

Andrzej Wilkowski

Uniwersytet Ekonomiczny we Wrocławiu

UWAGI O WSPÓŁCZYNNIKU KORELACJI

Streszczenie: W pracy omówiono wybrane własności klasycznego współczynnika korelacji i maksymalnego współczynnika korelacji. Wprowadzono także współczynnik zależności prostoliniowej oraz podano dowód asymptotycznej normalności jego próbkowego odpowiednika.

Słowa kluczowe: współczynnik korelacji, maksymalny współczynnik korelacji, współczynnik zależności prostoliniowej, asymptotyczna normalność.

Jednym z celów niniejszej pracy jest omówienie wybranych własności współczynnika korelacji (i jego estymatora), a także maksymalnego współczynnika korelacji (niektóre z nich zostały udowodnione w roku 2008, są zatem dosyć aktualne). W dalszej części artykułu podano współczynnik zależności prostoliniowej i jego związki ze współczynnikiem korelacji (są to wyniki własne, m.in. autora). Współczynnik ten jest teoretyczną podstawą konstrukcji innych miar zależności nieliniowej między zmiennymi losowymi.

Przypomnijmy, że **współczynnikiem korelacji liniowej r** zmiennych losowych X i Y nazywamy wielkość

$$r_{X,Y} = \frac{\text{Cov}X, Y}{\sqrt{\text{Var}X \text{Var}Y}}. \quad (1)$$

Oczywiście $-1 \leq r \leq 1$, $r_{X,Y} = r_{Y,X}$, $r_{mX+n,Y} = r_{mX+n, Y}$, o ile $m \neq 0$.

Przedstawia on ważną charakterystykę rozkładu wektora losowego (X, Y) . Główne jego własności są ściśle związane z dwiema prostymi regresji:

$$y - E(Y) = r(x - E(X)) + \sqrt{1-r^2}(Y - E(Y) - r(X - E(X))), \quad (2)$$

$$y - E(Y) = r(x - E(X)) + \sqrt{1-r^2}(Y - E(Y) - r(X - E(X))), \quad (3)$$

które są prostymi najlepszej zgodności, w sensie metody najmniejszych kwadratów, z masą prawdopodobieństwa w rozkładzie zmiennej (X, Y) [Cramer 1958]. Miarami zgodności tych prostych są następujące wyrażenia:

$$\min_a, b \in RE(Y - b - aX)^2 = \text{Var}Y_1 - r^2, \quad (4)$$

$$\min_a, b \in RE(X - b - aY)^2 = \text{Var}X_1 - r^2. \quad (5)$$

Widać z tego, że każda zmienna ma wariancję zmniejszoną w stosunku $(1 - r^2)$: 1, na skutek odjęcia od niej jej najlepszej, średniokwadratowej, liniowej oceny wyrażonej w zależności od drugiej zmiennej. Współczynnik r można zatem uważać za miarę stopnia liniowości wykazywanej przez rozkład wektora losowego (X, Y) . Stopień ten osiąga wartość największą, gdy $r = 1$, a cała masa prawdopodobieństwa jest rozpostarta na prostej. Przypadek przeciwny zachodzi, gdy $r = 0$, wtedy nie można zmniejszyć wariancji jakiegokolwiek zmiennej losowej przez odjęcie funkcji liniowej drugiej zmiennej.

Maksymalny współczynnik korelacji $R(X, Y)$ między zmiennymi losowymi X oraz Y został wprowadzony przez Gebeleina [Gebelein 1941]. Definiuje go wyrażenie:

$$R_{X,Y} = \sup_{f,g} \frac{\text{Cov}(f(X), g(Y))}{\sqrt{\text{Var}(f(X)) \text{Var}(g(Y))}} \quad (6)$$

gdzie supremum dotyczy wszystkich funkcji f, g , takich że $0 < \text{Var}(f(X), \text{Var}(g(Y))) < \infty$.

Wymieńmy kilka własności współczynnika R :

- jeśli wektor (X, Y) ma rozkład normalny, to $R(X, Y) = r(X, Y)$ [Lancaster 1957],
- jeśli niezdegenerowane zmienne losowe X_1, \dots, X_n są niezależne, o jednakowym rozkładzie, to

$$R_{S_m, S_n} = mn, \quad (7)$$

gdzie $m \leq n$ są naturalne oraz $S_k = \sum_{j=1}^k X_j$, $j = 1, \dots, n$ [Dembo, Kagan, Sheep 2001],

- jeśli niezdegenerowane zmienne losowe X_1, \dots, X_n są niezależne, o jednakowym rozkładzie, to

$$R(\sum_{j=1}^m X_j, \sum_{j=1}^n X_j) = m - \frac{m^2}{n}, \quad (8)$$

gdzie liczby naturalne l, m, n , spełniają warunek: $1 \leq l + 1 \leq m \leq n$ [Yaming 2008].

Warto wspomnieć o korelacyjnej nierówności Gaussa. Załóżmy, że A, B są wypukłymi, symetrycznymi podzbiorami przestrzeni R^n . Niech ν będzie gaussowską, centralną miarą na R^n . Wówczas **hipoteza korelacyjna Gaussa** stanowi, że:

$$\nu(A \cap B) \geq \nu(A)\nu(B). \quad (9)$$

Dowód tego znajduje się w pracy [He-Jing, Ze-Chun 2008].

Zauważmy, że jeśli mamy proste regresji zmiennych losowych X oraz Y :

$$y = a_1 x + b_1, \quad (10)$$

$$x = a_2 y + b_2, \quad (11)$$

możemy także wyznaczyć **współczynnik korelacji liniowej r** ; mianowicie:

$$r_{2X,Y} = a_1 a_2. \quad (12)$$

Zdefiniujemy **współczynnik zależności prostoliniowej** k zmiennych X, Y [Antoniewicz 1988]. Będziemy go rozumieli jako kosinus kąta, pod jakim przecinają się proste regresji. Po łatwych przekształceniach otrzymujemy:

$$k_{X,Y} = \cos \alpha = \frac{a_1 + a_2 a_{12}}{\sqrt{1 + a_2^2} \sqrt{1 + a_1^2}}, \quad (13)$$

gdzie α jest kątem przecięcia prostych regresji.

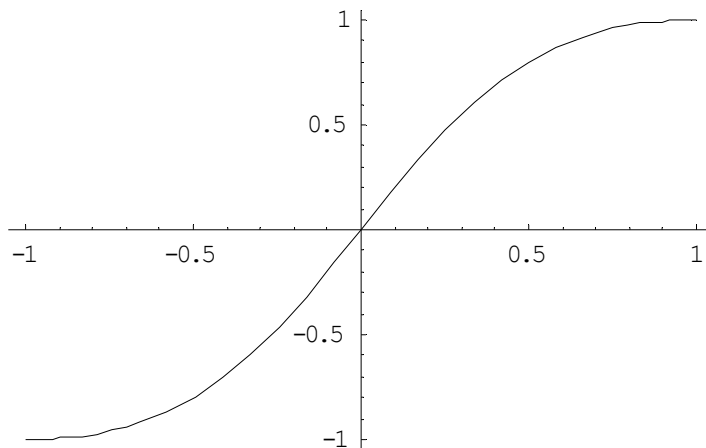
Możemy także napisać:

$$k_{X,Y} = \frac{r \sqrt{\text{Var} X \text{Var} Y}}{\sqrt{\text{Var} X + r^2 \text{Var} Y} \sqrt{\text{Var} Y + r^2 \text{Var} X}}. \quad (14)$$

Z tego wynika, że **współczynnik zależności prostoliniowej** k jest równy jeden, gdy między zmiennymi jest dokładna zależność liniowa, jeśli zaś $k = 0$, to takiej zależności nie ma. Oczywiście $k^2 = 1$ tylko wtedy, gdy $r^2 = 1$, oraz $k = 0$, gdy $r = 0$. Wartości pośrednie nie są jednak przyjmowane jednoznacznie. Może się zdarzyć, że przy ustalonej wielkości współczynnika r otrzymamy różne wartości współczynnika k (w zależności od wariancji). Rozpatrzmy teraz unormowane zmienne losowe (tzn. wariancja równa jeden, wartość oczekiwana równa zero). Wtedy wzór (14) przyjmie postać:

$$k = \frac{r}{\sqrt{1+r^2} \sqrt{1+r^2}}, \quad r \in [-1, 1]. \quad (15)$$

Rysunek 1 przedstawia wykres tej funkcji.



Rys. 1. Wykres funkcji $k(r)$

Źródło: opracowanie własne.

Można wyznaczyć maksymalną różnicę między współczynnikami k oraz r . Okazuje się, że:

$$\max_{r \in [-1, 1]} |k(r) - r| = \frac{1}{2}. \quad (16)$$

Maksimum jest osiągane dla $r = \pm 5 - 2$ oraz $Var(X) = Var(Y)$. Dowód znajduje się w pracy [Wilkowski 1994].

Zwróćmy jeszcze uwagę na fakt, że współczynnik korelacji liniowej r jest również kosinusem kąta, ale między innymi wektorami.

Jednym z ważniejszych rodzajów zbieżności według rozkładu jest zbieżność do rozkładu normalnego. Ciąg zmiennych losowych (X_n) zbiega według rozkładu do $N(m, s^2)$, $s > 0$, jeżeli równoważnie ciąg $((X_n - m)/s)$ zbiega według rozkładu do $N(0,1)$. Ogólniej mówimy, że **ciąg zmiennych losowych (X_n) jest asymptotycznie normalny o średniej m_n i wariancji sn^2** , jeżeli $sn^2 > 0$ dla dostatecznie dużych n oraz

$$X_n - mn \text{ sn } d N0,1. \quad (17)$$

Zapisujemy to jako: **X_n jest $AN(mn, sn^2)$** . Oczywiście ciągi (m_n) oraz (s_n) są ciągami stałych. Liczby te nie muszą być jednak średnią i odchyleniem standardowym zmiennej losowej X_n ; zmienna ta nie musi mieć ani średniej, ani odchylenia standardowego. Zauważmy, że jeżeli X_n jest $AN(mn, sn^2)$, to nie wynika z tego, że ciąg (X_n) w ogóle zbiega według rozkładu. Mamy jednak zawsze

$$\sup_t P(X_n \leq t) - P(N(mn, sn^2) \leq t) \rightarrow 0, n \rightarrow \infty. \quad (18)$$

Chcąc zatem obliczać prawdopodobieństwa, można traktować X_n jako zmienną losową $N(mn, sn^2)$ [Serfling 1999].

Niech $(X_1, Y_1), \dots, (X_n, Y_n)$ będą niezależnymi obserwacjami, o jednakowym rozkładzie, z pewnego rozkładu dwuwymiarowego (wektor (X_1, Y_1) ma taki sam rozkład jak wektor losowy (X, Y)). Jak pamiętamy, współczynnikiem korelacji liniowej zmiennych losowych X i Y jest wielkość

$$r_{X,Y} = \frac{Cov(X, Y)}{\sqrt{Var(X) Var(Y)}}. \quad (19)$$

Jego próbkowy odpowiednik ma postać:

$$r_n = \frac{1}{n} \sum_{i=1}^n (X_i - \bar{X})(Y_i - \bar{Y}) / \sqrt{\frac{1}{n} \sum_{i=1}^n (X_i - \bar{X})^2 \frac{1}{n} \sum_{i=1}^n (Y_i - \bar{Y})^2}, \quad (20)$$

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

Okazuje się, że r_n jest $AN(r, n^{-1} \mathbf{d} \mathbf{S} \mathbf{d}^T)$ [Serfling 1999], gdzie \mathbf{S} jest macierzą kowariancji wektora (X, Y, X^2, Y^2, XY) , a wektor $\mathbf{d} = (rE(X)Var(X) - E(Y)Var(X)Var(Y), rE(Y)Var(Y) - E(X)Var(X)Var(Y), -r2Var(X), -r2Var(Y), 1Var(X)Var(Y))$.

W rozważaniach dotyczących asymptotycznej normalności korzysta się z następującego twierdzenia.

Twierdzenie A. Przypuśćmy, że $X_n = X_{n1}, \dots, X_{nk}$ jest $AN(\mu, bn^2 \Sigma)$, gdzie Σ jest pewną macierzą kowariancji oraz $bn \rightarrow 0$. Niech $g(\mathbf{x}) = (g_1(\mathbf{x}), \dots, g_m(\mathbf{x}))$, $\mathbf{x} = (x_1, \dots, x_k)$, będzie taką funkcją wektorową, że każda jej składowa g_i jest funkcją rzeczywistą, a różniczka w punkcie $\mathbf{x} = \mu$ jest niezerowa (wystarczy, aby co najmniej jedna pochodna cząstkowa była ciągła i różna od zera w tym punkcie). Niech

$$\mathbf{D} = \partial g_i \partial x_j \text{ w } \mathbf{x} = \mu, \quad (21)$$

oczywiście \mathbf{D} jest macierzą wymiaru $m \times k$.

Wówczas

$$g(\mathbf{X}_n) \text{ jest } AN(g(\mu), bn2D\Sigma DT). \quad (22)$$

Dowód tego twierdzenia znajduje się w książce [Serfling 1999].

Poprzednio został wprowadzony współczynnik zależności prostoliniowej k zmiennych losowych X oraz Y rozumiany jako kosinus kąta, pod jakim przecinają się proste regresji tych zmiennych. W dalszym ciągu $(X_1, Y_1), \dots, (X_n, Y_n)$ będą niezależnymi obserwacjami o jednakowym rozkładzie z pewnego rozkładu dwuwymiarowego (wektor (X_i, Y_i) ma taki sam rozkład jak wektor losowy (X, Y)). Na podstawie wzorów (14) i (20) wnioskujemy, że próbkowy odpowiednik współczynnika k ma postać:

$$\begin{aligned} kn &= (i = 1n(X_i - X)^2 + i = 1n(Y_i - Y)^2) / (i = 1n(X_i - X)^2 + i = 1n(Y_i - Y)^2) \\ &= 1n(Y_i - Y)^2 / (i = 1n(X_i - X)^2 + i = 1n(Y_i - Y)^2). \end{aligned} \quad (23)$$

Twierdzenie B. Niech wektor $\mathbf{V} = \mathbf{X}, \mathbf{Y}$, $\mathbf{1n}i = \mathbf{1n}Xi^2$, $\mathbf{1n}i = \mathbf{1n}Yi^2$, $\mathbf{1n}i = \mathbf{1n}XiYi$, funkcja $g: R^5 \rightarrow R$ będzie określona wzorem

$$\begin{aligned} g(z_1, z_2, z_3, z_4, z_5) &= z_5 - z_1 z_2 (z_3 - z_1 z_2 z_4 - z_2 z_2 + z_4 - z_2 z_2 z_3 - z_1 z_2) z_3 - \\ & z_1 z_2 + z_5 - z_1 z_2 z_2 z_3 - z_1 z_2 z_4 - z_2 z_2 + z_5 - z_1 z_2 z_2 z_4 - z_2 z_2. \end{aligned} \quad (24)$$

Wówczas kn jest $AN(k, n^{-1}SST)$, gdzie \mathbf{S} jest macierzą kowariancji wektora (X, Y, X^2, Y^2, XY) , a wektor

$$\delta = \partial g / \partial z_1 z = EV, \dots, \partial g / \partial z_5 z = EV. \quad (25)$$

Dowód. Zauważmy, że wektor \mathbf{V} jest $AN(E(\mathbf{V}), n^{-1}\mathbf{S})$ [Serfling 1999]. Widać, że $kn = g(V)$, funkcja g spełnia założenia twierdzenia A (w naszym przypadku $b_n = n^{-1/2}$). Mamy zatem tezę twierdzenia B.

Na zakończenie zauważmy, że idee związane ze współczynnikiem zależności prostoliniowej mogą być wykorzystane przy konstrukcji innych miar zależności. Analogicznie jak w wypadku prostych regresji można definiować krzywe regresji (wprowadzając np. zaburzenie w części liniowej równania krzywej, raz przy jednej zmiennej, a raz przy drugiej). Następnie trzeba wyznaczyć punkt przecięcia krzywych regresji, który leży bliżej środka ciężkości zbioru danych. Kosinus kąta, pod jakim przecinają się krzywe w tym punkcie, będzie wtedy współczynnikiem zależności nieliniowej.

Literatura

- Antoniewicz R., *Metoda najmniejszych kwadratów dla zależności niejawnych i jej zastosowania w ekonomii*, Prace Naukowe Akademii Ekonomicznej we Wrocławiu nr 445, AE, Wrocław 1988.
 Brandt S., *Data Analysis. Statistical and Computational Methods for Statistics and Engineers*, ed. 3, Springer Verlag, New York 1999.

- Cieciura M., Zacharski J., *Metody probabilistyczne w ujęciu praktycznym*, Wydawnictwo Naukowe PWN, Warszawa 2007.
- Cramer H., *Metody matematyczne w statystyce*, PWN, Warszawa 1958.
- Dembo A., Kagan A., Sheep L.A., *Remarks on the maximum correlation coefficient*, „Bernoulli” 2001, 7.
- Gebelein H., *Das statistische Problem der Korrelation als Variations und Eigenwertproblem und sein Zusammenhang mit der Ausgleichsrechnung*, Z. Angew. Math. Mech. 1941, 21.
- He-Jing H., Ze-Chun H., *Gaussian Correlation Conjecture for Symmetric Convex Sets*, arXiv:0811.0488v1 [math.PR], 4 November 2008.
- Lancaster H.O., *Some properties of the bivariate normal distribution considered in the form of a contingency table*, „Biometrika” 44, s. 289-292.
- Renyi A., *On measures of dependence*, Acta Math. Hungar. 1959, 10.
- Serfling R.J., *Twierdzenia graniczne statystyki matematycznej*, PWN, Warszawa 1999.
- Yaming Y., *On the maximal correlation coefficient*, „Statistics and Probability Letters” 2008, 78.
- Wilkowski A., *Współczynnik zależności prostopolnołiniowej a współczynnik korelacji*, Prace Naukowe Akademii Ekonomicznej we Wrocławiu nr 667, AE, Wrocław 1994.

NOTE ON CORRELATION COEFFICIENT

Summary: The paper discusses correlation coefficient and the maximal coefficient of correlation. It also introduces the line dependent coefficient and gives the proof of the normal asymptotic distribution of its sample equivalent.