

# Naive Bayes Classifier

## Unit 1 확률기초.

• 조건부 확률  $P(B|A) = \frac{P(A \cap B)}{P(A)}$

독립 :  $P(A \cap B) = P(A)P(B)$

조건부 독립 :  $P(A, B|C) = P(A|C)P(B|C)$

→ 한 사건이 일어났다는 가정 하에서, 서로 다른 두 사건은 독립인 상황.

## Unit 2 베이즈정리.

$$P(H|D) = \frac{\overset{\text{likelihood}}{\downarrow} P(D|H) \overset{\text{Prior}}{\downarrow} P(H)}{\underset{\text{Posterior}}{\uparrow} P(D) \underset{\text{Normalizing constant}}{\uparrow}}$$

→ 두 확률변수의 prior와 posterior 사이의 관계를 나타내는 정리  
: Prior로부터 posterior를 구하고자 한다.

- Prior : 사건 확률, 과거 경험을 토대로 사후적으로 추정된 parameter  $H$ 의 확률
- Posterior : 사후 확률, 관측결과 사건  $D$ 가 일어난 조건 하의  $H$ 의 확률
- Likelihood : 모델 파라미터  $H$ 를 바탕으로 하는 관측결과 사건  $D$ 의 확률
- Normalizing Constant : 사건  $D$ 의 발생 가능성. 보트 상수로 취급.

## Unit 3. Naive Bayes Classification

• 가정 : 종속변수  $Y$ 가 주어졌을 때, 입력변수들이 모두 독립이다 (조건부 독립 가정)

$$f^*(x) = \underset{Y=y}{\operatorname{argmax}} P(X=x|Y=y) P(Y=y)$$

$$\approx \underset{Y=y}{\operatorname{argmax}} P(Y=y) \prod_{i=1, i \in d} P(X=x_i|Y=y_i)$$

↳ 장점 ① 알려야 할 라카미터의 수가 줄어든다.

③ text에서 강점.

② feature들의 종류로 바뀌면서 계산이 수월해진다.

④ input의 연속형일 때로 사용가능 (가우시안 나이브베이즈)

↳ 단점 ① 희귀한 확률이 나왔을 때

⑤ 조건부 독립 가정 자체가 비현실적.

• 라플라스 수정  $P_{LAP} = \frac{C(x) + 1}{\sum x [C(x) + 1]}$

→ likelihood  $P(D|H)$ 가 0이 되는 것을 방지하도록 최소한의 확률을 정해놓음