

Rozkład normalny wielowymiarowy

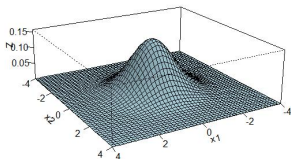
Joanna Czarnowska¹

¹Uniwersytet Gdański
Instytut Matematyki

Rozkład normalny standardowy dwuwymiarowy

Dwuwymiarowy Rozkład Normalny

$\rho = 0$



$$(X, Y) \sim N(0, 0, 1, 1, 0)$$

$$\text{Gęstość: } f(x, y) = \frac{1}{2\pi} e^{-\frac{x^2+y^2}{2}}, \quad (x, y) \in \mathbb{R}^2$$

Wartość oczekiwana $\mu = (EX, EY) = (0, 0)$.

Macierz kowariancji

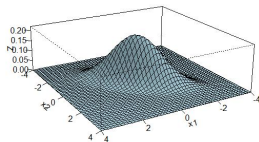
$$\begin{aligned} \Sigma &= \begin{bmatrix} \text{cov}(X, X) & \text{cov}(X, Y) \\ \text{cov}(Y, X) & \text{cov}(Y, Y) \end{bmatrix} \\ &= \begin{bmatrix} \text{Var}(X) & \text{cov}(X, Y) \\ \text{cov}(Y, X) & \text{Var}(Y) \end{bmatrix} = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} = I_2 \end{aligned}$$

Uwaga. Zobacz Przykład 1a) w skrypcie 07GaussWielowymiarowy-wykład.R

Rozkład normalny dwuwymiarowy $(X, Y) \sim N(\mu_1, \mu_2, \sigma_1, \sigma_2, \rho)$

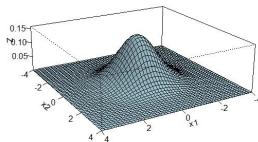
Dwuwymiarowy Rozkład Normalny

$\rho = -0.7$



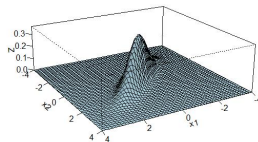
Dwuwymiarowy Rozkład Normalny

$\rho = 0$



Dwuwymiarowy Rozkład Normalny

$\rho = 0.9$



Gęstość

$$f(x, y) = \frac{1}{2\pi\sigma_1\sigma_2\sqrt{1-\rho^2}} \exp\left(-\frac{1}{2(1-\rho^2)}\left[\frac{(x-\mu_1)^2}{\sigma_1^2} - 2\rho\frac{(x-\mu_1)(y-\mu_2)}{\sigma_1\sigma_2} + \frac{(y-\mu_2)^2}{\sigma_2^2}\right]\right)$$

Wartość oczekiwana: $\mu = (EX, EY) = (\mu_1, \mu_2)$

Macierz kowariancji

$$\Sigma = \begin{bmatrix} \text{cov}(X, X) & \text{cov}(X, Y) \\ \text{cov}(Y, X) & \text{cov}(Y, Y) \end{bmatrix} = \begin{bmatrix} \sigma_1^2 & \rho\sigma_1\sigma_2 \\ \rho\sigma_1\sigma_2 & \sigma_2^2 \end{bmatrix}$$

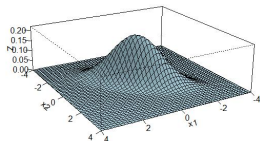
gdzie σ_1, σ_2 to odchylenia standardowe i ρ współczynnik korelacji

Rozkład normalny dwuwymiarowy

Gęstości i próby wygenerowane z rozkładów: $N(0, 0, 1, 1, -0.7)$, $N(0, 0, 1, 1, 0)$ i $N(0, 0, 1, 1, 0.9)$.

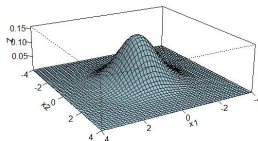
Dwuwymiarowy Rozkład Normalny

$\rho = -0.7$



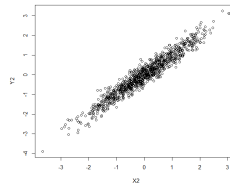
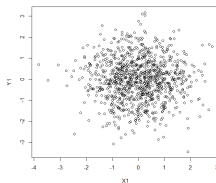
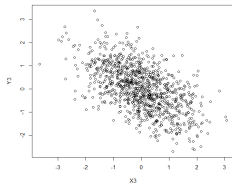
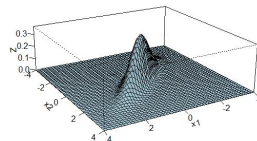
Dwuwymiarowy Rozkład Normalny

$\rho = 0$



Dwuwymiarowy Rozkład Normalny

$\rho = 0.9$



Uwaga. Zobacz Przykład 1b) w skrypcie 07GaussWielowymiarowy-wykład.R

Rozkłady brzegowe

Fakt. Jeśli $(X, Y) \sim N(\mu_1, \mu_2, \sigma_1, \sigma_2, \rho)$, to $X \sim N(\mu_1, \sigma_1)$ i $Y \sim N(\mu_2, \sigma_2)$

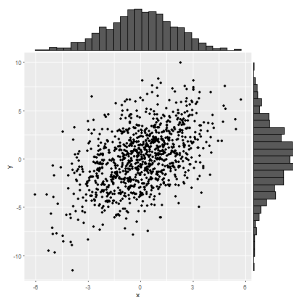
$$f_1(x) = \int_{-\infty}^{\infty} f(x, y) dy = \frac{1}{\sqrt{2\pi}\sigma_1} e^{-\frac{(x-\mu_1)^2}{2}}, \quad x \in \mathbb{R}$$

$$f_2(y) = \int_{-\infty}^{\infty} f(x, y) dx = \frac{1}{\sqrt{2\pi}\sigma_2} e^{-\frac{(y-\mu_2)^2}{2}}, \quad y \in \mathbb{R}$$

gdzie f_1 , f_2 , f są odpowiednio gęstościami zmiennych losowych X , Y oraz wektora (X, Y) .

Na wykresie obok próba licząca 1000, wygenerowana z rozkładu $N(0, 0, 2, 3, 0.5)$ wraz z histogramami rozkładów brzegowych.

Uwaga. Zobacz Przykład 2 w skrypcie 07GaussWielowymiarowy-wykład.R



Estymator kowariancji

Mamy próbę $(X_1, Y_1), \dots, (X_n, Y_n)$ prostą z rozkładu o nieznanej kowariancji i nieznanym współczynniku korelacji ρ .

Estymator kowariancji

$$\tilde{S}_{xy} = \frac{1}{n} \sum_{i=1}^n (X_i - \bar{X}_n)(Y_i - \bar{Y}_n),$$

gdzie $\bar{X}_n = \frac{1}{n} \sum_{i=1}^n X_i$, $\bar{Y}_n = \frac{1}{n} \sum_{i=1}^n Y_i$.

Estymator \tilde{S}_{xy} jest obciążonym estymatorem kowariancji, nieobciążonym jest

$$S_{xy} = \frac{n}{n-1} \tilde{S}_{xy} = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X}_n)(Y_i - \bar{Y}_n).$$

Estymator współczynnika korelacji

Estymator współczynnika korelacji

$$\hat{\rho} = R_{xy} = \frac{S_{xy}}{S_x \cdot S_y} = \frac{\sum_{i=1}^n (X_i - \bar{X}_n)(Y_i - \bar{Y}_n)}{\sqrt{\sum_{i=1}^n (X_i - \bar{X}_n)^2} \sqrt{\sum_{i=1}^n (Y_i - \bar{Y}_n)^2}},$$

gdzie S_x^2 i S_y^2 są nieobciążonymi estymatorami wariancji

$$S_x^2 = \hat{\sigma}_1^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X}_n)^2, \quad S_y^2 = \hat{\sigma}_2^2 = \frac{1}{n-1} \sum_{i=1}^n (Y_i - \bar{Y}_n)^2.$$

Rozkład normalny dwuwymiarowy – estymacja parametrów

Estymatorami największej wiarygodności wektora μ i macierzy kowariancji Σ , dwuwymiarowego rozkładu normalnego są $\hat{\mu} = (\overline{X}_n, \overline{Y}_n)$ oraz

$$\hat{\Sigma} = \begin{bmatrix} \tilde{S}_x^2 & \tilde{S}_{xy} \\ \tilde{S}_{xy} & \tilde{S}_y^2 \end{bmatrix}, \quad (1)$$

gdzie \tilde{S}_x^2 i \tilde{S}_y^2 są obciążonymi estymatorami wariancji

$$\tilde{S}_x^2 = \frac{1}{n} \sum_{i=1}^n (X_i - \overline{X}_n)^2, \quad \tilde{S}_y^2 = \frac{1}{n} \sum_{i=1}^n (Y_i - \overline{Y}_n)^2.$$

Jest to estymator obciążony ($E(\hat{\Sigma}) = \frac{n-1}{n} \Sigma$), nieobciążonym estymatorem jest

$$\hat{\Sigma} = \begin{bmatrix} S_x^2 & S_{xy} \\ S_{xy} & S_y^2 \end{bmatrix}. \quad (2)$$

Uwaga. Zobacz Przykład 3a) w skrypcie 07GaussWielowymiarowy-wyklad.R

Odległość Mahalanobisa

Dla danego wektora $x = (x_1, x_2)$, wektora średnich $\mu = (\mu_1, \mu_2)$ oraz macierzy kowariancji Σ , definiujemy **odległość Mahalanobisa** wektora x od wektora μ , jako

$$\begin{aligned} d(x, \mu) &= \left(\begin{bmatrix} x_1 - \mu_1 & x_2 - \mu_2 \end{bmatrix} \Sigma^{-1} \begin{bmatrix} x_1 - \mu_1 \\ x_2 - \mu_2 \end{bmatrix} \right)^{1/2} \\ &= \left((x - \mu)^T \Sigma^{-1} (x - \mu) \right)^{1/2}, \end{aligned}$$

gdzie Σ^{-1} jest macierzą odwrotną do macierzy Σ .

Jeśli macierz Σ jest macierzą jednostkową, to odległość Mahalanobisa jest zwykłą odległością euklidesową

$$d(x, \mu) = \left((x - \mu)^T (x - \mu) \right)^{1/2} = \sqrt{(x_1 - \mu_1)^2 + (x_2 - \mu_2)^2}.$$

Uwaga. Zobacz

- ▶ Przykład 3b) w skrypcie 07GaussWielowymiarowy-wyklad.R
- ▶ https://pl.wikipedia.org/wiki/Odległość_Mahalanobisa

Testowanie normalności rozkładu dwuwymiarowego

Mamy realizację $(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)$ próby z rozkładu dwuwymiarowego o nieznannej dystrybuancie F , wartości oczekiwanej $\mu = (\mu_1, \mu_2)$ i macierzy kowariancji Σ . Chcemy zweryfikować hipotezę o normalności rozkładu F . Jedną z metod weryfikacji wykorzystuje odległość Mahalanobisa.

Fakt. Jeśli $X = (X_1, X_2) \sim N(\mu, \Sigma)$, to kwadrat odległości Mahalanobisa wektora X od wektora μ ma rozkład chi-kwadrat o $n = 2$ stopniach swobody

$$d^2(X, \mu) = (X - \mu)^T \Sigma^{-1} (X - \mu) \sim \chi^2(2). \quad (3)$$

Mając dane $(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)$, obliczamy

- ▶ $\hat{\mu}$ oraz macierz kowariancji $\hat{\Sigma}$,
- ▶ d_1, d_2, \dots, d_n – kwadraty odległości Mahalanobisa punktów $(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)$ od $\hat{\mu}$.

Odległości d_1, d_2, \dots, d_n wykorzystujemy do utworzenia wykresów diagnostycznych oraz weryfikacji hipotezy, że (3) ma rozkład $\chi^2(2)$. Odrzucenie tej hipotezy, daje podstawy do odrzucenia hipotezy o normalności F . Brak odrzucenia – nie daje takich podstaw.

Zobacz Przykład 3c)d) w skrypcie 07GaussWielowymiarowy-wyklad.R

Testowanie normalności rozkładu – test Mardii

Klasycznym testem normalności, dla rozkładu dwuwymiarowego, jest test Mardii oparty na teście skośności i kurtozy rozkładu. Statystyki testowe to

$$A = \frac{1}{6n} \sum_{i=1}^n \sum_{j=1}^n D_{ij}^3, \quad B = \sqrt{n} \left(\frac{1}{8n} \sum_{i=1}^n D_i^2 - 1 \right),$$

gdzie

$$D_{ij} = \begin{bmatrix} x_i - \mu_1 & y_i - \mu_2 \end{bmatrix} \widehat{\Sigma}^{-1} \begin{bmatrix} x_j - \mu_1 \\ y_j - \mu_2 \end{bmatrix},$$
$$D_i = \begin{bmatrix} x_i - \mu_1 & y_i - \mu_2 \end{bmatrix} \widehat{\Sigma}^{-1} \begin{bmatrix} x_i - \mu_1 \\ y_i - \mu_2 \end{bmatrix}.$$

Jeśli hipoteza zerowa o normalności rozkładu dwuwymiarowego jest prawdziwa, to asymptotyczne rozkłady statystyk wynoszą

$$A \sim \chi^2(4), \quad B \sim N(0, 1).$$

Duże wartości statystyk świadczą przeciwko hipotezie zerowej.

Uwaga. Zobacz Przykład 4 w skrypcie 07GaussWielowym-wyklad.R