

< 7장 의문점, 문제 >

1. Empirical Distribution Function

Def. of CDF $\rightarrow F_X(x) = \mathbb{P}(X \leq x)$.

7.3 Theorem. At any fixed value of x ,

$$\begin{aligned}\mathbb{E}(\hat{F}_n(x)) &= F(x), \\ \mathbb{V}(\hat{F}_n(x)) &= \frac{F(x)(1-F(x))}{n}, \\ \text{MSE} &= \frac{F(x)(1-F(x))}{n} \rightarrow 0, \\ \hat{F}_n(x) &\xrightarrow{P} F(x).\end{aligned}$$

7.1 Definition. The **empirical distribution function** \hat{F}_n is the CDF that puts mass $1/n$ at each data point X_i . Formally,

$$\hat{F}_n(x) = \frac{\sum_{i=1}^n I(X_i \leq x)}{n} \quad (7.1)$$

where

$$I(X_i \leq x) = \begin{cases} 1 & \text{if } X_i \leq x \\ 0 & \text{if } X_i > x. \end{cases}$$

$\circ E(\hat{F}_n(x)) = F(x)$ ↖ fixed

$n \cdot \hat{F}_n(x) = \sum_{i=1}^n I(X_i \leq x) \sim \text{Bin}(n, F(x))$

$$\begin{aligned} &= E\left(\frac{\sum_{i=1}^n I(X_i \leq x)}{n}\right) = \frac{1}{n} E\left(\sum_{i=1}^n I(X_i \leq x)\right) \\ &= \frac{1}{n} \sum_{i=1}^n E[I(X_i \leq x)] \\ &= \frac{1}{n} \sum_{i=1}^n P(X_i \leq x) \\ &= \frac{1}{n} \sum_{i=1}^n F_{X_i}(x) \\ &= F(x) \quad \downarrow ? \end{aligned}$$

$$\textcircled{2} V(\hat{F}_n(x)) = \frac{F(x)(1-F(x))}{n}$$

$$\begin{aligned} &= E[\hat{F}_n(x)^2] - E[\hat{F}_n(x)]^2 \\ &= \frac{1}{n^2} \sum_{i=1}^n \sum_{j=1}^n P(X_i \leq x, X_j \leq x) - \frac{1}{n^2} \sum \sum P(X_i \leq x) P(X_j \leq x) \\ &= \frac{1}{n^2} \sum \sum [P(X_i \leq x, X_j \leq x) - P(X_i \leq x) P(X_j \leq x)] \\ &\quad \vdots \\ &= \frac{F(x)(1-F(x))}{n} \end{aligned}$$

③

$$\text{MSE} = \frac{F(x)(1-F(x))}{n} \rightarrow 0,$$

$$\hat{F}_n(x) \xrightarrow{P} F(x).$$

X_n converges to X in probability, written $X_n \xrightarrow{P} X$, if, for every $\epsilon > 0$,

$$\mathbb{P}(|X_n - X| > \epsilon) \rightarrow 0 \quad (5.1)$$

as $n \rightarrow \infty$.

$$\mathbb{P}(|X - \mu| \geq t) \leq \frac{\sigma^2}{t^2}$$

→ Chebyshev's inequality

$$P(|\hat{F}_n(x) - F(x)| \geq \epsilon) \leq \frac{\frac{F(x)(1-F(x))}{n}}{\frac{\epsilon^2}{1}} = \frac{F(x)(1-F(x))}{n\epsilon^2} \rightarrow 0 \text{ as } n \rightarrow \infty$$

$$\therefore \hat{F}_n(x) \xrightarrow{P} F(x)$$

추가정보는 여기에 ↓

<http://www.win.tue.nl/~rmcastro/AppStat2013/files/lecture1.pdf>

- ✓ 9. 100 people are given a standard antibiotic to treat an infection and another 100 are given a new antibiotic. In the first group, 90 people recover; in the second group, 85 people recover. Let p_1 be the probability of recovery under the standard treatment and let p_2 be the probability of recovery under the new treatment. We are interested in estimating $\theta = p_1 - p_2$. Provide an estimate, standard error, an 80 percent confidence interval, and a 95 percent confidence interval for θ .

$$\hat{p}_1 = 0.9$$

$$\hat{p}_2 = 0.85$$

$$se(\hat{p}_1) = \sqrt{\frac{v(\hat{p}_1)}{n}} = \sqrt{\frac{0.9 \times 0.1}{100}} = \sqrt{\frac{0.09}{100}}$$

$$se(\hat{p}_2) = \sqrt{\frac{v(\hat{p}_2)}{n}} = \sqrt{\frac{0.85 \times 0.15}{100}} = \sqrt{\frac{0.1275}{100}}$$

① estimate

$$\theta = p_1 - p_2 \Rightarrow \hat{p}_1 - \hat{p}_2 = 0.05$$

② standard error

$$Var(\hat{p}_1 - \hat{p}_2) = Var(\hat{p}_1) + Var(\hat{p}_2) - cov(\hat{p}_1, \hat{p}_2)$$

$$se(\hat{p}_1 - \hat{p}_2) = \sqrt{Var(\hat{p}_1) + Var(\hat{p}_2)} = 0.04663689 \approx 0.047$$

③ 80 percent interval

$$0.05 \pm 1.28 \times 0.047 \Rightarrow -0.01 \sim 0.11$$

④ 95 percent interval

$$0.05 \pm 2 \times 0.047 \Rightarrow -0.04 \sim 0.14$$