

무슨 3주차 강의 Naive Bayes Classifier 강의

Unit 02: 베이즈 정리

베이즈 정리란: 두 확률 변수 (prior) 과 사후확률 (posterior) 사이의 관계를 나타내는 정리

↳ 사전확률 (prior)로부터 사후확률 (posterior)을 얻고자한다.

(사전에 A에 대한 것들 비로써 P(A)와 P(B|A)를 알고있다면, 이를 통해서 P(A|B)를 추론할 수 있게 해준다.)

↳ 조건부확률!

P(B)를 생각해 보면, B라는 결과는 A의 독립항으로 볼 수 있고, $P(B) = P(A \cap B)$ 이고,
A의 독립항을 A_i 라 하면 $P(B) = P(A \cap B) = P(A_1 \cap B) + P(A_2 \cap B) + P(A_3 \cap B)$ ($i=1,2,3$)
∴ 앞 이항에 장황함 없다: $P(B) = P(B|A_1)P(A_1) + P(B|A_2)P(A_2) + P(B|A_3)P(A_3)$

$$\therefore P(A_i|B) = \frac{P(A_i \cap B)}{P(B)} = \frac{P(B|A_i)P(A_i)}{P(B|A_1)P(A_1) + P(B|A_2)P(A_2) + P(B|A_3)P(A_3)}$$

사전확률 '과거 값을 잘 설명한 것'
조건부 확률 (H)을 비로써 하는 확률과
사전(D)의 확률.

$$P(H|D) = \frac{P(D|H)P(H)}{P(D)}$$

Posterior

Normalizing Constant : 사건 D의 발생 확률.

Prior : 사전확률.

- 관측결과 D가 발생한 조건하의 파라미터의 확률

- 사후확률
- 과거 값을 보면 나중에 환경 파라미터의 확률.
- 확률적 (데이터, 실험 등)의 관찰 및 그에 불합치하면, 개인을 통해 '사후적'인 확률 값을.

- 조건부 확률 (H)과 무관한 관측결과 (D) 자체의 확률.
- 보통 상수 상자로 무시하고 계산.

⇒ but 안먹으려하면... 변수가 늘어날수록 연산량이 기하급수적으로 많아짐. $(2^d - 1)^k$ 개 (혹은 많긴해...)

⇒ ∴ solution : 조건부 독립을 가정. 그 때만...

Unit 3: Naive Bayes Classification

Naive Bayes Classification

↳ 계산량이 기하급수적으로 많아지는 것을 막기 위해 종속 변수 (Y)가 주어졌을 때, 입력 변수들이 모두 독립이라고 가정한다.

예측한 결과의 조건부 확률은 각 조건부 확률의 곱으로
추정할 수 있을 수 있는 단순한 가정을 가요.

알아야 할 파라미터 수가 대폭 줄어들게 된다. $(2^d - 1)^k$ 개 → dk 개
feature들의 곱으로 바뀌면서 계산이 수월해진다.

$$P^*(X) = \arg \max_{Y=q} P(X=q|Y=q) P(Y=q)$$

$$\approx \arg \max_{Y=q} P(Y=q) \prod_{i=1 \leq i \leq d} P(X=x_i|Y=q)$$

라플라스 스무딩

↳ 사후확률이 0일 때, likelihood가 0이 되면서 결과값이 0이 되는 경우 발생

이런 경우를 방지하기 위해 0.5번의 확률을 강제한다.

$$P_{LAP} = \frac{c(X)+1}{\sum c(X)+V}, \quad P(X|C) = \frac{\text{count}(X,C)+1}{\sum_{x \in V} \text{count}(X,C)+V}$$

나이브 베이즈 장/단점.

장점.	단점
- 예측중의 차원이 높을 때 유리	- 확률적 확률이 높을 때 (라플라스 스무딩)
- Test에서 강점.	- 조건부 독립이라는 가정이 비현실적.
- 가벼운 나이브 베이즈를 활용하면 input이 연속형: ok	