar{x} ?

$$\int_z p(X = \bar{x}, Z = z; \theta) dz = \int_z p(\bar{x}, z; \theta) dz$$

VAE에서는 z 가 연속형 벡터이므로 하나의 데이터 포인트 X 에 대한 VAE의 marginal likelihood 계산은 위와 같고, Z 에 대해 Marginalize하기 위한 computational cost가 엄청날 것임, Intractible함.

12. MLE Learning for Latent Variable Models

- Suppose that our joint distribution is

$$p(\mathbf{X}, \mathbf{Z}; \theta)$$

- We have a dataset \mathcal{D} , where for each datapoint the \mathbf{X} variables are observed (e.g., pixel values) and the variables \mathbf{Z} are never observed (e.g., cluster or class id.). $\mathcal{D} = \{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(M)}\}$.
- Maximum likelihood learning:

$$\log \prod_{\mathbf{x} \in \mathcal{D}} p(\mathbf{x}; \theta) = \sum_{\mathbf{x} \in \mathcal{D}} \log p(\mathbf{x}; \theta) = \sum_{\mathbf{x} \in \mathcal{D}} \log \sum_{\mathbf{z}} p(\mathbf{x}, \mathbf{z}; \theta)$$

- Evaluating $\log \sum_{\mathbf{z}} p(\mathbf{x}, \mathbf{z}; \theta)$ can be intractable. Suppose we have 30 binary latent features, $\mathbf{z} \in \{0, 1\}^{30}$. Evaluating $\sum_{\mathbf{z}} p(\mathbf{x}, \mathbf{z}; \theta)$ involves a sum with 2^{30} terms. For continuous variables, $\log \int_{\mathbf{z}} p(\mathbf{x}, \mathbf{z}; \theta) d\mathbf{z}$ is often **intractable**. Gradients ∇_θ also hard to compute.
- Need **approximations**. One gradient evaluation per training data point $\mathbf{x} \in \mathcal{D}$, so approximation needs to be cheap.

Intractible한 Marginal Likelihood of X 값을 근사할 방법이 필요함

13. Approximation 1: Naive Monte Carlo From Uniform Distribution of Z

z를 uniform한 분포로부터 sampling하여 근사하는 방법 -> 안좋은

Likelihood function $p_\theta(\mathbf{x})$ for Partially Observed Data is hard to compute:

$$p_\theta(\mathbf{x}) = \sum_{\text{All values of } \mathbf{z}} p_\theta(\mathbf{x}, \mathbf{z}) = |\mathcal{Z}| \sum_{\mathbf{z} \in \mathcal{Z}} \frac{1}{|\mathcal{Z}|} p_\theta(\mathbf{x}, \mathbf{z}) = |\mathcal{Z}| \mathbb{E}_{\mathbf{z} \sim \text{Uniform}(\mathcal{Z})} [p_\theta(\mathbf{x}, \mathbf{z})]$$

z가 가질 수 있는 모든 값에 대한 Uniform Distribution

We can think of it as an (intractable) expectation. Monte Carlo to the rescue:

- ① Sample $\mathbf{z}^{(1)}, \dots, \mathbf{z}^{(k)}$ uniformly at random
- ② Approximate expectation with sample average

$$\sum_{\mathbf{z}} p_\theta(\mathbf{x}, \mathbf{z}) \approx |\mathcal{Z}| \frac{1}{k} \sum_{j=1}^k p_\theta(\mathbf{x}, \mathbf{z}^{(j)})$$

Works in theory but not in practice. For most \mathbf{z} , $p_\theta(\mathbf{x}, \mathbf{z})$ is very low (most completions don't make sense). Some completions have large $p_\theta(\mathbf{x}, \mathbf{z})$ but we will never "hit" likely completions by uniform random sampling. Need a clever way to select $\mathbf{z}^{(j)}$ to reduce variance of the estimator.

Intractible한 expextation인 이유 : Z에 대한 Uniform Distribution이지만 Z가 연속형이기 때문

Monte Carlo 방식으로 Uniform Distribution에 대한 Expectation을 Approximation하여 X Marginal Likelihood를 Approximation할 수 있음

Estimator's Variance Issue

Uniform Sampling에서 Monte Carlo를 통한 Estimation의 문제점:대부분의 \mathbf{z} 경우에 $P(\mathbf{x}, \mathbf{z})$ 의 값이 작을 것임 (무수히 많은 \mathbf{z} 경우에 대해 \mathbf{x} 와 크게 연관이 있는 경우는 상대적으로 적으니)

이에 따라 \mathbf{z} 에 대해 uniform sampling하는 방식으로 추정치를 구하면 추정치에 대한 분산이 커짐. (대부분의 \mathbf{z} 에 대해 $P(\mathbf{x}, \mathbf{z})$ 가 대부분 작은 값을 가지다가 갑자기 큰 값을 가지게 되는 경우도 생기니까! 그런데 이런 경우가 1/N의 확률로 나타나니->작은 값을 가질 확률도 1/N이고 큰 값을 가질 확률도 1/N이니 추정치의 분산이 커질 수 밖에 없음) \Rightarrow 실제 log-marginal likelihood에 대해 불안정한 예측(high variance)이므로 unstable한 estimator임.

14. Approximation 2: Importance Sampling

Likelihood function $p_\theta(\mathbf{x})$ for Partially Observed Data is hard to compute:

$$p_\theta(\mathbf{x}) = \sum_{\text{All possible values of } \mathbf{z}} p_\theta(\mathbf{x}, \mathbf{z}) = \sum_{\mathbf{z} \in \mathcal{Z}} \frac{q(\mathbf{z})}{q(\mathbf{z})} p_\theta(\mathbf{x}, \mathbf{z}) = \mathbb{E}_{\mathbf{z} \sim q(\mathbf{z})} \left[\frac{p_\theta(\mathbf{x}, \mathbf{z})}{q(\mathbf{z})} \right]$$

Monte Carlo to the rescue:

- ① Sample $\mathbf{z}^{(1)}, \dots, \mathbf{z}^{(k)}$ from $q(\mathbf{z})$
- ② Approximate expectation with sample average

$$p_\theta(\mathbf{x}) \approx \frac{1}{k} \sum_{j=1}^k \frac{p_\theta(\mathbf{x}, \mathbf{z}^{(j)})}{q(\mathbf{z}^{(j)})}$$

What is a good choice for $q(\mathbf{z})$? Intuitively, frequently sample \mathbf{z} (completions) that are likely given \mathbf{x} under $p_\theta(\mathbf{x}, \mathbf{z})$.

- Uniform Distribution 말고 다른 Distribution에서 Z를 Sampling하는 방식으로 Marginal Likelihood에 대한 근사값을 구하려고함.
- Naive하게 uniform한 분포로부터 sampling하지 말고, $q(\mathbf{z})$ 라는 arbitrary한 분포로부터 importance sampling(각 \mathbf{x} 에 대해 importance한 \mathbf{z} 를 더 often하게 sampling할 수 있도록)을 하자

What is a good choice for $q(\mathbf{z})$? Intuitively, frequently sample \mathbf{z} (completions) that are likely given \mathbf{x} under $p_\theta(\mathbf{x}, \mathbf{z})$.

** X's Marginal Likelihood Estimation with $q(\mathbf{z})$ monte carlo is unbiased

- ③ This is an unbiased estimator of $p_\theta(\mathbf{x})$

$$\mathbb{E}_{\mathbf{z}^{(j)} \sim q(\mathbf{z})} \left[\frac{1}{k} \sum_{j=1}^k \frac{p_\theta(\mathbf{x}, \mathbf{z}^{(j)})}{q(\mathbf{z}^{(j)})} \right] = p_\theta(\mathbf{x})$$

→ Estimator(Monte Carlo Approximation with $q(\mathbf{z})$)의 Expectation 값은 실제 target값인 marginal likelihood of \mathbf{x} 이다 : **Unbiased Estimator**

16. **Log**(X's Marginal Likelihood) Estimation with $q(\mathbf{z})$ monte carlo is **biased**

Estimating log-likelihoods

Likelihood function $p_\theta(\mathbf{x})$ for Partially Observed Data is hard to compute:

$$p_\theta(\mathbf{x}) = \sum_{\text{All possible values of } \mathbf{z}} p_\theta(\mathbf{x}, \mathbf{z}) = \sum_{\mathbf{z} \in \mathcal{Z}} \frac{q(\mathbf{z})}{q(\mathbf{z})} p_\theta(\mathbf{x}, \mathbf{z}) = \mathbb{E}_{\mathbf{z} \sim q(\mathbf{z})} \left[\frac{p_\theta(\mathbf{x}, \mathbf{z})}{q(\mathbf{z})} \right]$$

Monte Carlo to the rescue:

- ① Sample $\mathbf{z}^{(1)}, \dots, \mathbf{z}^{(k)}$ from $q(\mathbf{z})$
- ② Approximate expectation with sample average (*unbiased estimator*):

Log가 붙지 않으면 실제 타겟값에 대해 unbiased estimator

$$p_\theta(\mathbf{x}) \approx \frac{1}{k} \sum_{j=1}^k \frac{p_\theta(\mathbf{x}, \mathbf{z}^{(j)})}{q(\mathbf{z}^{(j)})}$$

Recall that for training, we need the *log*-likelihood $\log(p_\theta(\mathbf{x}))$. We could estimate it as:

$$\log(p_\theta(\mathbf{x})) \approx \log \left(\frac{1}{k} \sum_{j=1}^k \frac{p_\theta(\mathbf{x}, \mathbf{z}^{(j)})}{q(\mathbf{z}^{(j)})} \right) \stackrel{k=1}{\approx} \log \left(\frac{p_\theta(\mathbf{x}, \mathbf{z}^{(1)})}{q(\mathbf{z}^{(1)})} \right)$$

However, it's clear that $\mathbb{E}_{\mathbf{z}^{(1)} \sim q(\mathbf{z})} \left[\log \left(\frac{p_\theta(\mathbf{x}, \mathbf{z}^{(1)})}{q(\mathbf{z}^{(1)})} \right) \right] \neq \log \left(\mathbb{E}_{\mathbf{z}^{(1)} \sim q(\mathbf{z})} \left[\frac{p_\theta(\mathbf{x}, \mathbf{z}^{(1)})}{q(\mathbf{z}^{(1)})} \right] \right)$

- MLE Learning을 위해서는 모델의 (1) \mathbf{x} 에 대한 그냥 likelihood estimation이 아니라 (2) **Log** likelihood estimation이 필요함
- (1)에 대한 Monte Carlo Estimation은 unbiased하지만
- (2)에 대한 Monte Carlo Estimation은 biased! → 로그 값에 대한 추정은 biased

의견공유 하고 싶어요!!

** Monte Carlo Estimation한 값 또한 확률 변수임 샘플링을 매번 할때마다 가지는 샘플들이 다르기 때문 이런 MCE에 대한 기댓값이 실제 값과 일치하다면 unbiased, 일치하지 않다면 biased(**Jensen's inequality**로 인해 항상 실제값보다 작은 값으로 biased된다고 합니다!)

- 따라서 **log**-likelihood에 대한 직접적인 monte carlo estimation은 biased함. 올바르지 않은 estimation
- 다른 방식으로의 안정된 estimation이 필요함

17. ELBO(Evidence Lower Bound) Estimation

Evidence Lower Bound

Log-Likelihood function for Partially Observed Data is hard to compute:

$$\log \left(\sum_{\mathbf{z} \in \mathcal{Z}} p_{\theta}(\mathbf{x}, \mathbf{z}) \right) = \log \left(\sum_{\mathbf{z} \in \mathcal{Z}} \frac{q(\mathbf{z})}{q(\mathbf{z})} p_{\theta}(\mathbf{x}, \mathbf{z}) \right) = \log \left(\mathbb{E}_{\mathbf{z} \sim q(\mathbf{z})} \left[\frac{p_{\theta}(\mathbf{x}, \mathbf{z})}{q(\mathbf{z})} \right] \right)$$

- $\log()$ is a concave function. $\log(px + (1-p)x') \geq p \log(x) + (1-p) \log(x')$.
- Idea: use Jensen Inequality (for concave functions)

$$\log(\mathbb{E}_{\mathbf{z} \sim q(\mathbf{z})} [f(\mathbf{z})]) = \log\left(\sum_{\mathbf{z}} q(\mathbf{z}) f(\mathbf{z})\right) \geq \sum_{\mathbf{z}} q(\mathbf{z}) \log f(\mathbf{z}) = \mathbb{E}_{\mathbf{z} \sim q(\mathbf{z})} [\log f(\mathbf{z})]$$

Choosing $f(\mathbf{z}) = \frac{p_{\theta}(\mathbf{x}, \mathbf{z})}{q(\mathbf{z})}$

$$\log \left(\mathbb{E}_{\mathbf{z} \sim q(\mathbf{z})} \left[\frac{p_{\theta}(\mathbf{x}, \mathbf{z})}{q(\mathbf{z})} \right] \right) \geq \mathbb{E}_{\mathbf{z} \sim q(\mathbf{z})} \left[\log \left(\frac{p_{\theta}(\mathbf{x}, \mathbf{z})}{q(\mathbf{z})} \right) \right]$$

Called Evidence Lower Bound (**ELBO**).

→ Jensen's Inequality로 인한 Log-Likelihood의 하한인 ELBO에 대한 Monte Carlo Estimation(Approximation)은 unbiased함.
ELBO항은 앞에 log가 붙지 않음

- ELBO는 실제 marginal likelihood의 하한임은 보장
- ELBO값은 불편추정을 할 수 있음
- 이 ELBO에 대한 Estimation을 구하는 방식으로 ELBO 값이 실제 Log-Marginal Likelihood와 가까워 지면서도(E step),이 ELBO값이 커지도록 업데이트를 하여 간접적으로 Log-Marginal Likelihood를 크게하자(M step) → EM Algorithm
- \mathbf{z} 를 샘플링할 q 를 어떻게 고르냐에 따라 lower bound의 tight 정도가 결정됨 (하한인 ELBO값이 실제 log-marginal probability랑 얼마나 가까운가)

18. Variational Inference

- Suppose $q(\mathbf{z})$ is **any** probability distribution over the hidden variables
- **Evidence lower bound (ELBO)** holds for any q

$$\begin{aligned}
 \log p(\mathbf{x}; \theta) &\geq \sum_{\mathbf{z}} q(\mathbf{z}) \log \left(\frac{p_{\theta}(\mathbf{x}, \mathbf{z})}{q(\mathbf{z})} \right) \\
 &= \sum_{\mathbf{z}} q(\mathbf{z}) \log p_{\theta}(\mathbf{x}, \mathbf{z}) - \underbrace{\sum_{\mathbf{z}} q(\mathbf{z}) \log q(\mathbf{z})}_{\text{Entropy } H(q) \text{ of } q} \\
 &= \sum_{\mathbf{z}} q(\mathbf{z}) \log p_{\theta}(\mathbf{x}, \mathbf{z}) + H(q)
 \end{aligned}$$

E step : Lower bound를 어떻게 우리가 실제로 관심이 있는 log-marginal likelihood에 tight하게 맞출까?(best approximation?
 => "training"시에 z 샘플링에 이용되는 Q를 model이 관측값 x에 대해 뱉어내는 z의 distribution인 posterior distribution으로 만들면됨)

- Equality holds if $q = p(\mathbf{z}|\mathbf{x}; \theta)$

$$\log p(\mathbf{x}; \theta) = \sum_{\mathbf{z}} q(\mathbf{z}) \log p(\mathbf{z}, \mathbf{x}; \theta) + H(q)$$

- (Aside: This is what we compute in the E-step of the EM algorithm)

E step : Lower bound를 어떻게 우리가 실제로 관심이 있는 log-marginal likelihood에 tight하게 맞출까?(best approximation?)

→ When doing ELBO Approximation, Best way for inferring Z variable(deciding the q distribution for z sampling) is to use true posterior distribution! Q(z)를 실제 모델에서 정의되는 Ptheta(z|x) distribution(True Posterior Distribution)으로 맞추어주면됨.

→ 그렇게 하면 부등식에서의 equality를 취함(ELBO값이 Log marginal likelihood값과 같아짐)

→ 그러나 이 Ptheta(z|x)를 구하는 것은 어렵다고 함.(Intractible)