

5장. 표본분포의 근사

5.1 중심극한정리

이항 확률의 정규 근사(부록 p.502~505):

$$\sum_{x: a \leq \frac{x-np}{\sqrt{np(1-p)}} \leq b} \binom{n}{x} p^x (1-p)^{n-x} \sim \int_a^b \frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2}z^2} dz, \quad n \rightarrow \infty$$

[Key]

(1)(스털링의 공식) $\Gamma(m+1) \sim m^{m+1/2} e^{-m} \sqrt{2\pi}, \quad m \rightarrow \infty$

(2)(이항 확률밀도의 근사) $a \leq \frac{x-np}{\sqrt{npq}} \leq b$ 인 x 에 대하여

$$\binom{n}{x} p^x (1-p)^{n-x} \sim \frac{1}{\sqrt{2\pi}} \exp\left\{-\frac{1}{2}\left(\frac{x-np}{\sqrt{npq}}\right)^2\right\} \frac{1}{\sqrt{npq}}, \quad n \rightarrow \infty \quad (q \equiv 1-p)$$

(3)(누적 이항 확률의 근사)

$$\begin{aligned} \sum_{x: a \leq \frac{x-np}{\sqrt{npq}} \leq b} \binom{n}{x} p^x (1-p)^{n-x} &\sim \frac{1}{\sqrt{npq}} \sum_{x: a \leq \frac{x-np}{\sqrt{npq}} \leq b} \phi\left(\frac{x-np}{\sqrt{npq}}\right) \\ &\sim \frac{b-a}{N-1} \sum_{j=1}^N \phi(z_j) \\ &\sim \int_a^b \phi(z) dz, \quad \phi(z) = \frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2}z^2} \end{aligned}$$

정리 5.1.1: 중심극한정리(中心極限整理 central limit theorem)

확률변수 X_1, \dots, X_n 이 서로 독립이고 동일한 분포를 따르며 $\text{Var}(X_1)$ 이 양의 실수일 때

$$E(X_1) = \mu, \text{Var}(X_1) = \sigma^2 \quad (0 < \sigma < +\infty)$$

이라고 하면, 표준정규분포 $N(0,1)$ 을 따르는 확률변수 Z 에 대하여 다음이 성립한다.

$$\lim_{n \rightarrow \infty} P\left(\frac{(X_1 + \dots + X_n) - n\mu}{\sigma/\sqrt{n}} \leq x\right) = P(Z \leq x) \quad \forall x: -\infty < x < +\infty$$

[증명의 개요]

$$mgf_{\sqrt{n}(\bar{X}_n - \mu)/\sigma}(t) = [mgf_{(X_1 - \mu)/\sigma}(t/\sqrt{n})]^n$$

여기에서 $m(s) = mgf_{(X_1 - \mu)/\sigma}(s)$ 라고 하면 적률생성함수의 성질과 테일러 정리로부터

$$m\left(\frac{t}{\sqrt{n}}\right) = 1 + \frac{1}{2} \frac{t^2}{n} + R_{n,t}, \quad \lim_{n \rightarrow \infty} nR_{n,t} = 0$$

$$\log mgf_{\sqrt{n}(\bar{X}_n - \mu)/\sigma}(t) = n \log\left(1 + \frac{1}{2} \frac{t^2}{n} + R_{n,t}\right), \quad \lim_{n \rightarrow \infty} nR_{n,t} = 0$$

$$= n \left\{ \frac{1}{2} \frac{t^2}{n} + r_{n,t} \right\}, \quad \lim_{n \rightarrow \infty} nr_{n,t} = 0$$

$$\therefore \lim_{n \rightarrow \infty} mgf_{\sqrt{n}(\bar{X}_n - \mu)/\sigma}(t) = \exp(t^2/2)$$

예 5.1.1

(a)(포아송분포의 정규근사)

서로 독립이고 포아송분포 $\text{Poisson}(\lambda)$ 를 따르는 X_1, \dots, X_n 에 대하여

$$E(X_1) = \text{Var}(X_1) = \lambda \quad (0 < \lambda < +\infty)$$

이므로

$$\lim_{n \rightarrow \infty} P\left(\frac{(X_1 + \dots + X_n)/n - \lambda}{\sqrt{\lambda/n}} \leq x\right) = P(Z \leq x), \quad Z \sim N(0, 1)$$

한편 포아송분포의 성질로부터 $X_1 + \dots + X_n \sim \text{Poisson}(n\lambda)$ 이므로

$$\lim_{n \rightarrow \infty} \sum_{k: \frac{k - n\lambda}{\sqrt{n\lambda}} \leq x} \frac{e^{-n\lambda} (n\lambda)^k}{k!} = \int_{-\infty}^x \frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2}z^2} dz$$

따라서 $Y_n \sim \text{Poisson}(n)$ 이고 n 이 충분히 크면 표준정규분포의 누적분포함수 Φ 를 이용하여 다음과 같이 근사계산을 할 수 있다.¹⁾

$$P(a < Y_n \leq b) \simeq \Phi\left(\frac{b-n}{\sqrt{n}}\right) - \Phi\left(\frac{a-n}{\sqrt{n}}\right)$$

(b)(감마분포 또는 카이제곱분포의 정규근사)

서로 독립이고 감마분포 $\text{Gamma}(\alpha, \beta)$ 를 따르는 X_1, \dots, X_n 에 대하여

$$E(X_1) = \alpha\beta, \quad \text{Var}(X_1) = \alpha\beta^2 \quad (\alpha > 0, \beta > 0)$$

이므로

$$\lim_{n \rightarrow \infty} P\left(\frac{(X_1 + \dots + X_n)/n - \alpha\beta}{\sqrt{\alpha\beta^2/n}} \leq x\right) = P(Z \leq x), \quad Z \sim N(0, 1)$$

한편 $Y_n \sim \chi^2(n)$ 이면

$$Y_n \stackrel{d}{=} X_1 + \dots + X_n, \quad X_i \stackrel{iid}{\sim} \chi^2(1) = \text{Gamma}(1/2, 2) \quad (i = 1, \dots, n)$$

이므로

$$\lim_{n \rightarrow \infty} P\left(\frac{Y_n - n}{\sqrt{2n}} \leq x\right) = P(Z \leq x), \quad Z \sim N(0, 1)$$

따라서 $Y_n \sim \chi^2(n)$ 이고 n 이 충분히 크면 (a)에서와 마찬가지로

$$P(a < Y_n \leq b) \simeq \Phi\left(\frac{b-n}{\sqrt{2n}}\right) - \Phi\left(\frac{a-n}{\sqrt{2n}}\right)$$

1) 이러한 근사계산은 중심극한정리에서 정규분포로의 수렴이 균등수렴(연습문제 5.9 참조)이기 때문에 가능한 것이다.

정리 5.1.2: 다차원 경우의 중심극한정리

다차원 확률변수 $X_1 = (X_{11}, \dots, X_{1k})^t, \dots, X_n = (X_{n1}, \dots, X_{nk})^t$ 이 서로 독립이고 동일한 분포를 따르며 분산행렬 $\text{Var}(X_1)$ 이 존재할 때

$$E(X_1) = \mu, \text{Var}(X_1) = \Sigma$$

라고 하면, $N_k(0, \Sigma)$ 를 따르는 $Z = (Z_1, \dots, Z_k)^t$ 에 대하여 다음이 성립한다.²⁾

$$\lim_{n \rightarrow \infty} P\{\sqrt{n}((X_{1j} + \dots + X_{nj})/n - \mu_j) \leq x_j, j = 1, \dots, k\} = P(Z_1 \leq x_1, \dots, Z_k \leq x_k) \quad \forall x_j$$

예 5.1.2 다항분포의 다변량 정규근사:

다항분포의 대의적 정의로부터 $X_n = (X_{n1}, \dots, X_{nk})^t \sim \text{Multi}(n, (p_1, p_2, \dots, p_k)^t)$ 이면

$$X_n \stackrel{d}{=} Z_1 + \dots + Z_n, Z_i = (Z_{i1}, Z_{i2}, \dots, Z_{ik})^t \stackrel{iid}{\sim} \text{Multi}(1, (p_1, p_2, \dots, p_k)^t) \quad (i = 1, \dots, n)$$

한편 p_1, \dots, p_k 를 대각원소로 하는 대각행렬을 $D(p_j)$ 라고 하면 정리 3.2.2로부터

$$E(X_1) = p = (p_1, \dots, p_k)^t, \quad \text{Var}(X_1) = D(p_j) - pp^t,$$

따라서 중심극한정리로부터

$$\lim_{n \rightarrow \infty} P\{(X_n - np)/\sqrt{n} \leq x\} = P(Z \leq x), \quad Z \sim N_k(0, D(p_j) - pp^t)$$

이 때 분산행렬 $D(p_j) - pp^t$ 은 특이행렬로서 역행렬을 갖지 않는 점에 유의하여야 한다.

2) 다차원 확률변수에 관한 이러한 결합확률을 간략히 다음과 같이 나타내기도 한다.

$$\lim_{n \rightarrow \infty} P\{\sqrt{n}((X_1 + \dots + X_n)/n - \mu) \leq x\} = P(Z \leq x), \quad \forall x \in R^k, \quad Z \sim N_k(0, \Sigma)$$

5.2 극한분포와 확률적 수렴

예 5.2.1(확률분포를 근사할 때 모든 점에서 수렴할 것을 요구하지 않아도 되는 예)

동전을 던져서 앞면이 나오면 0부터 $1+n^{-1}$ 사이의 수를 랜덤하게 선택하여 보여주고, 뒷면이 나오면 $1+n^{-1}$ 을 보여주는 실험에서 보게 될 숫자를 $X_n (n=1,2,\dots)$ 이라고 하자.

확률변수 $X_n (n=1,2,\dots)$ 의 누적분포함수가 다음과 같은 경우에 분포의 근사를 생각해보자.

$$cdf_{X_n}(x) = P(X_n \leq x) = \begin{cases} 0, & x < 0 \\ \frac{1}{2} \frac{1}{1+n^{-1}}x, & 0 \leq x < 1+n^{-1} \\ 1, & x \geq 1+n^{-1} \end{cases}$$

이러한 누적분포함수들의 극한으로 정의되는 함수는

$$G(x) = \lim_{n \rightarrow \infty} cdf_{X_n}(x) = \begin{cases} 0, & x < 0 \\ \frac{1}{2}x, & 0 \leq x \leq 1 \\ 1, & x > 1 \end{cases}$$

로서, $x=1$ 의 오른쪽에서 연속이 아닌 함수이다. 따라서 정리 1.5.1로부터 이 함수는 누적분포함수가 아니다.

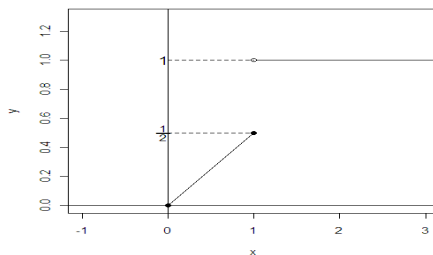


그림 1.5.1

누적분포함수의 극한이 연속이 아닌 경우

예 5.2.1에서 극한 함수가 연속이 아닌 점 $x=1$ 의 오른쪽에서 연속이도록 수정한 함수를

$$F(x) = \begin{cases} 0, & x < 0 \\ \frac{1}{2}x, & 0 \leq x < 1 \\ 1, & x \geq 1 \end{cases}$$

라고 하면 함수 $F(x)$ 는 누적분포함수이고, 예 5.2.1의 확률변수들 $X_n (n=1,2,\dots)$ 에 대하여

$$\lim_{n \rightarrow \infty} cdf_{X_n}(x) = F(x) \quad \forall x : x \neq 1$$

임을 알 수 있다. 이 경우에 동전을 던져서 앞면이 나오면 0부터 1사이의 수를 랜덤하게 선택하여 보여주고, 뒷면이 나오면 1을 보여주는 실험에서 보게 될 숫자를 X 라고 하면 함수 $F(x)$ 는 확률변수 X 의 누적분포함수이다. 따라서

$$\lim_{n \rightarrow \infty} cdf_{X_n}(x) = cdf_X(x) \quad \forall x : x \neq 1$$

이고 $X_n (n=1,2,\dots)$ 의 분포가 X 의 분포로 근사된다고 할 수 있을 것이다.

극한분포(極限分布 limiting distribution):

확률변수의 열 $X_n (n = 1, 2, \dots)$ 과 Z 에 대하여

$$\lim_{n \rightarrow \infty} P(X_n \leq x) = P(Z \leq x)$$

가 Z 의 누적분포함수 cdf_Z 가 연속인 모든 점 x 에서 성립할 때, Z 의 분포를 X_n 분포들의 극한분포 또는 점근분포 (漸近分布 asymptotic distribution)라고 하며 기호로는

$$X_n \xrightarrow[n \rightarrow \infty]{d} Z$$

로 나타낸다. 즉 cdf_Z 가 연속인 점들의 집합을 $\text{Conti}(cdf_Z)$ 라고 하면

$$X_n \xrightarrow[n \rightarrow \infty]{d} Z \Leftrightarrow \lim_{n \rightarrow \infty} cdf_{X_n}(x) = cdf_Z(x) \quad \forall x : x \in \text{Conti}(cdf_Z)$$

예 5.2.2 이항분포의 포아송근사:

이항분포 $B(n, \lambda/n) (0 < \lambda/n < 1)$ 을 따르는 확률변수 X_n 과 포아송분포 $\text{Poisson}(\lambda)$ 를 따르는 확률변수 X 에 대하여 다음이 성립하는 것은 3장 4절에서 소개되었다.

$$P(X_n = k) = \binom{n}{k} \left(\frac{\lambda}{n}\right)^k \left(1 - \frac{\lambda}{n}\right)^{n-k} \xrightarrow[n \rightarrow \infty]{} \frac{1}{k!} \lambda^k e^{-\lambda} = P(X = k), \quad k = 1, 2, \dots$$

한편 이들 확률변수들은 0 또는 자연수의 값만을 가지므로 모든 점 x 에서

$$cdf_{X_n}(x) = \sum_{k: 0 \leq k \leq x} P(X_n = k) \xrightarrow[n \rightarrow \infty]{} \sum_{k: 0 \leq k \leq x} P(X = k) = cdf_X(x)$$

임을 알 수 있다. 따라서 X 의 분포 $\text{Poisson}(\lambda)$ 는 X_n 의 분포 $B(n, \lambda/n)$ 의 극한분포이고 이를 기호로 다음과 같이 나타내기도 한다.

$$B(n, \lambda/n) \approx \text{Poisson}(\lambda), \quad n \rightarrow \infty$$

예 5.2.3

균등분포 $U(0,1)$ 에서의 랜덤표본 n 개에 기초한 순서통계량을 $U_{(1)} < \dots < U_{(n)}$ 이라고 할 때

$$P\{n(1 - U_{(n)}) \leq x\} = 1 - P\left\{U_{(n)} < 1 - \frac{x}{n}\right\}$$

이고

$$P\left\{U_{(n)} < 1 - \frac{x}{n}\right\} = \begin{cases} 0, & 1 - x/n < 0 \quad \Leftrightarrow \quad x > n \\ \left(1 - \frac{x}{n}\right)^n, & 0 \leq 1 - x/n < 1 \quad \Leftrightarrow \quad 0 < x \leq n \\ 1, & 1 \leq 1 - x/n \quad \Leftrightarrow \quad x \leq 0 \end{cases}$$

$$\lim_{n \rightarrow \infty} P\left\{U_{(n)} < 1 - \frac{x}{n}\right\} = \begin{cases} e^{-x}, & 0 < x < +\infty \\ 1, & x \leq 0 \end{cases}$$

이므로 $n(1 - U_{(n)})$ 의 극한분포를 다음과 같이 구할 수 있다.

$$\lim_{n \rightarrow \infty} P\{n(1 - U_{(n)}) \leq x\} = \begin{cases} 1 - e^{-x}, & x \geq 0 \\ 0, & x < 0 \end{cases}$$

$$\therefore n(1 - U_{(n)}) \xrightarrow[n \rightarrow \infty]{d} Z, \quad Z \sim \text{Exp}(1)$$

정리 5.2.1: 극한분포가 상수의 분포인 경우

확률변수의 열 $X_n (n=1,2,\dots)$ 의 극한분포가 상수 c 의 분포일 조건은 다음과 같다.

$$X_n \xrightarrow[n \rightarrow \infty]{d} X, P(X=c)=1 \Leftrightarrow \lim_{n \rightarrow \infty} P(|X_n - c| \geq \epsilon) = 0 \quad \forall \epsilon > 0$$

[Key] $P(X=c)=1$ 인 경우에 그 누적분포함수는

$$cdf_X(x) = \begin{cases} 0, & x < c \\ 1, & x \geq c \end{cases}$$

이므로 cdf_X 는 $x=c$ 이외의 모든 점에서 연속인 함수이다.

확률수렴(確率收斂 converge in probability): $X_n \xrightarrow[n \rightarrow \infty]{P} c$ 또는 $p \lim_{n \rightarrow \infty} X_n = c$

$$p \lim_{n \rightarrow \infty} X_n = c \Leftrightarrow \lim_{n \rightarrow \infty} P(|X_n - c| \geq \epsilon) = 0 \quad \forall \epsilon > 0$$

정리 5.2.2: 대수의 법칙(大數의 法則 law of large numbers)

확률변수 X_1, \dots, X_n 이 서로 독립이고 동일한 분포를 따르며 $E(X_1)$ 이 실수로 정의되면³⁾

$$p \lim_{n \rightarrow \infty} \frac{1}{n} \sum_{i=1}^n X_i = E(X_1)$$

[증명] 여기에서는 추가적인 조건으로서 $\text{Var}(X_1) < +\infty$ 인 경우만을 다루기로 한다.

체비셰프의 부등식을 표본평균 $\overline{X}_n = (X_1 + \dots + X_n)/n$ 에 적용하면, 임의의 양수 ϵ 에 대하여

$$\begin{aligned} 0 \leq P(|\overline{X}_n - E(X_1)| \geq \epsilon) &= P(|\overline{X}_n - E(\overline{X}_n)| \geq \epsilon) \leq \text{Var}(\overline{X}_n)/\epsilon^2 = \text{Var}(X_1)/n\epsilon^2 \\ \therefore \lim_{n \rightarrow \infty} P(|\overline{X}_n - E(X_1)| \geq \epsilon) &= 0 \end{aligned}$$

랜덤포본 X_1, \dots, X_n 을 관측할 때 관측값이 A 에 속하는 상대도수는

$$\frac{1}{n} \sum_{i=1}^n I_A(X_i)$$

로 주어진다. 이 상대도수에 대수의 법칙을 적용하면 다음이 성립하는 것을 알 수 있다.

$$p \lim_{n \rightarrow \infty} \frac{1}{n} \sum_{i=1}^n I_A(X_i) = E(I_A(X_1)) = P(X_1 \in A)$$

즉 시행 회수가 많아짐에 따라 상대도수가 확률에 한없이 가까이 가는 것을 뜻한다.

다차원 경우의 확률수렴:

다차원 확률변수 $X_n = (X_{n1}, \dots, X_{nk})^t$ 과 상수의 벡터 $c = (c_1, \dots, c_k)^t$ 에 대하여

$$\|X_n - c\| = \sqrt{(X_{n1} - c_1)^2 + \dots + (X_{nk} - c_k)^2}$$

라고 할 때

$$\lim_{n \rightarrow \infty} P(\|X_n - c\| \geq \epsilon) = 0 \quad \forall \epsilon > 0$$

이면 $X_n = (X_{n1}, \dots, X_{nk})^t$ 이 $c = (c_1, \dots, c_k)^t$ 로 확률수렴한다고 한다.

³⁾ 기댓값의 정의에 따라 이는 $E(|X_1|) < +\infty$ 임을 뜻한다.

정리 5.2.3

다차원 확률변수 $X_n = (X_{n1}, \dots, X_{nk})^t$ 이 상수의 벡터 $c = (c_1, \dots, c_k)^t$ 로 확률수렴하는 것 즉 $\text{p} \lim_{n \rightarrow \infty} X_n = c$

와 다음의 조건들은 같은 뜻이다.

$$(a) \lim_{n \rightarrow \infty} P(\max_{1 \leq i \leq k} |X_{ni} - c_i| \geq \epsilon) = 0$$

$$(b) \text{p} \lim_{n \rightarrow \infty} X_{n1} = c_1, \dots, \text{p} \lim_{n \rightarrow \infty} X_{nk} = c_k$$

[증명의 Key] (a) $\max_{1 \leq i \leq k} |X_{ni} - c_i| \leq \|X_n - c\| \leq k \max_{1 \leq i \leq k} |X_{ni} - c_i|$
 $\therefore P(\max_{1 \leq i \leq k} |X_{ni} - c_i| \geq \epsilon) \leq P(\|X_n - c\| \geq \epsilon) \leq P(\max_{1 \leq i \leq k} |X_{ni} - c_i| \geq \epsilon/k)$

$$(b) (|X_{ni} - c_i| \geq \epsilon) \subseteq (\max_{1 \leq i \leq k} |X_{ni} - c_i| \geq \epsilon) = \bigcup_{i=1}^k (|X_{ni} - c_i| \geq \epsilon) \quad (i = 1, \dots, k)$$

$$\max_{1 \leq i \leq k} P(|X_{ni} - c_i| \geq \epsilon) \leq P(\max_{1 \leq i \leq k} |X_{ni} - c_i| \geq \epsilon) \leq \sum_{i=1}^k P(|X_{ni} - c_i| \geq \epsilon)$$

정리 5.2.4 : 다차원 경우에 대수의 법칙

다차원 확률변수 $X_1 = (X_{11}, \dots, X_{1k})^t, \dots, X_n = (X_{n1}, \dots, X_{nk})^t$ 이 서로 독립이고 동일한 분포를 따르며 $E(X_1) = (E(X_{11}), \dots, E(X_{1k}))^t$ 이 정의될 수 있으면⁴⁾

$$\text{p} \lim_{n \rightarrow \infty} \frac{1}{n} \sum_{i=1}^n X_i = E(X_1)$$

[증명] 일차원 경우에 대수의 법칙과 정리 5.2.3으로부터 이 정리는 명백히 성립한다.

정리 5.2.5: 연속 함수와 확률수렴

다차원 확률변수 $X_n = (X_{n1}, \dots, X_{nk})^t$ 이 상수의 벡터 $c = (c_1, \dots, c_k)^t$ 로 확률수렴하고 실수값 함수 g 가 c 에서 연속이면 $g(X_n)$ 이 $g(c)$ 로 확률수렴한다. 즉 함수 g 가 c 에서 연속일 때

$$\text{p} \lim_{n \rightarrow \infty} X_n = c \text{ 이면 } \text{p} \lim_{n \rightarrow \infty} g(X_n) = g(c)$$

[증명] 함수 g 가 c 에서 연속이므로, 다음이 성립한다.

$$\forall \epsilon > 0, \exists \delta > 0 : (\|x - c\| < \delta \Rightarrow |g(x) - g(c)| < \epsilon)$$

따라서 임의의 양수 ϵ 에 대하여 이러한 양수 δ 를 선택하면

$$(\|X_n - c\| < \delta) \subseteq (|g(X_n) - g(c)| < \epsilon)$$

$$\therefore (|g(X_n) - g(c)| \geq \epsilon) \subseteq (\|X_n - c\| \geq \delta)$$

$$\therefore P(|g(X_n) - g(c)| \geq \epsilon) \leq P(\|X_n - c\| \geq \delta)$$

그런데 $\text{p} \lim_{n \rightarrow \infty} X_n = c$ 이므로

$$0 \leq \lim_{n \rightarrow \infty} P(|g(X_n) - g(c)| \geq \epsilon) \leq \lim_{n \rightarrow \infty} P(\|X_n - c\| \geq \delta) = 0$$

4) 기댓값의 정의에 따라 이는 $E(|X_{11}|) < +\infty, \dots, E(|X_{1k}|) < +\infty$ 임을 뜻한다.

정리 5.2.6: 확률수렴과 사칙연산⁵⁾

확률변수 $X_n, Y_n (n = 1, 2, \dots)$ 이 각각 실수 a, b 로 확률수렴할 때 다음이 성립한다.

- (a) $\text{p} \lim_{n \rightarrow \infty} (X_n + Y_n) = \text{p} \lim_{n \rightarrow \infty} X_n + \text{p} \lim_{n \rightarrow \infty} Y_n = a + b$
- (b) $\text{p} \lim_{n \rightarrow \infty} (X_n - Y_n) = \text{p} \lim_{n \rightarrow \infty} X_n - \text{p} \lim_{n \rightarrow \infty} Y_n = a - b$
- (c) $\text{p} \lim_{n \rightarrow \infty} (X_n \times Y_n) = \text{p} \lim_{n \rightarrow \infty} X_n \times \text{p} \lim_{n \rightarrow \infty} Y_n = a \times b$
- (d) $\text{p} \lim_{n \rightarrow \infty} (X_n \div Y_n) = \text{p} \lim_{n \rightarrow \infty} X_n \div \text{p} \lim_{n \rightarrow \infty} Y_n = a \div b \quad (b \neq 0)$

[증명] 정리 5.2.3으로부터 이차원 확률변수 $(X_n, Y_n)^t$ 가 벡터 $(a, b)^t$ 로 확률수렴한다. 즉

$$\text{p} \lim_{n \rightarrow \infty} (X_n, Y_n)^t = (a, b)^t$$

따라서 정리 5.2.5에 의해 연속인 이변수 함수를 적용하면 확률수렴성이 보존된다.

한편 사칙연산은 이변수 함수로서 그 정의역에서 연속인 함수이므로 이 정리가 성립한다. 구체적으로는 정리 5.2.5를 다음 함수들에 적용하면 이 정리의 결과를 얻게 된다.

$$g_1(x, y) = x + y, \quad g_2(x, y) = x - y, \quad g_3(x, y) = x \times y, \quad g_4(x, y) = x \div y$$

예 5.2.4 표본적률(標本積率 sample moment) 벡터의 확률수렴:

랜덤표본 X_1, \dots, X_n 의 함수로서

$$\widehat{m}_r = \frac{1}{n} \sum_{i=1}^n X_i^r$$

을 r 차 표본적률이라고 하며, 이는 모집단 적률인 $m_r = E(X_1^r)$ 의 추론에 사용된다. 한편 정리 1.6.2로부터 $E(|X_1^k|) < +\infty$ 이면 즉 k 차 적률 $E(X_1^k)$ 가 실수로 정의될 수 있으면 더 낮은 차의 적률이 모두 실수로 정의될 수 있는 것을 알고 있다. 따라서 대수의 법칙으로부터 $E(|X_1^k|) < +\infty$ 이면

$$\text{p} \lim_{n \rightarrow \infty} \frac{1}{n} \sum_{i=1}^n (X_i, \dots, X_i^k)^t = (E(X_1), \dots, E(X_1^k))^t$$

이로부터 표본적률의 벡터인

$$(\widehat{m}_1, \dots, \widehat{m}_k)^t = \left(\frac{1}{n} \sum_{i=1}^n X_i, \dots, \frac{1}{n} \sum_{i=1}^n X_i^k \right)^t = \frac{1}{n} \sum_{i=1}^n (X_i, \dots, X_i^k)^t$$

은 표본크기가 커짐에 따라 그 추측 대상인 모집단의 적률벡터에 한없이 가까이 간다는 것을 알 수 있다.

5) 이 정리에서 확률변수 X_n 이나 Y_n 이 각각 상수 x_n 이나 $y_n (n = 1, 2, \dots)$ 인 경우에도 그 결과가 성립한다. 이는 X_n 이나 Y_n 이 하나의 값만 갖는 상수의 확률변수들인 경우로 생각할 수 있기 때문이다.

예 5.2.5 표본분산과 표본표준편차의 확률수렴:

랜덤표본 $X_1, X_2, \dots, X_n (n \geq 2)$ 을 이용하여 모분산에 관한 추론에 사용되는 표본분산

$$S_n^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X})^2$$

은 다음과 같이 두 표본적률 $(\sum_{i=1}^n X_i/n, \sum_{i=1}^n X_i^2/n)^t$ 의 함수로 나타낼 수 있다.

$$S_n^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X})^2 = \frac{n}{n-1} \left\{ \frac{1}{n} \sum_{i=1}^n X_i^2 - \left(\frac{1}{n} \sum_{i=1}^n X_i \right)^2 \right\}$$

한편 예 5.2.4에서와 같이 대수의 법칙으로부터, $E(X_1^2) < +\infty$ 이면

$$\text{p} \lim_{n \rightarrow \infty} \frac{1}{n} \sum_{i=1}^n X_i^2 = E(X_1^2), \quad \text{p} \lim_{n \rightarrow \infty} \frac{1}{n} \sum_{i=1}^n X_i = E(X_1)$$

따라서 정리 5.2.6으로부터

$$\text{p} \lim_{n \rightarrow \infty} S_n^2 = \text{p} \lim_{n \rightarrow \infty} \frac{n}{n-1} \left\{ \frac{1}{n} \sum_{i=1}^n X_i^2 - \left(\frac{1}{n} \sum_{i=1}^n X_i \right)^2 \right\} = 1 \times \{E(X_1^2) - [E(X_1)]^2\}$$

그러므로 모분산 $\sigma^2 = \text{Var}(X_1)$ 이 양의 실수일 때

$$\text{p} \lim_{n \rightarrow \infty} S_n^2 = \sigma^2$$

또한 표본표준편차는 $S_n = \sqrt{S_n^2}$ 이고 제곱근함수는 연속함수이므로 정리 5.2.5로부터

$$\text{p} \lim_{n \rightarrow \infty} S_n = \sqrt{\text{p} \lim_{n \rightarrow \infty} S_n^2} = \sqrt{\sigma^2} = \sigma$$

정리 5.2.7

분산이 실수로 정의될 수 있는 확률변수 $X_n (n = 1, 2, \dots)$ 에 대하여 다음이 성립한다.

$$\lim_{n \rightarrow \infty} \text{Var}(X_n) = 0, \quad \lim_{n \rightarrow \infty} E(X_n) = a \quad \text{이면 } \text{p} \lim_{n \rightarrow \infty} X_n = a$$

[증명] 정리 1.6.3의 마코프 부등식으로부터, 임의의 양수 ϵ 에 대하여 다음이 성립한다.

$$P(|X_n - a| \geq \epsilon) \leq E[(X_n - a)^2] / \epsilon^2$$

한편

$$E[(X_n - a)^2] = \text{Var}(X_n) + \{E(X_n) - a\}^2$$

이므로, $\lim_{n \rightarrow \infty} \text{Var}(X_n) = 0, \lim_{n \rightarrow \infty} E(X_n) = a$ 이면

$$0 \leq \lim_{n \rightarrow \infty} P(|X_n - a| \geq \epsilon) \leq \lim_{n \rightarrow \infty} E[(X_n - a)^2] / \epsilon^2 = 0$$

예 5.2.6

균등분포 $U(0,1)$ 에서의 랜덤표본 n 개에 기초한 순서통계량을 $U_{(1)} < \dots < U_{(n)}$ 이라고 할 때 $U_{(n)}$ 의 확률밀도함수가

$$pdf_{U_{(n)}}(x) = n x^{n-1} I_{(0,1)}(x)$$

이므로

$$E(U_{(n)}) = \frac{n}{n+1}, \quad \text{Var}(U_{(n)}) = \frac{n}{n+2} - \left(\frac{n}{n+1} \right)^2$$

따라서

$$\lim_{n \rightarrow \infty} \text{Var}(U_{(n)}) = 0, \quad \lim_{n \rightarrow \infty} E(U_{(n)}) = 1$$

$$\therefore \text{p} \lim_{n \rightarrow \infty} U_{(n)} = 1$$

예 5.2.7 표본분위수(標本分位數 sample quantile)⁶⁾의 확률수렴:

모집단 분포가 연속형이고 그 누적분포함수 $F(x)$ 가 순증가함수⁷⁾일 때 누적확률이 α 가 되는 값 $F^{-1}(\alpha)$ ($0 < \alpha < 1$)를 모집단의 α 분위수(分位數 quantile)라고 한다. 한편 랜덤포본 n 개에 기초한 순서통계량을 $X_{(1)} < \dots < X_{(n)}$ 이라고 할 때

$$r_n \sim \alpha n, \quad \text{즉} \quad \lim_{n \rightarrow \infty} r_n/n = \alpha \quad (0 < \alpha < 1)$$

를 만족하는 자연수 r_n 에 대하여 $X_{(r_n)}$ 을 표본분위수라고 한다. 한편 정리 4.3.4로부터

$$X_{(r_n)} \stackrel{d}{=} h\left(\frac{1}{n}Z_1 + \dots + \frac{1}{n-r_n+1}Z_{r_n}\right), \quad Z_i \stackrel{iid}{\sim} \text{Exp}(1) \quad (i = 1, \dots, n)$$

$$h(y) = F^{-1}(1 - e^{-y}) \quad (y > 0)$$

한편 $Y_n \equiv \frac{1}{n}Z_1 + \dots + \frac{1}{n-r_n+1}Z_{r_n}$ 이라고 하면

$$E(Y_n) = \frac{1}{n} + \dots + \frac{1}{n-r_n+1}, \quad \text{Var}(Y_n) = \frac{1}{n^2} + \dots + \frac{1}{(n-r_n+1)^2}$$

$$\frac{1}{n} + \dots + \frac{1}{n-r_n+1} = \frac{1}{n} \sum_{k=0}^{r_n-1} \frac{1}{1-k/n} \sim \int_0^\alpha \frac{1}{1-x} dx = -\log(1-\alpha)$$

$$\frac{1}{n^2} + \dots + \frac{1}{(n-r_n+1)^2} = \frac{1}{n} \frac{1}{n} \sum_{k=0}^{r_n-1} \frac{1}{(1-k/n)^2} \sim \frac{1}{n} \int_0^\alpha \frac{1}{(1-x)^2} dx = \frac{1}{n} \frac{\alpha}{1-\alpha}$$

따라서

$$\lim_{n \rightarrow \infty} \text{Var}(Y_n) = \lim_{n \rightarrow \infty} \frac{1}{n} \frac{\alpha}{1-\alpha} = 0, \quad \lim_{n \rightarrow \infty} E(Y_n) = -\log(1-\alpha)$$

$$\therefore \text{p} \lim_{n \rightarrow \infty} Y_n = \text{p} \lim_{n \rightarrow \infty} \left(\frac{1}{n}Z_1 + \dots + \frac{1}{n-r_n+1}Z_{r_n} \right) = -\log(1-\alpha)$$

그러므로 함수 $h(y) = F^{-1}(1 - e^{-y})$ ($y > 0$)가 연속인 경우에 정리 5.2.5로부터

$$\text{p} \lim_{n \rightarrow \infty} h\left(\frac{1}{n}Z_1 + \dots + \frac{1}{n-r_n+1}Z_{r_n}\right) = h(-\log(1-\alpha)) = F^{-1}(\alpha)$$

이고, 정리 5.2.1의 필요조건으로부터 이러한 경우에

$$h\left(\frac{1}{n}Z_1 + \dots + \frac{1}{n-r_n+1}Z_{r_n}\right) \xrightarrow[n \rightarrow \infty]{d} F^{-1}(\alpha)$$

그런데 $X_{(r_n)} \stackrel{d}{=} h\left(\frac{1}{n}Z_1 + \dots + \frac{1}{n-r_n+1}Z_{r_n}\right)$ 이므로 이러한 경우에

$$X_{(r_n)} \xrightarrow[n \rightarrow \infty]{d} F^{-1}(\alpha)$$

따라서 누적분포함수의 역함수 F^{-1} 이 연속일 때 정리 5.2.1의 충분조건으로부터

$$\text{p} \lim_{n \rightarrow \infty} X_{(r_n)} = F^{-1}(\alpha)$$

6) 표본분위수에 대한 이러한 정의는 표본 크기가 충분히 큰 경우에 사용하는 정의이며, 이는 예 4.1.2의 표본 중앙값의 정의와 다른 것을 알 수 있다. 이러한 정의는 확률수렴과 같은 극한 성질을 간략히 밝히기 위한 것이며, 일반적인 정의에 따른 표본분위수의 경우에도 모집단의 분위수로 확률수렴하는 것이 알려져 있다.

7) 확률적분변환에 관한 정리 4.3.3에서와 마찬가지로 누적분포함수가 순증가가 아닌 경우에는 F 의 역변환을

$$F^{-1}(u) = \inf\{x : F(x) \geq u\} \quad (0 \leq u \leq 1)$$

로 정의하고, 일반적으로는 이 역변환을 이용하여 분위수를 정의한다.