

## 5장. 표본분포의 근사

### 5.3 극한분포의 계산

#### 정리 5.3.1: 슬릿츠키(Slutsky)의 정리

확률변수  $X_n, Y_n (n = 1, 2, \dots)$ ,  $Z$ 와 실수의 상수  $c$ 에 대하여

$$X_n \xrightarrow[n \rightarrow \infty]{d} Z, \quad \text{p} \lim_{n \rightarrow \infty} Y_n = c$$

이면 다음이 성립한다.

$$(a) X_n + Y_n \xrightarrow[n \rightarrow \infty]{d} Z + c$$

$$(b) X_n - Y_n \xrightarrow[n \rightarrow \infty]{d} Z - c$$

$$(c) Y_n X_n \xrightarrow[n \rightarrow \infty]{d} cZ$$

$$(d) X_n / Y_n \xrightarrow[n \rightarrow \infty]{d} Z / c \quad (c \neq 0)$$

[증명의 Key]  $\text{p} \lim_{n \rightarrow \infty} Y_n = c$  이므로  $|Y_n - c| < \epsilon$  인 경우와  $|Y_n - c| \geq \epsilon$  인 경우로 나누어

$$\lim_{n \rightarrow \infty} P(|Y_n - c| < \epsilon) = 1, \quad \lim_{n \rightarrow \infty} P(|Y_n - c| \geq \epsilon) = 0$$

임을 활용한다.

#### 예 5.3.1 스튜던트화된 표본평균의 극한분포:

모평균이  $\mu$ 이고 모표준편차가  $\sigma (0 < \sigma < +\infty)$ 인 모집단에서의 랜덤포본  $n$ 개로부터의 표본평균과 표본표준편차를 각각  $\bar{X}_n, S_n$ 이라고 할 때

$$\frac{\bar{X}_n - \mu}{S_n / \sqrt{n}}$$

를 스튜던트화(Studentized)된 표본평균이라고 한다. 한편 중심극한 정리와 예 5.2.5로부터

$$\frac{\bar{X}_n - \mu}{\sigma / \sqrt{n}} \xrightarrow[n \rightarrow \infty]{d} Z, \quad Z \sim N(0, 1) \quad \text{이고} \quad \text{p} \lim_{n \rightarrow \infty} S_n = \sigma$$

이므로 슬릿츠키의 정리로부터

$$\frac{\bar{X}_n - \mu}{S_n / \sqrt{n}} = \frac{\bar{X}_n - \mu}{\sigma / \sqrt{n}} / \frac{S_n}{\sigma} \xrightarrow[n \rightarrow \infty]{d} Z / 1 = Z, \quad Z \sim N(0, 1)$$

모집단의 분포가 무엇이든 양수의 분산이 정의될 수만 있으면

$$\lim_{n \rightarrow \infty} P(-z_{\alpha/2} < \frac{\bar{X}_n - \mu}{S_n / \sqrt{n}} \leq z_{\alpha/2}) = 1 - \alpha$$

$$\therefore P\{\bar{X}_n - z_{\alpha/2} S_n / \sqrt{n} \leq \mu < \bar{X}_n + z_{\alpha/2} S_n / \sqrt{n}\} \simeq 1 - \alpha, \quad n \rightarrow \infty$$

이므로, 모평균  $\mu$ 에 관한  $100(1 - \alpha)\%$  점근(漸近 asymptotic) 신뢰구간이 다음과 같이 주어진다.

$$\mu \in [\bar{X}_n - z_{\alpha/2} S_n / \sqrt{n}, \bar{X}_n + z_{\alpha/2} S_n / \sqrt{n})$$

### 예 5.3.2 표본분산의 극한분포:

모평균이  $\mu$ 이고 모표준편차가  $\sigma$  ( $0 < \sigma < +\infty$ )인 모집단에서의 랜덤포본  $n$ 개로부터의 표본분산  $S_n^2$ 의 표본분포에 대하여 생각해보자. 표본평균을  $\bar{X}_n$ 라고 하면

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

한편  $Y_i = (X_i - \mu)^2$  ( $i = 1, \dots, n$ )이라고 하면 이들은 서로 독립이고 동일한 분포를 따르므로, 중심극한정리로부터  $E[(X_1 - \mu)^4] < +\infty$  일 때

$$\sqrt{n} \left\{ \frac{1}{n} \sum_{i=1}^n (X_i - \mu)^2 - \sigma^2 \right\} \xrightarrow[n \rightarrow \infty]{d} W, \quad W \sim N(0, E[(X_1 - \mu)^4] - \sigma^4)$$

또한 중심극한정리와 대수의 법칙으로부터

$$\sqrt{n}(\bar{X}_n - \mu) \xrightarrow[n \rightarrow \infty]{d} Z, \quad Z \sim N(0, \sigma^2), \quad \text{p} \lim_{n \rightarrow \infty} (\bar{X}_n - \mu) = 0$$

이므로 슬릿츠키의 정리와 정리 5.2.1로부터

$$\begin{aligned} \sqrt{n}(\bar{X}_n - \mu)(\bar{X}_n - \mu) &\xrightarrow[n \rightarrow \infty]{d} 0 \times Z = 0, \quad Z \sim N(0, \sigma^2) \\ \therefore \text{p} \lim_{n \rightarrow \infty} \sqrt{n}(\bar{X}_n - \mu)^2 &= 0 \end{aligned}$$

따라서  $E[(X_1 - \mu)^4] < +\infty$  일 때,  $W \sim N(0, E[(X_1 - \mu)^4] - \sigma^4)$ 라고 하면

$$\sqrt{n} \left( \frac{1}{n} \sum_{i=1}^n (X_i - \bar{X}_n)^2 - \sigma^2 \right) = \sqrt{n} \left\{ \frac{1}{n} \sum_{i=1}^n (X_i - \mu)^2 - \sigma^2 \right\} - \sqrt{n}(\bar{X}_n - \mu)^2 \xrightarrow[n \rightarrow \infty]{d} W - 0 = W$$

그러므로  $E[(X_1 - \mu)^4] < +\infty$  일 때  $\rho_4 = E\left[\left(\frac{X_1 - \mu}{\sigma}\right)^4\right] - 3$ 이라고 하면<sup>1)</sup>

$$\sqrt{n}(S_n^2 - \sigma^2) = \sqrt{n} \left\{ \frac{1}{n} \sum_{i=1}^n (X_i - \bar{X})^2 - \sigma^2 \right\} + \frac{1}{\sqrt{n}} S_n^2 \xrightarrow[n \rightarrow \infty]{d} W, \quad W \sim N(0, (\rho_4 + 2)\sigma^4)$$

### 다차원 경우의 극한분포:

다차원 확률변수의 열  $X_n = (X_{n1}, \dots, X_{nk})^t$  ( $n = 1, 2, \dots$ )과  $Z = (Z_1, \dots, Z_k)^t$ 에 대하여

$$\lim_{n \rightarrow \infty} P(X_{n1} \leq x_1, \dots, X_{nk} \leq x_k) = P(Z_1 \leq x_1, \dots, Z_k \leq x_k)^2$$

가  $Z = (Z_1, \dots, Z_k)^t$ 의 결합누적분포함수  $cdf_Z(x_1, \dots, x_k)$ 가 연속인 모든 점  $x = (x_1, \dots, x_k)^t$

에서 성립할 때,  $Z = (Z_1, \dots, Z_k)^t$ 의 분포를  $X_n = (X_{n1}, \dots, X_{nk})^t$  분포들의 극한분포 또는 점근분포라고 하며 기호로는 일차원의 경우와 마찬가지로 다음과 같이 나타낸다.

$$X_n \xrightarrow[n \rightarrow \infty]{d} Z \quad \text{또는} \quad (X_{n1}, \dots, X_{nk})^t \xrightarrow[n \rightarrow \infty]{d} (Z_1, \dots, Z_k)^t$$

1) 여기에서 정의된  $\rho_4$ 를 모집단의 첨예도(尖銳度 kurtosis)라고 하며 이는 모집단 분포가 평균 부근에 밀집되어 있는 정도를 나타내는 측도이고, 모집단의 분포가 정규분포인 경우에 첨예도는 0이다.

2) 이러한 결합확률을 간략히  $P(Z_1 \leq x_1, \dots, Z_k \leq x_k) = P(Z \leq x)$ 로 나타내기로 한다.

### 정리 5.3.2: 연속함수와 극한분포

다차원 확률변수  $X_n = (X_{n1}, \dots, X_{nk})^t$  ( $n = 1, 2, \dots$ )과  $Z = (Z_1, \dots, Z_k)^t$ 에 대하여

$$(X_{n1}, \dots, X_{nk})^t \xrightarrow[n \rightarrow \infty]{d} (Z_1, \dots, Z_k)^t$$

일 때, 연속함수  $g$ 에 대하여

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

[증명] 이 정리의 증명은 이 책의 수준을 넘으므로 생략한다.

[유의] 정리 5.3.1이 이 정리의 특별한 경우임에 유의...

### 예 5.3.3 이항분포와 카이제곱근사:

이항분포의 정규근사로부터  $X_n \sim B(n, p)$  ( $0 < p < 1$ )일 때

$$\frac{X_n - np}{\sqrt{np(1-p)}} \xrightarrow[n \rightarrow \infty]{d} Z, \quad Z \sim N(0, 1)$$

또한 제곱 함수는 연속함수이므로 정리 5.3.2로부터

$$\begin{aligned} \left( \frac{X_n - np}{\sqrt{np(1-p)}} \right)^2 &\xrightarrow[n \rightarrow \infty]{d} Z^2, \quad Z \sim N(0, 1) \\ \therefore \frac{(X_n - np)^2}{np(1-p)} &\xrightarrow[n \rightarrow \infty]{d} W, \quad W \sim \chi^2(1) \end{aligned}$$

한편,  $X_{n1} = X_n, X_{n2} = n - X_n, p_1 = p, p_2 = 1 - p$ 라고 하면

$$(X_{n1}, X_{n2})^t \sim \text{Multi}(n, (p_1, p_2)^t)$$

이고 다음 등식이 성립함을 알 수 있다.

$$\frac{(X_{n1} - np_1)^2}{np_1} + \frac{(X_{n2} - np_2)^2}{np_2} = \frac{(X_n - np)^2}{np} + \frac{(n - X_n - n(1-p))^2}{n(1-p)} = \frac{(X_n - np)^2}{np(1-p)}$$

따라서  $(X_{n1}, X_{n2})^t \sim \text{Multi}(n, (p_1, p_2)^t)$ 이고

$$\sum_{j=1}^2 \frac{(X_{nj} - np_j)^2}{np_j} \xrightarrow[n \rightarrow \infty]{d} W, \quad W \sim \chi^2(1)$$

### 예 5.3.4 다항분포와 카이제곱근사:

예 5.3.3을 일반화 하여

$$(X_{n1}, \dots, X_{nk})^t \sim \text{Multi}(n, (p_1, \dots, p_k)^t) (p_1 + \dots + p_k = 1, p_i > 0, i = 1, \dots, k)$$

일 때 다음과 같은 카이제곱근사가 성립하는 것을 알아보자.

$$\sum_{j=1}^k \frac{(X_{nj} - np_j)^2}{np_j} \xrightarrow[n \rightarrow \infty]{d} W, \quad W \sim \chi^2(k-1)$$

여기에서  $X_n = (X_{n1}, \dots, X_{nk})^t, p = (p_1, \dots, p_k)^t$ 라고 하면 예 5.1.2에서와 같이 다차원 경우의 중심극한정리로부터

$$(X_n - np) / \sqrt{n} \xrightarrow[n \rightarrow \infty]{d} X, \quad X \sim N_k(0, D(p_j) - pp^t)$$

이제 정리 5.3.2를 이용하여 이 결과를 밝히기 위하여

$$Y_{nj} = X_{nj}, \mu_j = p_j (j = 1, \dots, r), Y_n = (Y_{n1}, \dots, Y_{nr})^t, \mu = (\mu_1, \dots, \mu_r)^t, r = k-1$$

이라고 하고,  $\mu_1, \dots, \mu_r$  을 대각원소로 하는 대각행렬을  $D(\mu_j)$ 라고 하면

$$(Y_n - n\mu) / \sqrt{n} \xrightarrow[n \rightarrow \infty]{d} Z, \quad Z \sim N_r(0, \Sigma), \quad \Sigma = D(\mu_j) - \mu\mu^t$$

또한 이차형식  $z^t \Sigma^{-1} z$ 은 연속함수이므로 정리 5.3.2로부터

$$\{(Y_n - n\mu) / \sqrt{n}\}^t \Sigma^{-1} \{(Y_n - n\mu) / \sqrt{n}\} \xrightarrow[n \rightarrow \infty]{d} Z^t \Sigma^{-1} Z, \quad Z \sim N_r(0, \Sigma)$$

정리 4.4.5로부터  $Z \sim N_r(0, \Sigma)$ 이면  $Z^t \Sigma^{-1} Z \sim \chi^2(r), r = k-1$ 이므로

$$\{(Y_n - n\mu) / \sqrt{n}\}^t \Sigma^{-1} \{(Y_n - n\mu) / \sqrt{n}\} \xrightarrow[n \rightarrow \infty]{d} W, \quad W \sim \chi^2(k-1)$$

한편 예 4.4.4에서와 같이 특수한 형태의 행렬에 관한 연산<sup>3)</sup>을 이용하여 계산하면

$$\Sigma^{-1} = (D(\mu_j) - \mu\mu^t)^{-1} = (I - D^{-1}(\mu_j)\mu\mu^t)^{-1} D^{-1}(\mu_j) = D^{-1}(\mu_j) + \frac{1}{p_k} 11^t$$

이므로

$$\begin{aligned} \{(Y_n - n\mu) / \sqrt{n}\}^t \Sigma^{-1} \{(Y_n - n\mu) / \sqrt{n}\} &= \sum_{j=1}^{k-1} \frac{(Y_{nj} - n\mu_j)^2}{n\mu_j} + \frac{(X_{nk} - np_k)^2}{np_k} \\ \therefore \sum_{j=1}^k \frac{(X_{nj} - np_j)^2}{np_j} &\xrightarrow[n \rightarrow \infty]{d} W, \quad W \sim \chi^2(k-1) \end{aligned}$$

3) 벡터  $a, b$ 에 대하여  $1 + b^t a \neq 0$ 이면

$$(I + ab^t)^{-1} = I + c ab^t, \quad c = -1 / (1 + b^t a)$$

### 정리 5.3.3: 일차근사를 이용한 극한분포 계산

다차원 확률변수  $X_n = (X_{n1}, \dots, X_{nk})^t$  ( $n = 1, 2, \dots$ ),  $Z = (Z_1, \dots, Z_k)^t$  와 벡터  $\theta = (\theta_1, \dots, \theta_k)^t$  에 대하여

$$\sqrt{n}(X_n - \theta) \xrightarrow[n \rightarrow \infty]{d} Z$$

이고 함수  $g(\theta)$ 의 일차편도함수<sup>4)</sup>  $\dot{g}(\theta)$ 가 연속함수이면 다음이 성립한다.

$$\sqrt{n}(g(X_n) - g(\theta)) \xrightarrow[n \rightarrow \infty]{d} (\dot{g}(\theta))^t Z$$

[증명] 다차원의 경우도 일차원 경우와 마찬가지로 일차원의 경우만 증명하기로 한다.

함수  $g(\theta)$ 가 미분가능하므로 다음을 만족하는  $r(x) (= r_\theta(x))$ 가 존재한다.

$$g(x) = g(\theta) + \{\dot{g}(\theta) + r(x)\}(x - \theta), \quad r(x) \xrightarrow{x \rightarrow \theta} 0$$

여기에서 잉여항의  $r(x)$ 에 대한 조건으로부터 다음이 성립하는 것을 알 수 있다.

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

따라서 임의의 양수  $\epsilon$ 에 대하여 이러한 양수  $\delta$ 를 선택하면

$$\begin{aligned} (|X_n - \theta| < \delta) &\subseteq (|r(X_n)| < \epsilon), \quad \text{즉} \quad (|r(X_n)| \geq \epsilon) \subseteq (|X_n - \theta| \geq \delta) \\ \therefore P(|r(X_n)| \geq \epsilon) &\leq P(|X_n - \theta| \geq \delta) \end{aligned}$$

한편  $\sqrt{n}(X_n - \theta) \xrightarrow[n \rightarrow \infty]{d} Z$  이므로 슬릿츠키의 정리와 정리 5.2.1로부터

$$\begin{aligned} X_n - \theta &= \frac{1}{\sqrt{n}} \times \sqrt{n}(X_n - \theta) \xrightarrow[n \rightarrow \infty]{d} 0 \times Z = 0, \\ \therefore \text{p} \lim_{n \rightarrow \infty} X_n &= \theta, \quad \text{즉} \quad \lim_{n \rightarrow \infty} P(|X_n - \theta| \geq \delta) = 0 \\ \therefore 0 &\leq \lim_{n \rightarrow \infty} P(|r(X_n)| \geq \epsilon) \leq \lim_{n \rightarrow \infty} P(|X_n - \theta| \geq \delta) = 0, \quad \text{즉} \quad \text{p} \lim_{n \rightarrow \infty} r(X_n) = 0 \end{aligned}$$

따라서

$$\sqrt{n}(g(X_n) - g(\theta)) = (\dot{g}(\theta) + r(X_n))\sqrt{n}(X_n - \theta), \quad \text{p} \lim_{n \rightarrow \infty} (\dot{g}(\theta) + r(X_n)) = \dot{g}(\theta)$$

그러므로 슬릿츠키의 정리로부터  $\sqrt{n}(g(X_n) - g(\theta)) \xrightarrow[n \rightarrow \infty]{d} \dot{g}(\theta)Z$

### 예 5.3.5 표본표준편차의 극한분포:

모평균이  $\mu$ 이고 모표준편차가  $\sigma$  ( $0 < \sigma < +\infty$ )인 모집단에서의 랜덤표본  $n$ 개로부터의 표본분산을  $S_n^2$  이라고 할 때, 예 5.3.2로부터  $E[(X_1 - \mu)^4] < +\infty$  라는 조건하에서

$$\sqrt{n}(S_n^2 - \sigma^2) \xrightarrow[n \rightarrow \infty]{d} W, \quad W \sim N(0, (\rho_4 + 2)\sigma^4), \quad \rho_4 = E\left[\left(\frac{X_1 - \mu}{\sigma}\right)^4\right] - 3$$

표본표준편차  $S_n$ 은 표본분산의 제곱근  $\sqrt{S_n^2}$  이고 제곱근함수  $g(x) = \sqrt{x}$ 의 도함수가  $\dot{g}(x) = 1/(2\sqrt{x})$  이므로, 정리 5.3.3으로부터  $E[(X_1 - \mu)^4] < +\infty$  라는 조건하에서

$$\begin{aligned} \sqrt{n}(\sqrt{S_n^2} - \sqrt{\sigma^2}) &\xrightarrow[n \rightarrow \infty]{d} \frac{1}{2\sqrt{\sigma^2}} W, \quad W \sim N(0, (\rho_4 + 2)\sigma^4) \\ \therefore \sqrt{n}(S_n - \sigma) &\xrightarrow[n \rightarrow \infty]{d} Z, \quad Z \sim N(0, (\rho_4 + 2)\sigma^2/4) \end{aligned}$$

4) 일차원 경우에는  $\dot{g}(\theta) = dg(\theta)/d\theta$ , 다차원 경우에는  $\dot{g}(\theta) = (\partial g(\theta)/\partial \theta_1, \dots, \partial g(\theta)/\partial \theta_k)^t$

### 예 5.3.6 표본상관계수의 극한분포:

모상관계수가  $\rho$  ( $-1 < \rho < 1$ )인 이변량의 모집단에서의 랜덤포본을  $(X_1, Y_1)^t, \dots, (X_n, Y_n)^t$  ( $n > 2$ )이라고 할 때, 모상관계수의 추측에 사용되는 표본상관계수(標本相關係數 sample correlation coefficient)는 다음과 같이 정의된다.

$$\hat{\rho}_n = \frac{\sum_{i=1}^n (X_i - \bar{X})(Y_i - \bar{Y})}{\sqrt{\sum_{i=1}^n (X_i - \bar{X})^2} \sqrt{\sum_{i=1}^n (Y_i - \bar{Y})^2}}$$

표본상관계수의 표본분포를 구할 때에는

$$\mu_1 = 0, \mu_2 = 0, \sigma_1 = 1, \sigma_2 = 1$$

이라고 가정해도 좋다. 또한 간략한 표현을 위하여

$$\bar{X}_n = \sum_{i=1}^n X_i/n, \bar{Y}_n = \sum_{i=1}^n Y_i/n, \overline{(XY)}_n = \sum_{i=1}^n X_i Y_i/n, \overline{(X^2)}_n = \sum_{i=1}^n X_i^2/n, \overline{(Y^2)}_n = \sum_{i=1}^n Y_i^2/n$$

라고 하면, 함수  $g(t_1, t_2, t_3, t_4, t_5) = (t_3 - t_1 t_2)(t_4 - t_1^2)^{-1/2}(t_5 - t_2^2)^{-1/2}$  에 대하여

$$\hat{\rho}_n = \frac{\overline{(XY)}_n - \bar{X}_n \bar{Y}_n}{\sqrt{\overline{(X^2)}_n - (\bar{X}_n)^2} \sqrt{\overline{(Y^2)}_n - (\bar{Y}_n)^2}} = g(\bar{X}_n, \bar{Y}_n, \overline{(XY)}_n, \overline{(X^2)}_n, \overline{(Y^2)}_n)$$

와 같이 표본상관계수를  $(X_i, Y_i, X_i Y_i, X_i^2, Y_i^2)^t$  ( $i = 1, \dots, n$ )의 평균의 함수로 나타낼 수 있다.

또한  $Z_i = (X_i, Y_i, X_i Y_i, X_i^2, Y_i^2)^t$  ( $i = 1, \dots, n$ )이라고 하면  $Z_i$  ( $i = 1, \dots, n$ )는 서로 독립이고 동일한 분포를 따르는 5차원의 확률변수이므로 그 분산행렬이 존재할 때 다음이 성립한다.

$$\sqrt{n} \left( \frac{1}{n} \sum_{i=1}^n Z_i - E(Z_1) \right) \xrightarrow[n \rightarrow \infty]{d} V, \quad V \sim N_5(0, \text{Var}(Z_1))$$

한편 추가된 가정으로부터

$$E(Z_1) = (0, 0, \rho, 1, 1)^t, \quad \rho = g(0, 0, \rho, 1, 1) = g(E(Z_1))$$

이고 함수  $g(t_1, t_2, t_3, t_4, t_5) = (t_3 - t_1 t_2)(t_4 - t_1^2)^{-1/2}(t_5 - t_2^2)^{-1/2}$ 의 편도함수들이 연속이므로, 정리 5.3.3으로부터 분산행렬  $\text{Var}(Z_1)$ 이 존재할 때<sup>5)</sup>  $\theta = E(Z_1) = (0, 0, \rho, 1, 1)^t$ 라고 하면

$$\sqrt{n}(\hat{\rho}_n - \rho) = \sqrt{n} \left( g\left(\frac{1}{n} \sum_{i=1}^n Z_i\right) - g(\theta) \right) \xrightarrow[n \rightarrow \infty]{d} (\dot{g}(\theta))^t V, \quad V \sim N_5(0, \text{Var}(Z_1))$$

여기에서 함수  $g$ 의 일차편도함수를 구하여 계산하면

$$\begin{aligned} \dot{g}(\theta) &= \dot{g}(0, 0, \rho, 1, 1) = (0, 0, 1, -\rho/2, -\rho/2)^t \\ (\dot{g}(\theta))^t V &\sim N(0, (\dot{g}(\theta))^t \text{Var}(Z_1) (\dot{g}(\theta))) = N(0, \text{Var}[(\dot{g}(\theta))^t Z_1]) \\ (\dot{g}(\theta))^t Z_1 &= X_1 Y_1 - \frac{\rho}{2} X_1^2 - \frac{\rho}{2} Y_1^2 \end{aligned}$$

이므로,  $E(X_1^4) < +\infty, E(Y_1^4) < +\infty$  일 때

$$\sqrt{n}(\hat{\rho}_n - \rho) \xrightarrow[n \rightarrow \infty]{d} W, \quad W \sim N(0, \text{Var}(X_1 Y_1 - \frac{\rho}{2} X_1^2 - \frac{\rho}{2} Y_1^2))$$

5)  $Z_1 = (X_1, Y_1, X_1 Y_1, X_1^2, Y_1^2)^t$  이므로 정리 1.6.2 으로부터  $E(X_1^4) < +\infty, E(Y_1^4) < +\infty$  이면  $Z_1$ 의 분산행렬이 존재하는 것을 알 수 있다.

### 예 5.3.7 표본상관계수와 분산안정변환:

예 5.3.6에서 이변량 정규분포  $N(\mu_1, \mu_2; \sigma_1^2, \sigma_2^2, \rho)$  ( $\sigma_1 > 0, \sigma_2 > 0, -1 < \rho < 1$ )가 모집단의 분포일 때 표본상관계수의 극한분포에 대하여 알아보자. 이 경우에는 예 4.4.3으로부터

$$Y_1 - \rho X_1 | X_1 = x_1 \sim N(0, 1 - \rho^2)$$

으로서 조건부 분포가 조건  $X_1 = x_1$ 에 의존하지 않으므로  $Y_1 - \rho X_1$ 은  $X_1$ 과 서로 독립이다.

따라서  $\mu_1 = \mu_2 = 0, \sigma_1 = \sigma_2 = 1$  인 경우에

$$T = \frac{Y_1 - \rho X_1}{\sqrt{1 - \rho^2}}$$

이라고 하면,  $X_1$ 과  $T$ 는 서로 독립이고 각각 표준정규분포  $N(0, 1)$ 을 따른다.

한편  $Y_1 = \rho X_1 + \sqrt{1 - \rho^2} T$  를 대입하여 정리하면

$$X_1 Y_1 - \frac{\rho}{2} X_1^2 - \frac{\rho}{2} Y_1^2 = \frac{\rho}{2} (1 - \rho^2) X_1^2 + (1 - \rho^2)^{3/2} X_1 T - \frac{\rho}{2} (1 - \rho^2) T^2$$

이고, 이 표현을 이용하여 계산하면  $\mu_1 = \mu_2 = 0, \sigma_1 = \sigma_2 = 1$  인 경우에

$$\text{Var}(X_1 Y_1 - \frac{\rho}{2} X_1^2 - \frac{\rho}{2} Y_1^2) = (1 - \rho^2)^2$$

$$\therefore \sqrt{n}(\hat{\rho}_n - \rho) \xrightarrow[n \rightarrow \infty]{d} W, \quad W \sim N(0, (1 - \rho^2)^2)$$

여기에서 표본상관계수  $\hat{\rho}_n$ 의 함수  $g(\hat{\rho}_n)$ 의 극한분포를 구해보면 정리 5.3.3으로부터

$$\sqrt{n}(g(\hat{\rho}_n) - g(\rho)) \xrightarrow[n \rightarrow \infty]{d} \dot{g}(\rho)W, \quad \dot{g}(\rho)W \sim N(0, [\dot{g}(\rho)]^2 (1 - \rho^2)^2)$$

이 때 극한분포의 분산을  $\rho$ 에 의존하지 않게 하는  $g(\hat{\rho}_n)$ 을  $\hat{\rho}_n$ 의 분산안정변환(分散安定變換 variance stabilizing transformation)이라고 한다. 특히

$$g(\hat{\rho}_n) = \frac{1}{2} \log \frac{1 + \hat{\rho}_n}{1 - \hat{\rho}_n}$$

은 피셔의 변환이라고 불리우며  $[\dot{g}(\rho)]^2 (1 - \rho^2)^2 = 1$ 이 되어 극한분포가  $N(0, 1)$ 이 된다. 즉

$$\sqrt{n}(g(\hat{\rho}_n) - g(\rho)) \xrightarrow[n \rightarrow \infty]{d} Z, \quad Z \sim N(0, 1)$$

### 예 5.3.8 표본분위수의 극한분포:

예 5.2.7에서는 표본분위수  $X_{(r_n)}$  ( $r_n \sim \alpha n$ ,  $0 < \alpha < 1$ )의 확률수렴에 대하여 알아보았다. 여기에서는 예 5.2.7에서의 조건이 만족될 때 표본분위수의 극한분포에 대하여 생각해 보자.

모집단의 누적분포함수의 역함수를  $F^{-1}$ 이라 하고  $h(y) = F^{-1}(1 - e^{-y})$  ( $y > 0$ )라고 하면 표본분위수의 분포는 정리 4.3.4로부터

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

한편

$$Y_n = \frac{1}{n}Z_1 + \cdots + \frac{1}{n - r_n + 1}Z_{r_n}$$

이라고 하면 그 평균과 분산이 다음과 같이 근사되는 것을 알고 있다.

$$E(Y_n) = \frac{1}{n} + \cdots + \frac{1}{n - r_n + 1} \underset{n \rightarrow \infty}{\sim} -\log(1 - \alpha)$$

$$\text{Var}(Y_n) = \frac{1}{n^2} + \cdots + \frac{1}{(n - r_n + 1)^2} \underset{n \rightarrow \infty}{\sim} \frac{1}{n} \frac{\alpha}{1 - \alpha}$$

이를 이용하여

$$W_n = \sqrt{n} \frac{Y_n - (-\log(1 - \alpha))}{\sqrt{\alpha/(1 - \alpha)}}$$

의 적률생성함수를 근사해보면<sup>6)</sup> 다음이 성립하는 것을 알 수 있다.

$$\lim_{n \rightarrow \infty} \text{mgf}_{W_n}(t) = \exp\left(\frac{1}{2}t^2\right) = \text{mgf}_W(t), \quad W \sim N(0, 1)$$

$$\therefore W_n \xrightarrow[n \rightarrow \infty]{d} W, \quad W \sim N(0, 1)$$

$$\therefore \sqrt{n}(Y_n - (-\log(1 - \alpha))) \xrightarrow[n \rightarrow \infty]{d} \sqrt{\alpha/(1 - \alpha)} W$$

따라서 정리 5.3.3으로부터 함수  $h$ 가 미분가능할 때<sup>7)</sup>

$$\sqrt{n}\{h(Y_n) - h(-\log(1 - \alpha))\} \xrightarrow[n \rightarrow \infty]{d} \dot{h}(-\log(1 - \alpha)) \sqrt{\alpha/(1 - \alpha)} W$$

그런데 이러한 조건하에서  $\dot{h}(-\log(1 - \alpha)) = (1 - \alpha)/f(F^{-1}(\alpha))$ ,  $f = \dot{F}$  이므로

$$\sqrt{n}\{X_{(r_n)} - F^{-1}(\alpha)\} \xrightarrow[n \rightarrow \infty]{d} Z, \quad Z \sim N(0, \alpha(1 - \alpha)/[f(F^{-1}(\alpha))]^2)$$

6) 연습문제 5.15에 이러한 적률 생성함수를 근사하는 과정이 소개되어 있다.

7) 확률밀도함수가  $f(x) = dF(x)/dx$  로 주어지고  $f(F^{-1}(\alpha)) > 0$  이면 이 조건이 만족된다.



## 5.4 모의실험을 이용한 근사

특정한 분포를 따르는 확률변수의 관측값을 흔히 난수(亂數 random number)라고 부르고, 최근에는 이러한 난수들을 생성할 수 있는 기능이 많은 통계패키지에 주어져 있다.

이러한 난수의 생성에 기본이 되는 것은 균등분포  $U(0,1)$ 에서의 관측값으로서 이를 균등난수(均等亂數 uniform random number)라고 하며, 부록 III에서는 패키지 R에서 균등난수를 생성하는 기능이 소개되어 있다. 또한 이러한 균등난수로부터 확률적분변환에 관한 정리 4.3.3을 이용하여 임의의 분포로부터의 난수를 생성할 수 있다.

### 예 5.4.1 로지스틱분포에서의 난수 생성:

예 4.1.5에서 소개된 로지스틱분포  $L(0,1)$ 의 확률밀도함수와 누적분포함수는 각각

$$f(z) = \frac{e^z}{(1+e^z)^2}, \quad F(z) = 1 - \frac{1}{(1+e^z)}, \quad -\infty < z < +\infty$$

로 주어진다. 이로부터 누적분포함수의 역함수는 다음과 같이 주어지는 것을 알 수 있다.

$$F^{-1}(u) = \log \frac{u}{1-u}$$

따라서 정리 4.3.3으로부터 균등분포  $U(0,1)$ 을 따르는 확률변수  $U$ 에 대하여

$$Z = \log \frac{U}{1-U} \sim L(0,1)$$

임을 알 수 있고, 이를 이용하여 로지스틱분포  $L(0,1)$ 에서의 난수를 구할 수 있다. 또한

$$\sigma Z + \mu = \sigma \log \frac{U}{1-U} + \mu \sim L(\mu, \sigma)$$

임을 이용하여 일반적인 로지스틱분포  $L(\mu, \sigma)$ 에서의 난수를 생성할 수 있다.

예 5.4.1에서와 마찬가지로, 균등분포  $U(0,1)$ 을 따르는 확률변수  $U$ 에 대하여

$$-\log(1-U) \sim \text{Exp}(1), \quad \sigma\{-\log(1-U)\} \sim \text{Exp}(\sigma)$$

임을 이용하여 지수분포  $\text{Exp}(\sigma)$ 에서의 난수를 생성할 수 있고, 표준정규분포의 누적분포함수의 역함수  $\Phi^{-1}(u)$ 에 대하여

$$\Phi^{-1}(U) \sim N(0,1), \quad \sigma\Phi^{-1}(U) + \mu \sim N(\mu, \sigma^2)$$

임을 이용하여 정규분포  $N(\mu, \sigma^2)$ 에서의 난수를 생성할 수 있다. 이러한 방법을 이용하여 여러 가지 분포에서의 난수를 생성하는 기능이 부록 III에 소개된 R에 패키지화 되어 있다.

### 정리 5.4.1: 난수를 이용한 정적분의 근사

서로 독립이고 균등분포  $U(0,1)$ 을 따르는  $U_1, \dots, U_n$ 과 구간  $[a,b]$ 에서 연속인 함수  $g(x)$ 에 대하여 다음이 성립한다.

$$X_i = (b-a)U_i + a \stackrel{iid}{\sim} U(a,b) \quad (i = 1, \dots, n)$$

$$\text{p} \lim_{n \rightarrow \infty} \frac{b-a}{n} \sum_{i=1}^n g(X_i) = \int_a^b g(x) dx$$

[증명]  $X_1, \dots, X_n$ 이 서로 독립이고 확률밀도함수가

$$pdf_{X_1}(x) = \frac{1}{b-a} I_{(a,b)}(x)$$

로 주어지는 균등분포  $U(a,b)$ 를 따르고

$$E[g(X_1)] = \frac{1}{(b-a)} \int_a^b g(x) dx$$

로 주어진다. 따라서 대수의 법칙으로부터

$$\text{p} \lim_{n \rightarrow \infty} \frac{b-a}{n} \sum_{i=1}^n g(X_i) = (b-a) E[g(X_1)] = \int_a^b g(x) dx$$

정리 5.4.1에서의 확률수렴을 이용하여 정적분의 근사값을 구하는 방법을 **몬테칼로적분(Monte Carlo integration)**이라고 한다. 즉 서로 독립적으로 생성된 균등난수  $u_1, \dots, u_n$ 을 이용하여

$$y_i = (b-a)g((b-a)u_i + a) \quad (i = 1, \dots, n)$$

$$\int_a^b g(x) dx \simeq \frac{1}{n} \sum_{i=1}^n y_i, \quad n \rightarrow \infty$$

와 같이 정적분의 근사값을 구하는 방법을 몬테칼로적분이라고 한다. 이는

$$Y_i = (b-a)g((b-a)U_i + a) \quad (i = 1, \dots, n)$$

의 관측값을 이용하여  $E(Y_1)$ 에 대한 추측을 하는 것이므로  $E(Y_1)$ 에 대한 점근신뢰구간으로 이러한 근사의 정밀도를 나타낼 수 있다.

### 예 5.4.2

적분 공식을 이용하면 정적분  $\int_1^3 x^2 dx$ 의 값이 26/3임은 잘 알고 있다. 한편 독립적으로 생성된  $U(0,1)$ 에서의 난수  $u_1, \dots, u_n$ 에 대하여

$$x_i = 2u_i + 1, y_i = 2x_i^2 \quad (i = 1, \dots, n)$$

이라고 하여

$$\int_1^3 x^2 dx \simeq \frac{2}{n} \sum_{i=1}^n x_i^2 = \frac{1}{n} \sum_{i=1}^n y_i$$

와 같이 정적분의 근사값을 구할 수 있다.

R 패키지를 이용하여 생성된 난수들을 이용하여 이 정적분의 근사값과 95% 점근신뢰구간을 구해보면 다음과 같이 주어진다. 이 결과에서 볼 수 있듯이 난수의 개수  $n$ 이 커질수록 이러한 근사의 정밀도가 좋아진다.

$n$	100	1000	10000	100000
$\bar{y}$	8.7612	8.5633	8.6515	8.6473
$\bar{y} - 1.96 s_y / \sqrt{n}$	7.8241	8.2728	8.5604	8.6185
$\bar{y} + 1.96 s_y / \sqrt{n}$	9.6983	8.8538	8.7426	8.6762

표 5.4.1 몬테칼로적분에 의한 정적분  $\int_1^3 x^2 dx (= 8.666 \dots)$ 의 근사값과 95% 오차한계

### 예 5.4.3 로지스틱분포의 분산:

로지스틱분포  $L(0,1)$ 은  $x=0$ 에 관해 대칭인 분포로서 그 분산은<sup>8)</sup>

$$\int_{-\infty}^{+\infty} x^2 \frac{e^x}{(1+e^x)^2} dx = \frac{\pi^2}{3}$$

임이 알려져 있다. 한편 예 5.4.1로부터 균등분포  $U(0,1)$ 을 따르는 확률변수  $U$ 에 대하여

$$X = \log \frac{U}{1-U} \sim L(0,1)$$

따라서 독립적으로 생성된  $U(0,1)$ 에서의 난수  $u_1, \dots, u_n$ 에 대하여

$$x_i = \log \frac{u_i}{1-u_i}, \quad y_i = x_i^2 \quad (i=1, \dots, n)$$

이라고 하여 다음과 같이 로지스틱분포  $L(0,1)$ 의 분산의 근사값을 구할 수 있다.

$$\int_{-\infty}^{+\infty} x^2 \frac{e^x}{(1+e^x)^2} dx = E(X^2) \simeq \frac{1}{n} \sum_{i=1}^n x_i^2 = \frac{1}{n} \sum_{i=1}^n y_i$$

R 패키지를 이용하여 생성된 난수들을 이용하여 이 정적분의 근사값과 95% 점근신뢰구간을 구해보면 다음과 같이 주어진다.

$n$	100	1000	10000	100000
$\bar{y}$	2.9775	3.3375	3.3201	3.2949
$\bar{y} - 1.96 s_y / \sqrt{n}$	2.0043	2.9799	3.2034	3.2584
$\bar{y} + 1.96 s_y / \sqrt{n}$	3.9506	3.6951	3.4369	3.3313

표 5.4.2 로지스틱분포  $L(0,1)$ 의 분산  $\pi^2/3 (= 3.289 \dots)$ 의 근사값과 95% 오차한계

### 예 5.4.4 표본비율의 극한분포:

모비율이  $p$  ( $0 < p < 1$ )인 베르누이시행을 독립적으로  $n$ 번 관측한 결과를  $X_1, \dots, X_n$ 이라고 할 때, 대수의법칙과 중심극한정리로부터 표본비율  $\hat{p}_n = (X_1 + \dots + X_n)/n$ 에 대하여 다음이 성립하는 것을 알고 있다.

$$p \lim_{n \rightarrow \infty} \hat{p}_n = p, \quad \frac{\hat{p}_n - p}{\sqrt{p(1-p)/n}} \xrightarrow[n \rightarrow \infty]{d} Z, \quad Z \sim N(0,1)$$

따라서 슬렛츠키의 정리와 정리 5.3.2로부터  $\frac{\hat{p}_n - p}{\sqrt{\hat{p}_n(1-\hat{p}_n)/n}} \xrightarrow[n \rightarrow \infty]{d} Z, \quad Z \sim N(0,1)$  이므로

$$\lim_{n \rightarrow \infty} P(-z_{\alpha/2} \leq \frac{\hat{p}_n - p}{\sqrt{\hat{p}_n(1-\hat{p}_n)/n}} \leq z_{\alpha/2}) = 1 - \alpha$$

$$\therefore \lim_{n \rightarrow \infty} P\{\hat{p}_n - z_{\alpha/2} \sqrt{\hat{p}_n(1-\hat{p}_n)/n} \leq p \leq \hat{p}_n + z_{\alpha/2} \sqrt{\hat{p}_n(1-\hat{p}_n)/n}\} = 1 - \alpha$$

이로부터 모비율  $p$ 에 관한  $100(1-\alpha)\%$  점근(漸近 asymptotic) 신뢰구간을 다음과 같이 나타낼 수 있다.

$$p \in [\hat{p}_n - z_{\alpha/2} \sqrt{\hat{p}_n(1-\hat{p}_n)/n}, \hat{p}_n + z_{\alpha/2} \sqrt{\hat{p}_n(1-\hat{p}_n)/n}]$$

8) 부분적분을 이용하고 무한급수로 나타내어 적분하면

$$\int_{-\infty}^{+\infty} x^2 \frac{e^x}{(1+e^x)^2} dx = 4 \int_0^{+\infty} \frac{e^{-x}}{1+e^{-x}} x dx = 4 \int_0^{+\infty} \sum_{n=1}^{\infty} (-e^{-x})^{n-1} e^{-x} x dx = 4 \sum_{n=1}^{\infty} \frac{(-1)^{n-1}}{n^2}$$

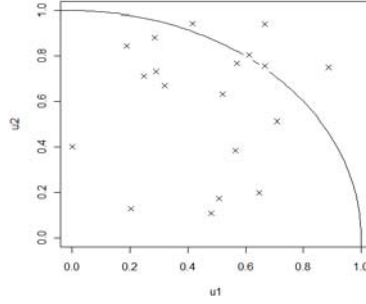
한편 후리에급수를 이용한 다음의 항등식에서  $x=0$ 을 대입하면 위의 무한급수의 값을 얻을 수 있다.

$$x^2 = \frac{1}{3} \pi^2 + 4 \sum_{n=1}^{\infty} \frac{(-1)^n}{n^2} \cos(nx) \quad (-\pi \leq x \leq \pi)$$

예 5.4.5 표본비율의 극한분포를 이용한 원주율  $\pi$ 의 근사:

서로 독립이고 균등분포  $U(0,1)$ 을 따르는 두 확률변수  $U_1, U_2$ 에 대하여 이들을 좌표로 하는 점  $(U_1, U_2)$ 을 관측한다는 것은 각 변의 길이가 1인 정사각형 내에 한 점을 랜덤하게 떨어뜨리는 것으로 생각할 수 있다. 이 실험에서 떨어뜨린 점이 사분원내에 떨어질 확률은

$$P(U_1^2 + U_2^2 \leq 1) = \pi/4$$



이제 이러한 실험을 독립적으로 반복하여 랜덤 그림 5.4.1 랜덤하게 떨어뜨린 점들  $(u_{i1}, u_{i2})$  하게 떨어뜨린 점들이 사분원내에 떨어지는 상대도수가 이 사건의 확률  $\pi/4$ 에 가까워진다는 것이 대수의 법칙의 뜻이다. 즉 서로 독립이고 균등분포  $U(0,1)$ 을 따르는 두 확률변수의 순서쌍들인  $(U_{i1}, U_{i2}) (i = 1, \dots, n)$ 를 독립적으로 관측하는 실험에서 상대도수

$$\hat{p}_n = \frac{1}{n} \sum_{i=1}^n I_{(U_{i1}^2 + U_{i2}^2 \leq 1)}$$

은  $p = P(U_1^2 + U_2^2 \leq 1) = \pi/4$ 에 확률수렴한다.

따라서 독립적으로 생성된  $U(0,1)$ 에서의 난수  $u_{11}, u_{12}; \dots; u_{n1}, u_{n2}$ 에 대하여

$$x_i = u_{i1}^2 + u_{i2}^2 \quad (i = 1, \dots, n)$$

이라고 하면 표본비율의 관측값

$$\hat{p}_n = \frac{1}{n} \sum_{i=1}^n I_{(x_i \leq 1)}$$

은 확률  $p = P(U_1^2 + U_2^2 \leq 1) = \pi/4$ 의 근사값이며  $4\hat{p}_n$ 을  $4p = \pi$ 의 근사값으로 사용할 수 있다. 또한 예 5.4.4로부터  $4p$ 에 관한 점근신뢰구간으로 근사의 정밀도를 나타낼 수 있다.

R 패키지를 이용하여 생성된 난수들을 이용하여  $4p = \pi$ 의 근사값과 95% 점근신뢰구간을 구해보면 다음과 같이 주어진다.

$n$	100	1000	10000	100000
$4\hat{p}_n$	3.3600	3.1320	3.1344	3.1452
$4\hat{p}_n - 4 \times 1.96 \sqrt{\hat{p}_n(1-\hat{p}_n)/n}$	3.0755	3.0308	3.1024	3.1351
$4\hat{p}_n + 4 \times 1.96 \sqrt{\hat{p}_n(1-\hat{p}_n)/n}$	3.6445	3.2331	3.1664	3.1552

표 5.4.3  $4p = \pi (= 3.14159 \dots)$ 의 근사값과 95% 오차한계