

2018/10/06/7/4 기계학습과 강의

[ML-15] Probability & Information Theory

Probability

- experiment : 모든 가능한 결과

outcome : experiment의 결과.

deterministic : 동일한 조건에서 반복하여 항상 동일한 결과.

random : 동일한 조건에서 반복해서 항상 동일한 결과 X

sample space : possible outcome들의 집합

event space : possible event들의 집합

event : sample space의 부분집합.

Axioms of probability

1. $P(A) \geq 0$ 2. $P(S) = 1$ 3. $P(A \cup B) = P(A) + P(B)$ ($P(MB) = \emptyset$)

Properties of probability

1. $P(A^c) = 1 - P(A)$

4. $P(A) \leq 1$

2. $P(\emptyset) = 0$

5. $P(A \cup B) = P(A) + P(B) - P(AB)$

3. $P(A) \leq P(B)$ (if $A \subseteq B$)

Conditional Probability : event A given event B

$$P(A|B) = \frac{P(AB)}{P(B)}$$

B가 주어지면 A의 확률

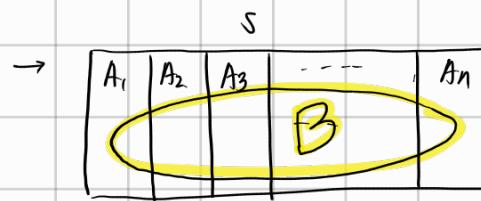
Bayes' Rule

$$\rightarrow P(A|B) = P(A|B)P(B) = P(B|A)P(A)$$

$$\rightarrow P(A|B) = \frac{P(B|A)P(A)}{P(B)}$$

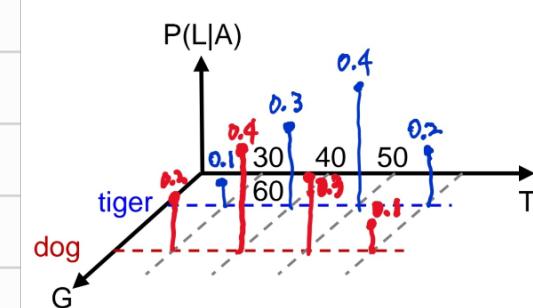
Bayes' Theorem

Suppose the events, A_1, A_2, \dots, A_n : partition of S



$$\rightarrow P(B) = \sum_{i=1}^n P(B|A_i)P(A_i) = \sum_{i=1}^n P(B|A_i)P(A_i)$$

Example)



$L=50$ 일 때 dog인지 tiger인지

ML test

$$\rightarrow P(L=50|\text{dog}) = 0.1 \quad P(L=50|\text{Tiger}) = 0.4$$

→ dog일 때 $L=50$ 의 확률은 0.1
→ tiger일 때 $L=50$ 의 확률은 0.4

MAP test

$$\rightarrow P(L=50|\text{dog})P(\text{dog}) = 0.24 \quad P(L=50|\text{Tiger})P(\text{Tiger}) = 0.08$$

→ 즉, tiger와 dog 각각에서 $L=50$ 의 확률은 Tiger가 높지만

Tiger와 dog의 개체수에 대한 확률이 적용된다

$L=50$ 일 때 전체에서 tiger와 dog의 확률은 dog가 더 높음.

ML classification.

$$\rightarrow k^* = \arg \max P(\mathbf{z}_{\text{new}} | c_k)$$

MAP classification.

$$\rightarrow k^* = \arg \max P(\mathbf{z}_{\text{new}} | c_k)P(c_k)$$

Independence

$$\rightarrow P(A \cap B) = P(A) P(B)$$

$$P(\bigcap_{i=1}^n A_i) = \prod_{i=1}^n P(A_i)$$

Independent rv

$$\rightarrow P(X=x, Y=y) = P(X=x) P(Y=y)$$

Discrete : $P_{XY}(x_i, y_j) = P_X(x_i) P_Y(y_j)$

Continuous : $f_{XY}(x, y) = f_X(x) f_Y(y)$

Naive Bayes classifiers

$$P(\mathbf{z}|C_k) = \frac{b}{\prod_{i=1}^d P(z_i|C_k)}$$

↳ 각 판정변수 (z)들의 발생 확률은 통상 독립으로 가정한다.

But 판정변수를 구성하는 요소들의 특징은 보장할 수 없다.

→ Naive Bayes는 1步도 독립이라고 가정하는 것.

Expectations

$$E[X] = \int \sum_k x_k p_X(x_k)$$

$$\begin{aligned} \text{Var}[X]^2 &= \text{Var}(X) = E[(X - E(X))^2] \\ &= \int (x - \mu_X)^2 f_X(x) \end{aligned}$$

$$\text{Var}(X) = E(X^2) - (E(X))^2$$

Conditional Expectation

$$E[Y|X] = \int \sum_k y_k p(Y_k|X)$$

$$E[g(X)] = \int g(x_k) p(x_k)$$

Correlation & Covariance → Independent $\xrightarrow{x} \xleftarrow{\text{O}} \text{uncorrelated}$

- Correlation : $E(X)$

[Orthogonal] : $E(XY) = 0$

[uncorrelated] : $E(XY) = E(X)E(Y)$

- Covariance : $\text{Cov}(X, Y) = \text{Cov}_{XY} = E[(X - E(X))(Y - E(Y))]$

↳ 평균 차이

$$= E(XY) - E(X)E(Y)$$

uncorrelated $\text{Cov}_{XY} = 0$

- Correlation coefficient : a normalized covariance

$$\rightarrow \rho_{XY} = \frac{\text{Cov}_{XY}}{\sqrt{\text{Var}[X]\text{Var}[Y]}} \quad |\rho_{XY}| \leq 1$$

- Correlation coefficient & Linear Dependence

let $Y = aX + b$

$$\Rightarrow ① E(Y) = E(aX + b) = aE(X) + b$$

$$E(Y) = E(aX + b) = aE(X) + b, \quad \text{Var}[Y] = a^2 \text{Var}[X]$$

$$\text{Cov}(X, Y) = \text{Cov}_{XY} = E(XY) - E(X)E(Y)$$

$$= aE(X^2) + bE(X) - a(E(X))^2 - bE(X)$$

$$= a \{ E(X^2) - (E(X))^2 \} = a \text{Var}[X]$$

$$② \rho_{XY} = \frac{\text{Cov}_{XY}}{\sqrt{\text{Var}[X]\text{Var}[Y]}} = \frac{a \text{Var}[X]}{\sqrt{\text{Var}[X](a^2 \text{Var}[X])}} = \begin{cases} 1 & a > 0 \\ -1 & a < 0 \end{cases}$$

Gaussian Distribution

$$f_X(\underline{x}) = \frac{1}{(2\pi)^{n/2} |\det K|^{1/2}} \exp \left[-\frac{1}{2} (\underline{x} - \underline{\mu})^T K^{-1} (\underline{x} - \underline{\mu}) \right]$$

$$\underline{x} = \begin{bmatrix} x_1 \\ \vdots \\ x_n \end{bmatrix}, \quad \underline{\mu} = \begin{bmatrix} \mu_1 \\ \vdots \\ \mu_n \end{bmatrix} = \begin{bmatrix} E(X_1) \\ \vdots \\ E(X_n) \end{bmatrix}, \quad K = \begin{bmatrix} \sigma_{11} & \cdots & \sigma_{1n} \\ \vdots & \ddots & \vdots \\ \sigma_{n1} & \cdots & \sigma_{nn} \end{bmatrix}$$

$$\sigma_{ij} = \text{cov}(X_i, X_j)$$

Information Theory

Information

→ 정보이론의 기본원리 : 학습의 핵심은 (불확실성이 감소) 정보 얻기

자기정보 : 사건 A의 정보량 → 불확실성을 계량화.

$$I(x_i) = \log_b \frac{1}{P(x_i)} = -\log_b P(x_i)$$

property : ① $I(x_i) = 0$ for $P(x_i) = 1$

② $I(x_i) \geq 0$

③ $I(x_i) > I(x_j)$ If $P(x_i) < P(x_j)$

④ $I(x_i, x_j) = I(x_i) + I(x_j)$ If x_i and x_j are Independent

Entropy (불확실성)

$$\begin{aligned} H(X) &= E[I(x_i)] = \sum_{i=1}^m P(x_i) I(x_i) \\ &= -\sum_{i=1}^m P(x_i) \log_2 P(x_i) \end{aligned}$$

Property : $0 \leq H(X) \leq \log_2 m$

$\Rightarrow H(X)$ 는 각 symbol이 일어날 확률에 동일할 때 최대

각 symbol이 일어날 확률이 동일하다는 것은 어느 한 조의 확률이 끝 때보다

불확실성이 크다는 것

1) 동일한 : 70%, 틀리 : 30%

2) 동일한 : 50% 틀리 : 50%

동일한 확률에 일어나면서 유가 나온다 더 예측 불가능한건 그거

\Rightarrow 즉 확률이 동일할 때 불확실성 $\uparrow \Rightarrow I(X) \uparrow \Rightarrow H(X) \uparrow$

Source Coding Theorem

- Average code length

$$L = \sum_{i=1}^m P(x_i) \overbrace{n_i}^{\rightarrow 각 symbol의 개수}$$

$$L \geq H(X) = -\sum_{i=1}^m P(x_i) \log_2 P(x_i)$$

Cross-Entropy, KL Divergence

$p(x)$: 실제값 분포, $q(x)$: 추정 학습 분포

polynomial q 의 교차엔트로피

$$\begin{aligned} H(p, q) &= E_p [I_q(x)] = E_p [-\log(q(x))] \\ &= -\sum_{i=1}^m p(x_i) \log(q(x_i)) \end{aligned}$$

분산 p의 KL Divergence (상대엔트로피)

$$D_{KL}(p||q) = \sum_{i=1}^m p(x_i) \log \left(\frac{p(x_i)}{q(x_i)} \right) = H(p, q) - H(p, p) \geq 0$$

$$\Rightarrow H(p, q) = H(p, p) + D_{KL}(p||q)$$

Cross Entropy loss

$$\begin{aligned} &y \in \{0, 1\} : \text{true} \quad \left. \right) \rightarrow H(y_i, h(x_i)) \\ &h(x) : \text{predicted} \quad \left. \right) = -\sum_{i=1}^k y_i \log(h(x_i)) \end{aligned}$$