

독립성 검정

정의

- 개요

1. X 와 Y 가 평균 μ_1, μ_2 , 분산이 σ_1^2, σ_2^2 이고 상관계수가 ρ 인 이변량 정규분포를 따른다고 하자.

2. $H_0 : \text{가설 } \rho = 0 \text{ VS } H_1 : \rho \neq 0$ 를 검정하기 위한 우도비 검정을 정의하면

1) $L(\theta; x_n, y_n) = f(x_1, y_1; \theta)f(x_2, y_2; \theta) \cdots f(x_n, y_n; \theta)$ 이고(단, $f(x_n, y_n; \theta)$ 은 이변량 정규분포의 PDF)

2) $\Lambda = \frac{L(\rho=0; x_n, y_n)}{L(\rho \neq 0; x_n, y_n)} = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum (x_i - \bar{x})^2 \sum (y_i - \bar{y})^2}} = R$ 이다.

3. 따라서 $P[|R| \geq c]$ 에 대해서, 이를 알려진 분포로 변환하는 $g(|R|)$ 이 요구된다.

1) Y 에서 iid인 확률표본 $Y_1 \cdots Y_n$ 을 선출하고,
이 때 X 의 확률표본 $X_1 \cdots X_n$ 이 실현값 $x_1 \cdots x_n$ 을 각각 갖는다고 하자.

2) 이 가정하에서 $Y_i | (X_1 = x_1 \cdots X_n = x_n)$ 의 pdf를 구하면

(1) $\rho = 0$ 라는 가정 하에, Y_i 와 X_i 는 두 확률변수의 PDF의 단순결합으로 분리 가능하므로

(2) 이 조건부 PDF는 **단지 Y_i 의 PDF로만 도출**된다.

정의

- 개요

3. 한편 Y_i 도 iid임을 가정했으므로, 이 결합 pdf는 $\left[\frac{1}{\sqrt{2\pi}\sigma_y}\right]^n \exp\left(-\frac{\sum(y_i-\mu_y)^2}{2\sigma_y^2}\right)$ 이다.

1) 즉, $\rho = 0$ 일 경우 $Y_i|(X_1=x_1 \cdots X_n=x_n) = Y_i$ 이다.

2) 이 때, $(X_1=x_1 \cdots X_n=x_n)$ 이 주어졌을 때의 조건부 상관계수 R_c 는

$$(1) R_c = \frac{\sum(X_i-\bar{X})\{[Y_i|(X_1=x_1 \cdots X_n=x_n)]-\bar{Y}\}}{\sqrt{\sum(X_i-\bar{X})^2 \sum([Y_i|(X_1=x_1 \cdots X_n=x_n)]-\bar{Y})^2}} \text{ 이고,}$$

$$(2) \rho = 0 \text{ 일 경우 } \frac{\sum(X_i-\bar{X})\{[Y_i|(X_1=x_1 \cdots X_n=x_n)]-\bar{Y}\}}{\sqrt{\sum(X_i-\bar{X})^2 \sum([Y_i|(X_1=x_1 \cdots X_n=x_n)]-\bar{Y})^2}} = \frac{\sum(X_i-\bar{x})(Y_i-\bar{Y})}{\sqrt{\sum(X_i-\bar{x})^2 \sum(Y_i-\bar{Y})^2}} \text{ 이다.}$$

정의

- 개요

4. 이제, $Y_i | (X_1 = x_1 \cdots X_n = x_n) = Y_i$ 를 통해 Y_i 는 X_i 와 확률적으로 무관한 독립 분포가 되었으며

$$1) \quad g(|R|) = \frac{R_c \sqrt{\sum (Y_i - \bar{Y})^2}}{\sqrt{\sum (x_i - \bar{x})^2}} = \frac{\sum (x_i - \bar{x})(Y_i - \bar{Y})}{\sum (X_i - \bar{x})^2} = \frac{\sum (x_i - \bar{x}) Y_i}{\sum (X_i - \bar{x})^2} \text{ 를 도출할 수 있다.}$$

2) 이 때, $\frac{\sum (x_i - \bar{x}) Y_i}{\sum (X_i - \bar{x})^2}$ 는 회귀분석에서 파라미터 β 의 불편추정량 $\hat{\beta}$ 와 같으며,

(1) 회귀분석에서 전개한 논리를 가져오면

$$(2) \quad \frac{R_c \sqrt{\sum (Y_i - \bar{Y})^2} / \sqrt{\sum (x_i - \bar{x})^2}}{\sqrt{R_c \sum \{Y_i - \bar{Y} - \left[\frac{R_c \sqrt{\sum (Y_j - \bar{Y})^2} / \sqrt{\sum (x_j - \bar{x})^2} \right] (x_i - \bar{x})\}^2} / (n-2) \sqrt{\sum (x_j - \bar{x})^2}}} \sim T(n-2) \text{ 이다.}$$

정의

- 개요

4. 이제, $Y_i | (X_1 = x_1 \cdots X_n = x_n) = Y_i$ 를 통해 Y_i 는 X_i 와 확률적으로 무관한 독립 분포가 되었으며

1) $g(|R|) = \frac{R_c \sqrt{\sum (Y_i - \bar{Y})^2}}{\sqrt{\sum (x_i - \bar{x})^2}} = \frac{\sum (x_i - \bar{x})(Y_i - \bar{Y})}{\sum (X_i - \bar{x})^2} = \frac{\sum (x_i - \bar{x})Y_i}{\sum (X_i - \bar{x})^2}$ 를 도출할 수 있다.

2) 이 때, $\frac{\sum (x_i - \bar{x})Y_i}{\sum (X_i - \bar{x})^2}$ 는 회귀분석에서 파라미터 β 의 불편추정량 $\hat{\beta}$ 와 같으며,

(1) 회귀분석에서 전개한 논리를 가져오면

$$(2) \frac{R_c \sqrt{\sum (Y_i - \bar{Y})^2} / \sqrt{\sum (x_i - \bar{x})^2}}{\sqrt{R_c^2 \sum \{Y_i - \bar{Y} - \left[\frac{R_c \sqrt{\sum (Y_j - \bar{Y})^2} / \sqrt{\sum (x_j - \bar{x})^2} \right] (x_i - \bar{x})\}^2} / (n-2) \sqrt{\sum (x_j - \bar{x})^2}}} = \frac{R_c \sqrt{n-2}}{\sqrt{1-R_c^2}} \sim T(n-2) \text{ 이다.}$$