

* Gaussian Distribution (정규 분포, 가우시안 분포, 가우스 분포, 가우시안 확률 분포)

$$\text{Gaussian (Normal) function} \quad f(x) = \frac{1}{\sqrt{2\pi}} e^{-\frac{x^2}{2}}$$

$$X \sim N(\mu, \sigma^2)$$

μ 는 분포값, σ^2 는 분산 또는 분포 폭

$$\text{Probability Density Function (PDF)} \quad f(x) = \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{(x-\mu)^2}{2\sigma^2}} \quad (-\infty < x < \infty)$$

* Markov Chain, MC (= Markov Process, MP) : Markov property를 지니는 이산 시간 (discrete time) 확률 과정 (stochastic process)

- Markov property : 과거 상태들 (s_1, s_2, \dots, s_{t-1})과 현재 상태 (s_t)가 주어졌을 때, 미래 상태 (s_{t+1})는 과거 상태에 독립적으로 현재 상태에 대해서만 결정

$$P[s_{t+1}|s_t] = P[s_{t+1}|s_1, \dots, s_t] \quad (\text{다른 state도 같 확률도 항상 같다})$$

$$* \text{Bayesian Rule} : P(A|B) = \frac{P(B|A)P(A)}{P(B)}$$

$$P(A|B) = \frac{P(A \cap B)}{P(B)} = \frac{P(B|A)P(A)}{P(B)} \quad (= \frac{P(B|A)P(A)}{\sum_i P(B|A_i)P(A_i)})$$

$P(A|B)$: 사후확률 (posterior) 사건 B가 발생한 후 갱신된 사건 A의 확률

$P(A)$: 사전확률 (prior) 사건 B가 발생하기 전에 가지고 있던 사건 A의 확률

$P(B|A)$: 가능도 (likelihood) 사건 A가 발생한 경우 사건 B의 확률

$P(B)$: 정규화 상수 (normalizing constant) 또는 증거 (evidence). 확률의 크기 조정

• 방법 1

전체 확률의 분해

$$P(A|B) = \frac{P(B|A)P(A)}{P(B)} = \frac{P(B|A)P(A)}{\sum_i P(B|A_i)P(A_i)} = \frac{P(B|A_1)P(A_1)}{\sum_i P(B|A_i)P(A_i)}$$

• 방법 2

$$\text{추가적인 사건 C가 발생했다면} \quad P(A|B, C) = \frac{P(C|A, B)P(A|B)}{P(C|B)}$$

$$P(A, B, C) = P(A|B, C)P(B, C) = P(A|B, C)P(C|B)P(B)$$

$$P(A, B, C) = P(C|A, B)P(A, B) = P(C|A, B)P(A|B)P(B)$$

$$\therefore P(A|B, C)P(C|B)P(B) = P(C|A, B)P(A|B)P(B)$$

$$\therefore P(A|B, C) = \frac{P(C|A, B)P(A|B)}{P(C|B)}$$

$$(+ P(A|B \cap C) = \frac{P(A \cap B \cap C)}{P(B \cap C)} = \frac{P(A \cap B \cap C)P(C)}{P(B \cap C)P(C)} = \frac{P(A|C)P(B|C)P(C)}{P(B \cap C)P(C)} = \frac{P(A|C)P(B|C)}{P(B|C)} = P(A|C))$$

* Likelihood

Likelihood: 고정된 관측값이 어떠한 확률 분포에서 어느 정도의 확률로 나타나는 지에 대한 확률 (관측된 사건 고정, 확률 분포 변화)

$$L(\theta|x) \quad (\theta: \text{확률 분포 구성 parameter}, x: \text{관측값, data})$$

Probability: 고정된 확률 분포에서 어떠한 관측값이 나타나는 지에 대한 확률 (확률 분포 고정, 관측된 사건 변화)

$$P(x|\theta)$$

- 단순 사건 $L(\theta|x) = p(x|\theta) = p_{\theta}(x) = p_{\theta}(X=x)$

단순사건에서 확률값 - 확률 질량 함수 PMF (Probability Mass Function)

$$L(\theta|x) \rightarrow p(x|\theta) \quad \text{사건 } x \text{에 대한 확률값 } y, y > 0, \sum f(x) = 1$$

$$\text{여러개 사건 동시에 발생 확률 } P(x|\theta) = \prod_{k=1}^n P(x_k|\theta)$$

$L(\theta|x)$: 관측값이 주어질 때, 변화되는 확률분포에서 주어진 관측값이 나올 확률 / $p(x|\theta)$: 확률분포가 주어져, 변화되는 관측값이 나올 확률

- 연속 사건

연속사건에서 확률 분포 - 확률 밀도 함수 PDF (Probability Density function)

$$L(\theta|x) \rightarrow p(x|\theta) \quad \text{사건 } x \text{에 대한 확률값 } y, y > 0, \int f(x) dx = 1$$

$L(\theta|x)$: 특정 구간 사건이 발생할 확률 \rightarrow PDF의 y값

$$p(x|\theta) \quad \text{사건 발생이 발생할 확률} \rightarrow \text{PDF의 면적} \quad p(a \leq x \leq b) = \int_a^b f(x) dx$$

* MLE 최대 우도 추정법 Maximum Likelihood Estimation: 관측되는 데이터를 가장 잘 설명하는 확률분포 parameter를 찾는 방법

- 확률 분포: 관측값 x 를 얼마나 확률로 나올 확률하는 확률 밀도 함수 PDF 정의
- parameter (θ)를 변수 두었을 때, peak (마지막에 0이 되는 값) 1개
- Likelihood를 계산하기 위한 확률분포의 iid를 사용 (독립 항등 분포, variable에 확률 분포가 independent, identically distribution)

• Exponential Distribution

확률 분포 parameter 1개 / $y = \lambda e^{-\lambda x}$ (λ : 확률 분포 parameter)

$$\begin{aligned} \lambda \text{에 대한 } L \text{의 수학적 표현: } L(\lambda|x_1, x_2, \dots, x_n) &= L(\lambda|x_1)L(\lambda|x_2) \dots L(\lambda|x_n) \quad (\text{iid}) \\ &= \lambda e^{-\lambda x_1} \lambda e^{-\lambda x_2} \dots \lambda e^{-\lambda x_n} = \lambda^n [e^{-\lambda x_1} e^{-\lambda x_2} \dots e^{-\lambda x_n}] \\ &= \lambda^n [e^{-\lambda(x_1+x_2+\dots+x_n)}] \end{aligned}$$

Likelihood의 최댓값 (peak) \rightarrow 미분했을 때 0인 지점

$$\frac{d}{d\lambda} L(\lambda|x_1, x_2, \dots, x_n) = \frac{d}{d\lambda} \lambda^n [e^{-\lambda(x_1+x_2+\dots+x_n)}]$$

$$\begin{aligned} \frac{d}{d\lambda} \log(\lambda^n [e^{-\lambda(x_1+x_2+\dots+x_n)}]) &= \frac{d}{d\lambda} (\log(\lambda^n) + \log[e^{-\lambda(x_1+x_2+\dots+x_n)}]) \\ &= \frac{d}{d\lambda} (n \log(\lambda) - \lambda(x_1+x_2+\dots+x_n)) \\ &= n \frac{1}{\lambda} - (x_1+x_2+\dots+x_n) \end{aligned}$$

$$(x_1+x_2+\dots+x_n) = n \frac{1}{\lambda}$$

$$\lambda(x_1+x_2+\dots+x_n) = n$$

$$\lambda = \frac{n}{(x_1+x_2+\dots+x_n)}$$

$$\mu = \frac{1}{\lambda}, \quad \sigma^2 = \frac{1}{\lambda^2}$$

Gaussian Distribution

이름 붙은 parameter 2개 : 평균 (μ), 분산 (σ) / $P(x|\mu, \sigma) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$

$$\ln [L(\mu, \sigma|x)] = -\frac{n}{2} \ln(2\pi) - n \ln(\sigma) - \frac{(x_1 - \mu)^2}{2\sigma^2} - \dots - \frac{(x_n - \mu)^2}{2\sigma^2}$$

$$\mu, \sigma \text{에 대해 각각 편미분 } \frac{\partial}{\partial \mu} \ln [L(\mu, \sigma|x_1, \dots, x_n)], \mu = \frac{(x_1 + \dots + x_n)}{n}$$

$$\frac{\partial}{\partial \sigma} \ln [L(\mu, \sigma|x_1, \dots, x_n)], \sigma = \sqrt{\frac{(x_1 - \mu)^2 + \dots + (x_n - \mu)^2}{n}}$$

- Log Likelihood : 미분할 때 0이 되는 지점에 해당 X, 미분 수학 계산 용이
- NLL (Negative Log Likelihood) : 최소값을 찾는 것

* 정보이론 기초

Entropy

- 정보량 $H = n \log(s) = \log(s^n)$ (s : 각 선택에서 가능한 결과의 개수, $\log s$: 결과에 대한 정보 개수, n : 결과의 개수)

- Entropy: 불확실성의 측정, H (bits)

처음에 전략을 가면서 그 사건을 예측하는 데 필요한 정보 개수. 결과 예측에 대한 기대값
모든 사건이 같은 확률로 발생할 때 그 불확실성을 갖는다
entropy 감소 = 정보 개수 감소 = 정보량 감소

이전 확률 분포에서, (결과 개수는 각 사건 발생 확률에 역반)

$$H = \sum (\text{사건 발생 확률}) \cdot \log\left(\frac{1}{\text{사건 발생 확률}}\right) = \sum p_i \log\left(\frac{1}{p_i}\right) = -\sum p_i \log(p_i)$$

Cross Entropy

- Cross Entropy: 어떤 문제에 대해 특정 전략을 쓸 때 예상되는 정보 개수에 대한 기대값 (전통: 학습 분포)
학습 분포로 된 어떤 문제 p에 대해 학습 분포로 된 어떤 전략 q를 사용할 때 정보 개수의 기대값

이산형

$$H(p, q) = \sum p_i \log\left(\frac{1}{q_i}\right) = -\sum p_i \log(q_i)$$

연속형

$$-\int p(x) \log q(x) dx$$

(여기서 p_i : 특정 확률에 대한 한도 또는 분포 확률, q_i : 현재 학습한 확률)

- Binary classification (1과 0 두 가지 결과만 존재)

$$-y \log \hat{y} - (1-y) \log(1-\hat{y}) \quad p = [y, 1-y], q = [\hat{y}, 1-\hat{y}] \text{인 cross entropy의}$$

→ logistic regression에서 쓰는 cost function.

- Cross Entropy → Log loss → negative log likelihood : Cross Entropy 한 것과 같은 것은 log likelihood 이다

$$\text{이런 데이터가 0 또는 1로 예측된 확률은 } \hat{y}, 1-\hat{y}, \text{ likelihood 식 } \rightarrow \hat{y}^y (1-\hat{y})^{(1-y)}$$

$y=1$ 일 때 \hat{y} 이다, $y=0$ 일 때 $(1-\hat{y})$ 이다. 이를 쓰면

$$\text{maximize } y \log \hat{y} + (1-y) \log(1-\hat{y}) = \text{minimize } -y \log \hat{y} - (1-y) \log(1-\hat{y})$$

- KL-divergence (Kullback-Leibler divergence) : 두 확률분포의 엔트로피 차이 계산

$$\begin{aligned}
 H(p, q) &= -\sum_i p_i \log q_i \\
 &= -\sum_i p_i \log p_i - \sum_i p_i \log \frac{p_i}{q_i} + \sum_i p_i \log p_i \\
 &= H(p) + \sum_i p_i \log \frac{p_i}{q_i} - \sum_i p_i \log p_i \\
 &= H(p) + \sum_i p_i \log \frac{p_i}{q_i}
 \end{aligned}$$

→ per unit q 에 정보량 차이 → KL-divergence

$$\sum_i p_i \log \frac{p_i}{q_i} = H(p, q) - H(p) \quad \text{per } q \text{에 cross entropy에서 } p \text{의 entropy가 빠져있기 때문. 두 분포의 차이}$$

$$KL(p||q) = H(p, q) - H(p)$$

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

- $KL(p||q) \geq 0$
- $KL(p||q) \neq KL(q||p)$

* generative model vs discriminative model

- 생성모델 : posterior를 간접적으로 도출 / 레이블 변수의 분포를 학습
레이블 X 가 생성되는 과정을 두 개의 확률 모형 $P(Y)$, $P(X|Y)$ 로 정의하여 $P(Y|X)$ 를 간접적으로 도출
- 판별모델 : posterior를 직접적으로 도출 / 결정경계 (decision boundary)를 학습
레이블 X 가 주어졌을 때 레이블 Y 가 나타날 조건부 확률 $P(Y|X)$ 를 직접적으로 반환 EX. 선형 회귀, 코지스택 회귀

