

Badanie zgodności z rozkładem dla kopuł

Praca magisterska pisana pod kierunkiem dra G. Wyłupka

Magdalena Trafidło

Uniwersytet Wrocławski, 2023

Badanie zgodności z rozkładem

Niech $X = (X_1, \dots, X_n) \sim P$.

Testujemy

$$H_0 : P \in \mathcal{P}_0 \text{ vs } H_A : P \notin \mathcal{P}_0,$$

lub równoważnie

$$H_0 : F \in \mathcal{F}_0 \text{ vs } H_A : F \notin \mathcal{F}_0,$$

gdzie

F – dystrybuanta rozkładu P ,

$\mathcal{P}_0 = \{\mathcal{P}_\theta : \theta \in \Theta\}$ – rodzina rozkładów,

$$\mathcal{F}_0 = \{\mathcal{F}_\theta : \theta \in \Theta\},$$

F_n – dystrybuanta empiryczna.

Statystyka Kolmogorova-Smirnova:

$$D_n(X) = \sup_{u \in \mathbb{R}} |F_n(u, X) - \hat{F}_0(u)|.$$

Statystyka Craméra-von Misesa:

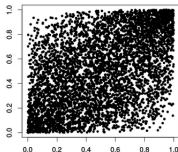
$$W^2 = n \int_{\mathbb{R}} (F_n(x) - \hat{F}_0(x))^2 dF_0(x).$$

Dystrybuanta empiryczna:

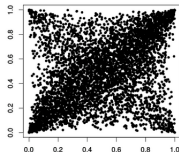
$$F_n(u) = \frac{1}{n} \sum_{i=1}^n 1(X_i \leq u).$$

Definicja

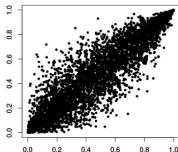
Dystrybuante C o rozkładach brzegowych jednostajnych na $(0,1)$ nazywamy kopułą.



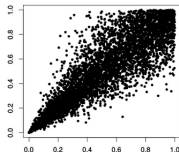
Bivariate Gaussian copula with $\rho = 0.5$



Bivariate Student-t copula with $\rho = 0.5$ and dof = 1



Bivariate Gumbel copula with $\alpha = 4$



Bivariate Clayton copula with $\alpha = 5$

Niech $X = (X_1, \dots, X_d)$, gdzie H - dystrybuanta rozkładu łącznego X , a F_1, \dots, F_d - dystrybuanty rozkładów brzegowych.

Twierdzenie (Sklar)

Jeśli rozkłady brzegowe są ciągłe, to istnieje jednoznacznie wyznaczona kopuła C taka, że

$$H(x_1, \dots, x_d) = C(F_1(x_1), \dots, F_d(x_d))$$

dla dowolnych x_i , $i = 1, \dots, d$.

Niech $X = (X_1, \dots, X_n)$, gdzie $X_i = (X_{1i}, \dots, X_{di}), i = 1, \dots, n$.

Testujemy

$$H_0 : C \in C_0 \text{ vs } H_A : C \notin C_0,$$

gdzie $C_0 = \{C_\theta : \theta \in \Theta\}$ to pewna ustalona rodzina kopuł.

Zamiast korzystać z obserwacji X_1, \dots, X_n , skorzystamy z pseudo-obserwacji U_1, \dots, U_n , gdzie

$$U_{ij} = R_{ij}/(n+1) = n\hat{F}_j(X_{ij})/(n+1),$$

$$\hat{F}_j(t) = \frac{1}{n} \sum_{i=1}^n 1(X_{ij} \leq t) - \text{dystrybuanta empiryczna.}$$

Testy oparte na empirycznych kopułach

Niech

$$\mathbb{C}_n = \sqrt{n}(C_n - C_{\theta_n}),$$

gdzie

C_{θ_n} – estymacja kopuły C przy założeniu H_0 ,

$$C_n(\mathbf{u}) = \frac{1}{n} \sum_{i=1}^n 1(U_{i1} \leq u_1, \dots, U_{id} \leq u_d) \text{ – empiryczna kopuła.}$$

Wtedy

$$S_n = \int_{[0,1]^d} \mathbb{C}_n(\mathbf{u})^2 dC_n(\mathbf{u}),$$

$$T_n = \sup_{\mathbf{u} \in [0,1]^d} |\mathbb{C}_n(\mathbf{u})|.$$

Rozważamy następujące przekształcenie:

$$X \mapsto V = H(X) = C(\mathcal{U}_1, \dots, \mathcal{U}_d), \text{ gdzie } \mathcal{U}_i = F_i(X_i).$$

Niech K - jednowymiarowa dystrybuanta V . K można wyestymować przy użyciu $V_1 = C_n(\mathcal{U}_1), \dots, V_n = C_n(\mathcal{U}_n)$:

$$K_n(v) = \frac{1}{n} \sum_{i=1}^n 1(V_i \leq v), \quad v \in [0, 1].$$

Przy H_0 wektor $\mathcal{U} = (\mathcal{U}_1, \dots, \mathcal{U}_d) \sim C_\theta$, a $C_\theta(\mathcal{U}) \sim K_\theta$.

Testy oparte o transformacje Kendalla

Niech $H'_0 : K \in K_0 = \{K_\theta : \theta \in \Theta\}$. Zauważmy, że $H_0 \subset H'_0$, tak więc nieodrzućenie H'_0 nie jest równoważne z akceptacją H_0 .

Niech $\mathbb{K}_n = \sqrt{n}(K_n - K_{\theta_n})$. Wtedy:

$$S_n^{(K)} = \int_0^1 \mathbb{K}_n(v)^2 dK_{\theta_n}(v),$$

$$T_n^{(K)} = \sup_{v \in [0,1]} |\mathbb{K}_n(v)|.$$

Definicja

Transformacja Rosenblatta kopuły C to funkcja $\mathcal{R} : (0, 1)^d \rightarrow (0, 1)^d$, która dla każdego $\mathbf{u} = (u_1, \dots, u_d) \in (0, 1)^d$ przypisuje wektor $\mathcal{R}(\mathbf{u}) = (e_1, \dots, e_d)$, gdzie $e_1 = u_1$ oraz

$$e_i = C(u_i | u_1, \dots, u_{i-1}) \text{ dla } i \in \{2, \dots, d\},$$

Kluczowa własność transformacji Rosenblatta: \mathcal{U} ma rozkład C
 $\iff \mathcal{R}(\mathcal{U})$ to d -wymiarowa niezależna kopuła C_\perp :

$$C_\perp(\mathbf{u}) = u_1 \times \dots \times u_d.$$

To znaczy, że testowanie $H_0 : \mathcal{U} \sim C \in \mathcal{C}_0$ jest równoważne z testowaniem $H_0^* : \mathcal{R}_\theta(\mathcal{U}) \sim C_\perp$.

Przy H_0 zmienne $E_1 = \mathcal{R}_{\theta_n}(U_1), \dots, E_n = \mathcal{R}_{\theta_n}(U_n) \sim C_{\perp}$ i mają w przybliżeniu jednostajny rozkład na $(0, 1)^d$.

Niech

$$\chi_i = \sum_{j=1}^d [\Phi^{-1}(E_{ij})]^2 \sim G, i \in \{1, \dots, n\},$$

gdzie G to dystrybuanta rozkładu χ^2 z d -stopniami swobody,

$$\text{a } G_n(t) = \frac{1}{n} \sum_{i=1}^n 1(\chi_i \leq t) - \text{odpowiedni rozkład empiryczny.}$$

Niech

$$\mathbb{G}_n = \sqrt{n}(G_n - G).$$

Przy założeniu, że \mathbb{G} zachowuje się asymptotycznie tak, jakby E_1, \dots, E_n faktycznie miały rozkład jednostajny, możemy skonstruować statystykę testową:

$$A_n = -n - \frac{1}{n} \sum_{i=1}^n (2i - 1) [\log\{G(\chi_{(i)})\} + \log\{1 - G(\chi_{(n+1-i)})\}].$$

Niech

$$D_n(u) = \frac{1}{n} \sum_{i=1}^n 1(E_i \leq u), \quad u \in [0, 1]^d.$$

Przy H_0 dystrybuanta empiryczna D_n powinna być blisko C_{\perp} . Niech

$$S_n^{(C)} = n \int_{[0,1]^d} \{D_n(u) - C_{\perp}(u)\}^2 dD_n(u) = \sum_{i=1}^n \{D_n(E_i) - C_{\perp}(E_i)\}^2,$$

$$S_n^{(B)} = n \int_{[0,1]^d} \{D_n(u) - C_{\perp}(u)\}^2 du,$$

gdzie $C_{\perp}(u) = u_1 \times \cdots \times u_d$ - niezależna kopuła.

Dysponujemy $X = (X_1, \dots, X_n)$.

Testowanie z wykorzystaniem bootstrapu parametrycznego:

- 1 Liczymy statystykę S na podstawie próby X .
- 2 Estymujemy na podstawie X parametr θ rozkładu $P \in \mathcal{P}_0$.
- 3 Dla $i = 1, \dots, N$:
 - 1 Generujemy $Y^i = (Y_1, \dots, Y_m)$, gdzie $Y_i \sim P(\hat{\theta})$, $m > n$.
 - 2 Na podstawie Y^i liczymy statystykę S_b^i .
- 4 Liczymy *bootstrap p-value*:

$$p_b = \sum_{i=1}^N 1(S_b^i > S) / N.$$