

Paper summary template

AI VISION Lab

1. 공부한 논문의 제목, 게재된 학회 혹은 저널 등 논문 기본 정보를 적으세요.
 - A. 제목: Estimation of Non-Normalized Statistical Models by Score Matching
 - B. 저널: JMLR(Journal of Machine Learning)
 - C. 저자: Aapo Hyvärinen
 - D. 출판연도: 2005
 - E. Keywords: non-normalized densities, Score Function

2. 논문에서 제안한 알고리즘 및 프레임워크에 대해 본인이 이해한대로 다이어그램을 그려보세요. **논문 Figure를 그대로 따라 그리면 안됩니다.**
 - A. $P(x, \theta) = \frac{1}{z(x)} q(x, \theta)$, 우리는 $z(x)$ 를 알아야, $\text{PDF}(P(x, \theta))$ 를 구할 수 있음
 - B. 그러나 $z(x)$ 는 $q(x, \theta)$ 를 적분한 값으로, $q(x, \theta)$ 의 dimension이 커질수록 구하기 어렵다는 문제가 있음 (n =dimension 크기라고 할 때, $n>2$ 인 경우 문제발생): **문제점 1**
 - C. 즉, 이는 $P(x, \theta)$ 값을 구하는 것을 불가능하게 함.
 - D. 과거 논문에서는 해당 과정을 MCMC, 혹은 근사하는 방식으로 우회
 - E. 이 논문에서는 확률모델을 log-density의 미분 값으로 표현 (Score Function)
 - F. 해당 과정의 경우, gradient of log-density의 distance를 최소화하는 것을 목적으로 함 (즉, model – data 사이의 분포 차이를 최소화함)
 - G. 이 과정에서, log를 취해주고, $\delta\theta$ 에 대해서 미분하는 과정에서, $z(x)$ 값을 고려하지 않고 수식을 전개해 나갈 수 있음
 - H. 해당과정에서 data 1개는 PDF가 아닌 단편적인 정보(=0)이고, p_{data} 에 대한 Score Function을 구하는 과정이 불가능함: **문제점 2**
 - I. p_{data} 의 Score Function 값을 $q(x)$ 의 Score Function에 대한 미분으로 구하는 과정을 증명함으로써, 결과적으로 Object Function을 q 에 대한 Score Function만으로 표현함
 - J. 해당과정을 통해 얻은 Score Object Function을 Multivariate Gaussian Density함수에 적용해서 확인(문제 X) ICA(Independent Component Analysis) Model에 적용하는 과정을 통해 Maximum likelihood estimation과의 차이를 비교(차이 존재)함
 - K. Pseudo-Likelihood Estimation과의 비교를 통해 일관성 측면에서 한계를 가지는 것을 보여줌

- L. 결과적으로는 Contrastive Divergence의 계산이 상대적으로 까다롭지만 Intractable models에도 사용이 가능하다는 장점 때문에 계속해서 사용될 것이라는 주장과, MCMC 기반의 목적함수보다 Consistency측면에서 우수하다는 장점을 가졌음을 언급함.
3. 본인이 생각하는 이 논문의 장점이 무엇이라고 생각하나요? **논문 Contribution bullet을 그대로 따라 적으면 안됩니다.**
- A. Normalized Density Function을 얻기 위해 필요한 Z (Normalized Constant)의 계산이 불가능할 정도로 어렵다는 것은 2005년 이전에도 존재하던 사실로, 이러한 한계를 해결하기 위해 선행 연구에서는 MCMC, 나 Gibbs sampling 같은 기법들을 통해 $P(x)$ 를 추정해 왔음.
- B. 그런데 이 논문에서는 Z 를 구하지 않고도 Non-normalized(논문에서는 $q(x)$)만으로 $P(x)$ 를 추정함으로써 새로운 방법론을 제안하고 있다는 점에서 상당히 우수하다고 생각함.
4. 이 논문을 읽으면서 느낀 점, 혹은 배운 점이 있으면 적어보세요.
- A. 문제해결을 우회하던 기존 연구와 달리 문제자체를 없애는 모습에서 상당히 인상깊었음.
5. 이 논문의 한계점이 있다면 무엇이라고 생각하나요?
- A. 미분을 함으로써 Z 에 대한 점을 고려하지 않아도 된다는 장점도 있지만, Local Consistency (Score 정보, 즉 미분 값)을 기준으로 식을 전개한다는 점에서, $P(x, \theta)$ 의 정확한 정보를 파악하지 못한다는 한계가 있음.
- B. 이 논문에서는 T 가 Inf에 수렴하게 된다면 $P(x, \theta)$ 와 일치함을 보이지만, 이는 실질적으로 이론에 불과함 (무한개의 데이터 = 불가능)
6. 본인의 연구에 접목시켜볼 점이 있을지 생각하고 적어보세요.
- A. Diffusion 기초 논문
7. 본 Summary를 작성하는 과정에서 생성형AI를 사용했나요?
- A. 아니요

날짜: 2025.06.24

이름: 신준원

서명: 

6.24. 화 Gcore matching (2008)

Abstract

확률분포의 추정 $\rightarrow P(X)$ 는 알고있는데 \leftarrow 는 몰라 (즉, normalization constant)
 성가시게는 구하지 못하는데 가려져 있는 정보가 있을 문제 발생.

\rightarrow 가려져 있는 정보.. Markov chain Monte Carlo (MCMC)이라 하는 행렬곱과 파싱
실행 방법 방법.

이 논문에서는 model의 "log-density 정보". 데이터의 log-density 정보의 차이를
최소화 시키는 방법으로 확률 분포의 형태를 추정함.

objective function

이론상 사실 log-density 정보는 매우 어려움 (non-parametric, 비선형)
 \rightarrow Simple formula for this object function. \rightarrow contribution.

denotation

$P_X(\cdot)$ = data. $P(\cdot; \theta)$ = model. $\hat{\theta}$ = estimated parameter value

Introduction

문제 $P(\mathcal{E}; \theta) = \frac{1}{Z(\theta)} \cdot \underbrace{f(\mathcal{E}; \theta)}_{\text{가려져 있는 정보}}$ $Z(\theta) \rightarrow$ 안지.

$$Z(\theta) = \int_{\mathcal{E} \in \mathcal{R}^n} f(\mathcal{E}; \theta) d\mathcal{E} \quad (n \geq 2 \text{ 인 경우, 보충!})$$

\rightarrow PDF 공식 문제.

따라서, or 또 MCMC 사용 \rightarrow 확률이 높은 정보 있음. (Gaussian \rightarrow normalized 가능)

또는, $\mathcal{E} \times \rightarrow$ non-normalized model 사용 방법을 역시 존재함.

\rightarrow 아, normalization 필요하 \Rightarrow 이러 서 만 방법.

$$\min \left(\text{Score function}(\mathcal{X}) - \text{score function}(\text{model}) \right)^2$$

이제 score function = gradient of log-density.

with simple formula

estimation by score function

함수 보정에 대한 loss function은 이이 구성 좋음.

그러면 $p(z; \theta)$ 에 대해 log 후 미분하면.

$$p(z; \theta) = \frac{1}{Z(\theta)} g(z; \theta) \quad \text{--- ①}$$

$$\log(p(z; \theta)) = \log g(z; \theta) - \log(Z(\theta)) \quad \text{--- ②}$$

여기서 g 는 θ 에 대한 항. 즉. 1. data와 무관하므로

미분할 수 있는 항 $\rightarrow Z(\theta)$ 에 의존하지 않아도 됨!

또는 data $z_n \in \mathcal{X}^n$

$$\psi(z; \theta) = \begin{pmatrix} \frac{\partial \log p(z; \theta)}{\partial \theta_1} \\ \vdots \\ \frac{\partial \log p(z; \theta)}{\partial \theta_n} \end{pmatrix} = \begin{pmatrix} \psi_1(z; \theta) \\ \vdots \\ \psi_n(z; \theta) \end{pmatrix} = \nabla_{\theta} \log p(z; \theta)$$

결과적으로 $\psi(z; \theta) = \nabla_{\theta} \log p(z; \theta) = \nabla_{\theta} \log g(z; \theta)$

data

$$\psi_x(\cdot) = \nabla_{\theta} \log p_x(\cdot)$$

\Rightarrow data 분포 기반 계산 가능.

With non-parametric estimation of PDF.

\rightarrow 그러나 해당 과정 필요X (여기서)

loss function

$$J(\theta) = \frac{1}{n} \int_{\mathcal{X}} p_x(z) \|\psi(z; \theta) - \psi_x(z)\|^2 dz$$

$$\rightarrow \hat{\theta} = \arg \min_{\theta} J(\theta) \quad J \text{를 최소화 인 } \theta = \hat{\theta}$$

해당하는 θ 를 찾기 어렵고 부근을 탐색할 수 있다는 장점이 있음.

그러나, 여전히 ψ_x 를 구하기 위해서는

non-parametric estimation problem \Rightarrow 해결해야 함

\rightarrow (제안) \rightarrow Theorem 1*

간단하지!

Theorem 1 가설 1) model score function $\psi(z; \theta)$ 미분가능
 가설 2) some weak regularity conditions
 → 부분 적분법을 사용해서 유도 가능
 $\rightarrow p_X(z) \cdot \nabla_z \psi(z; \theta) \rightarrow$

$$J(\theta) = \frac{1}{n} \int_{z \in \mathcal{R}^n} p_X(z) \|\psi(z; \theta) - \psi_X(z)\|^2 dz \quad \text{--- ①}$$

$$= \frac{1}{n} \int_{z \in \mathcal{R}^n} p_X(z) \|\psi(z; \theta)\|^2 + \boxed{\psi_X(z)^2} - 2 \cdot \psi(z; \theta) \cdot \psi_X(z) dz \quad \text{--- ②}$$

$J(\theta)$, θ 와 관련없으므로 상수항 추가. → 제1항

$$J(\theta) = \frac{1}{n} \int_{z \in \mathcal{R}^n} \boxed{p_X(z) \|\psi(z; \theta)\|^2} - 2 \psi(z; \theta) \cdot \psi_X(z) dz \quad \text{--- ③}$$

제1항 상수항

제2항의 항 → $\psi_X(z)$ 를 구하러가야 불가능

제1항

$$\frac{1}{n} \cdot (-2) \sum_i \int p_X(z) \cdot \boxed{\psi_{X,i}(z)} \cdot \psi_i(z; \theta) dz \quad \text{--- ④}$$

$$\rightarrow -2 \sum_i \int p_X(z) \cdot \frac{\partial \log p_X(z)}{\partial z_i} \cdot \psi_i(z; \theta) dz \Rightarrow -2 \int \boxed{p(z) (\log p(z))'} f(z) dz$$

$\frac{p'(z)}{p(z)}$
이항식

$$\int p(x) \cdot \frac{p'(x)}{p(x)} \cdot f(x) dx = \int p'(x) f(x) dx = \boxed{[p(x)f(x)]} - \int p(x) \cdot f'(x) dx$$

즉, 정리하면 $\sum_{i=1}^n \int p(x) \cdot \psi_i(z; \theta) dz \xrightarrow{\text{부분적분}} \psi_X(z)$ 항이 $J(\theta)$ 표현됨

$$J(\theta) = \int_{z \in \mathcal{R}^n} p_X(z) \sum_{i=1}^n \left[\frac{\partial \psi_i(z; \theta)}{\partial z_i} + \frac{1}{2} \psi_i(z; \theta) \right] dz + \text{const} \quad \text{--- ⑤}$$

$$\psi_i(z; \theta) = \frac{\partial \log g(z; \theta)}{\partial z_i} \quad \left| \quad \frac{\partial \psi_i(z; \theta)}{\partial z_i} = \frac{\partial^2 \log g(z; \theta)}{\partial z_i^2} \right|$$

정리하면 score function 항은 0으로 수렴 (코사키 법칙)

$$\textcircled{1} J(\theta) \approx \text{const}$$

$\theta = \theta_{\text{ML}}$, $\pm = \pm$ 미미한 값

Theorem 2 $J(\theta) = 0 \rightarrow \nabla_{\theta} p_1(\cdot) = \nabla_{\theta} p(\cdot; \theta)$ where $\theta \in \Theta^*$

where Θ^* is the set of θ such that $\nabla_{\theta} p_1(\cdot) = \nabla_{\theta} p(\cdot; \theta)$ (where $p_1, p_{\theta} \in \text{pdf}$, $\int p_1 = 1$)

where $p_1 = p(\cdot; \theta)$ at θ .

Example \rightarrow find the MLE of μ and Σ

① Gaussian Density (multivariate)

$$p(x; \mu, \Sigma) = \frac{1}{Z(\mu, \Sigma)} \exp\left(-\frac{1}{2}(x - \mu)^T \Sigma^{-1}(x - \mu)\right)$$

Σ : symmetric positive-definite matrix $\rightarrow \sim N(x - \mu) \left(-\frac{1}{2}\right)$

$$\eta(x) = \exp\left(-\frac{1}{2}(x - \mu)^T \Sigma^{-1}(x - \mu)\right)$$

$$\nabla \eta(x) = \psi(x; \mu, \Sigma) = -\Sigma^{-1}(x - \mu)$$

$$\partial \psi(x; \mu, \Sigma) = -\Sigma^{-1}$$

$$\tilde{J}(\mu, \Sigma) = \frac{1}{T} \sum_{t=1}^T \left[\sum_{i=1}^n -m_{ii} + \frac{1}{n} (x^{(t)} - \mu)^T \Sigma^{-1} (x^{(t)} - \mu) \right]$$

where \tilde{J} is the information \Rightarrow maximum likelihood estimation of θ .

$\frac{1}{T}$, Consistency of \tilde{J} $\frac{1}{T} \rightarrow 0$ (Theorem 1)

Intuitive Interpretation

ix) non-normalized log-likelihood of maximizing

↓
————→ 가장 좋은 $\mu=0$ / variance = infinite.
————→ which is \rightarrow - min of σ^2 .

Estimation of Basic Independent Component Analysis Model (ICA)

$$\text{가장 좋은 } \log p(x) = \sum_{k=1}^n G(w_k^T x) + f(w_1, \dots, w_n)$$

normalization constant = $-\log |\det W|$ (matrix W)

여기서 \log 는 로그, \det 는 행렬의 행렬식 = sampling 이 된다.