

## 0.3 카테고리 분포와 다항 분포

비크로이 분포  $\rightarrow$  이진 분류 문제

카테고리 분포  $\rightarrow$  다중 분류 문제 (multi-class classification)

### 카테고리 확률 분포

동전이 여러 개를 던지면?

카테고리 확률 분포  $\rightarrow$  1부터  $k$ 까지  $k$ 개 (중복) 중 1

$\hookrightarrow$  벙글벙글, 카테고리, 클래스

$\hookrightarrow$  주사위 던지기  $\Rightarrow k=6$ 인 카테고리 분포

"주사위 던지기" (one-hot encoding)

$$x=1 \rightarrow x = (1, 0, 0, 0, 0, 0)$$

$$x=2 \rightarrow x = (0, 1, 0, 0, 0, 0)$$

$\vdots$

$$x=6 \rightarrow x = (0, 0, 0, 0, 0, 1)$$

확률 분포가 표시

$$x = (x_1, x_2, x_3, x_4, x_5, x_6)$$

$$x_i = \begin{cases} 0 \\ 1 \end{cases}$$

constraint

$$\sum_{k=1}^K x_k = 1$$

1원 한 번 등장한다는 의미.

카테고리 분포의 평균

$$\mu = (\mu_1, \mu_2, \dots, \mu_K)$$

$0 \leq \mu_i \leq 1$  ... 0부터 1 사이 어떤 실수 값도 가능.

$$\sum_{k=1}^K \mu_k = 1 \quad \dots \text{예) } 0.2 + 0.2 + 0.2 + 0.2 + 0.2 = 1$$

# 카테리리 확률분포

카테리리 확률분포의 확률분포인 카테리리 확률분포

$$\text{Cat}(x_1, x_2, \dots, x_k; \mu_1, \dots, \mu_k)$$

3. 표기거나 확률벡터  $x = (x_1, x_2, \dots, x_k)$ ,  
 확률벡터  $\mu = (\mu_1, \mu_2, \dots, \mu_k)$   $\sum \mu_i = 1$  조건

$$\text{Cat}(x; \mu)$$

2. 간단히 표기할 수 있다.

확률질량함수를 다음처럼 표기한다.

one-hot-encoding  
 행으로  
 간략하게

$$\text{Cat}(x; \mu) = \begin{cases} \mu_1 & \text{if } x = (1, 0, 0, \dots, 0) \\ \mu_2 & \text{if } x = (0, 1, 0, \dots, 0) \\ \mu_3 & \text{if } x = (0, 0, 1, \dots, 0) \\ \vdots & \vdots \\ \mu_k & \text{if } x = (0, 0, 0, \dots, 1) \end{cases}$$

$$\text{Cat}(x; \mu) = \mu_1^{x_1} \mu_2^{x_2} \dots \mu_k^{x_k} = \prod_{k=1}^k \mu_k^{x_k}$$

예제 8.3.1

$k=2$  인 카테리리 분포의 pmf (확률질량함수)가 베르누이 분포의 pmf와 같음을 보이자.

$$k=2, \text{Cat}(x; \mu) = \begin{matrix} x_1 & x_2 \\ \mu_1 & \mu_2 \end{matrix} \quad \text{제한조건 } \mu_1 + \mu_2 = 1 \\ \therefore \mu_2 = 1 - \mu_1$$

$$\text{Cat}(x; \mu) = \mu_1^{x_1} (1 - \mu_1)^{x_2}$$

$$\checkmark \quad x_1 + x_2 = 1 \text{ or } x_2 = 1 - x_1 \\ \text{Cat}(x; \mu) = \mu_1^{x_1} (1 - \mu_1)^{1-x_1} = \text{Bern}(x; \mu)$$

# 카테고리 분포의 모멘트

• 기대값

$$E[x_k] = \mu_k$$

• 분산

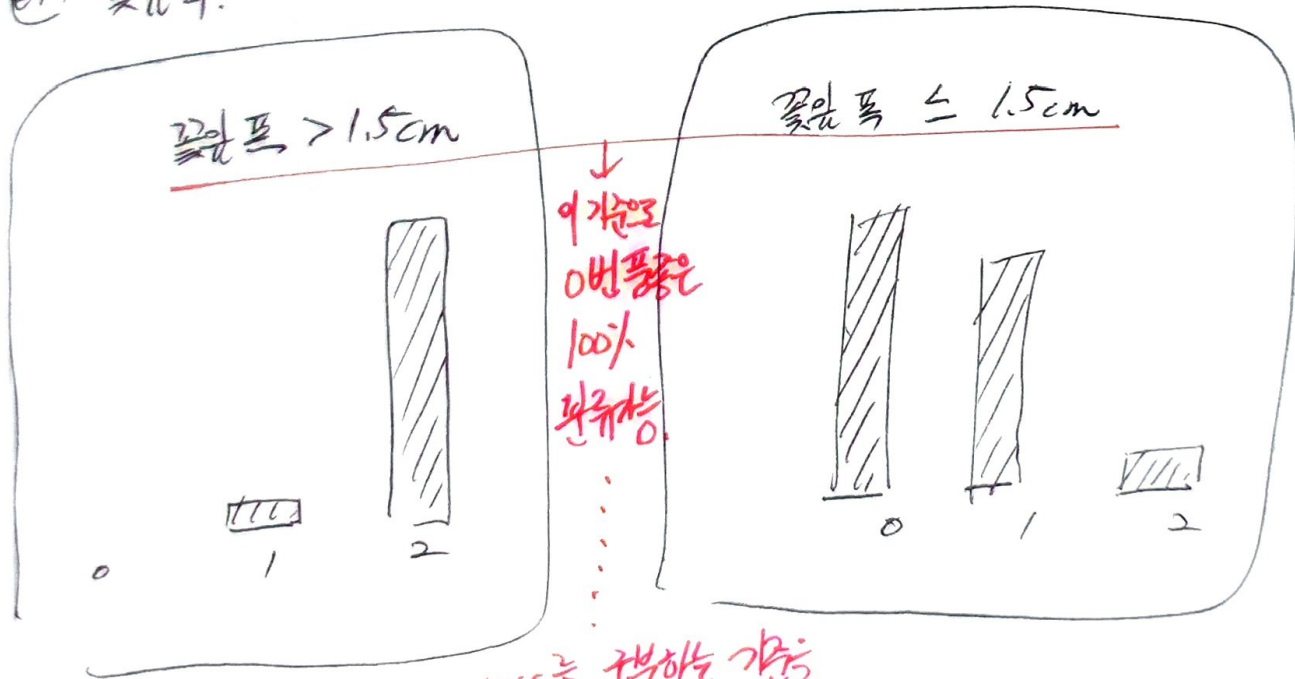
$$\text{Var}[x_k] = \mu_k(1 - \mu_k)$$

## 다중 분류 문제

"multi-class classification"

예측 범주가 2개 이상인 다중 분류 문제 — "카테고리 분포를 사용함"

④ 꽃잎 폭.



class를 구별하는 기준을  
탐색하고 선택

decision tree

random forest



# 다항분포

Multinomial distribution

비례수 확률변수 데이터가 복수이면, 데이터의 합은 이항분포를 이룬다  
카테고리 확률변수 데이터가 복수이면, 데이터의 합은 다항분포를 이룬다.

독립터지기  $X^N$  양면분포 확률변수  $\sim$  이항분포

주변변량기  $X^N$  각변이 다른 확률변수의 분포  $\sim$  다항분포

[다항분포의 확률질량함수]

$$Mu(x; N, \mu)$$

$$Mu(x; N, \mu) = \binom{N}{x} \prod_{k=1}^K \mu_k^{x_k} = \binom{N}{x_1, \dots, x_K} \prod_{k=1}^K \mu_k^{x_k}$$

이 식에서 조합 계수는 다음과 같이 정의된다.

$$\binom{N}{x_1, \dots, x_K} = \frac{N!}{x_1! \dots x_K!}$$

## 예제 8.8.3

$K=2$ 인 다항분포의 확률질량함수가 이항분포의 확률질량함수와 같음을 보이자.

$K=2$ 일때 다항분포의 확률질량함수

$$Mu(x; N, \mu) = \binom{N}{x_1, \dots, x_K} \prod_{k=1}^K \mu_k^{x_k} = \binom{N}{x_1, x_2} \prod_{k=1}^2 \mu_k^{x_k}$$

분류데이터의 총합이  $N$ 개이므로  $x_1 + x_2 = N \dots \textcircled{1}$

분류의 총합이 1이므로  $\mu_1 + \mu_2 = 1 \dots \textcircled{2}$

①, ②를 전제에 대입하면

$$Mu(x; N, \mu) = \binom{N}{x_1, x_2} \prod_{k=1}^2 \mu_k^{x_k} (1 - \mu_1)^{N - x_1} = Bin(x; N, \mu)$$

# 다항분포의 모멘트

다항분포의 기댓값과 분산

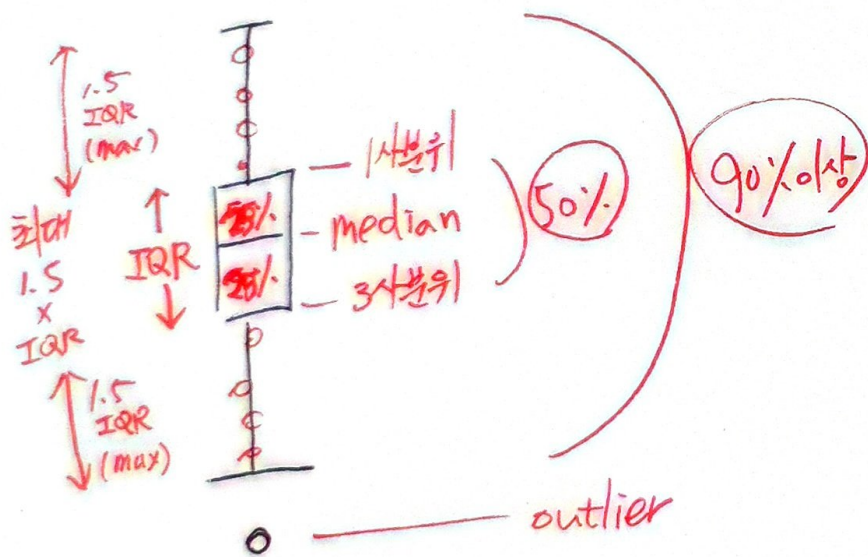
○ 기댓값

$$E[X_k] = N\mu_k$$

○ 분산

$$Var[X_k] = N\mu_k(1-\mu_k)$$

\* Boxplot



**\*\* ✗ boxplot은 분포의 형태가 대체로 정규분포를 갖는 가정을**  
**(대중 중앙이 많이 몰려있는 것이다)**