

• 분류 (Classification) 는 독립변수 값이 주어졌을 때, 그 독립변수 값과 가장 연관성이 큰 종속변수 (카테고리)를 찾는다는 일이다. 선택해야할 카테고리 혹은 클래스가 미리 주어졌다는 점에서 분류가 주어진 객체식 일치를 푸는 것과 같다.

## - 분류 모형의 종류

• 분류모형을 푸는 방법은 크게 두가지로 나눌 수 있다.

### 판별함수 모형 (discriminant function)

→ 카테고리가 둘로 서로 다른 영역을 나누는 boundary를 찾는 다음, 주어진 데이터가 어디에 속하는지 찾아내는 판별함수를 이용하는 방법

• Perceptron / Support Vector Machine / Neural network

### 회귀분석 모형

→ 2개 이상의 독립변수 값을 가지고 1개의 카테고리를 계산한다. 2개의 독립변수 값을 찾는 방법이기도 크게 두가지로 나뉜다. 2개의 독립변수

• Linear, Quadratic Discriminant Analysis / Naive Bayes (회귀분석 "1회" 모형)

• Logistic regression / Decision tree (회귀분석 "2회" 모형)

### 회귀분석 모형의 사용

훈련데이터가  $n$ 개,  $k$ 개의 클래스  $1, \dots, k$  중의 하나의 값을 가진다고 가정한다.

회귀분석모형은 다음과 같은 함수로  $x$ 에 대한 클래스를 예측한다.

(1) 입력 값이 주어졌을 때,  $n$ 개 클래스  $k$ 가 될 확률  $P(y=k | x)$ 을 모두 계산하고,

$$\begin{aligned} P_1 &= P(y=1 | x) \\ &\vdots \\ P_k &= P(y=k | x) \end{aligned}$$

(2) 다음으로 가장 확률이 큰 클래스를 선택하는 방법이다.

$$y = \text{arg max}_k P(y=k | x)$$

Scikit-learn에서 2개의 독립변수 모형은 `Predict_proba()`와 `Predict_log_proba()`를 제공한다.

### 회귀분석 Naive 모형의 사용

- Naive 모형은 서로 각 클래스별의 데이터 확률인  $P(x | y=k)$ 로 주어진 다음 베이즈 정리를 사용하여  $P(y=k | x)$ 를 계산하는 방법이다.

$$P(y=k | x) = \frac{P(x | y=k) P(y=k)}{P(x)}$$

- Naive모형에서는 전체 확률법칙을 이용하여 특징데이터  $x$ 의 수직적 확률들  $P(x)$ 를 구할 수 있다. (또는 클래스별 확률을 구할 수 있다)

$$P(x) = \sum_{k=1}^K P(x | y=k) P(y=k)$$

↳ 계산량이 많은 단점 /

### Naive 베이즈 (Naive Bayesian)

- 2개의 독립변수 기반 Naive 모형의 성능은 클래스가 3개인 경우에도 비로써 사용할 수 있다는 점이다. Naive 베이즈 모형은 2개의 독립변수 모형과 유사하다.

### 회귀분석 판별모형

- QAD나 Naive 베이즈 모형은 베이즈 모형의 성능을 하는 2개의 독립변수 기반 판별 모형이다. 2개의 독립변수 기반 판별 모형은 2개의 독립변수  $P(y | x)$ 를 구하기 위해 Likelihood  $P(x | y)$ 를 구하고 베이즈 정리를 사용하여 2개의 독립변수를 계산한다.

- 여기서 확률과 판별모형은 2개의 독립변수  $P(y=1 | x)$ 이  $x$ 에 대한 함수  $f(x) > 0$ 로 판별될 수 있다고 가정하고 그 함수를 직접 찾아내는 방법이다.

$$P(y=k | x) = f(x)$$