

Korelacja i regresja

Wojciech Kotłowski

Statystyka i analiza danych 2019/2020

5.05.2020

Miary populacyjne

- **Kowariancja:**

$$C(X, Y) = E[(X - \mu_X)(Y - \mu_Y)].$$

Mierzy **zależność liniową** dwóch zmiennych losowych.

Szczególny przypadek: wariancja $D^2[X] = C(X, X)$.

- **Korelacja** – unormowana kowariancja:

$$\rho(X, Y) = \frac{C(X, Y)}{D[X]D[Y]}, \quad \rho(X, Y) \in [-1, 1]$$

Jeśli X, Y – niezależne, to $\rho(X, Y) = 0$ (ale nie odwrotnie!)

Miary próbkowe

Dla zbioru n par $(X_1, Y_1), \dots, (X_n, Y_n)$:

- **Kowariancja:**

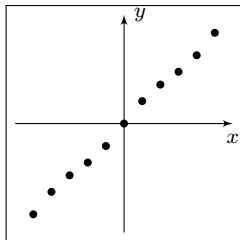
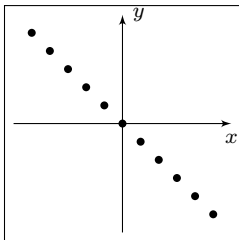
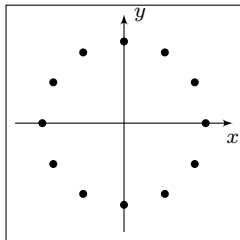
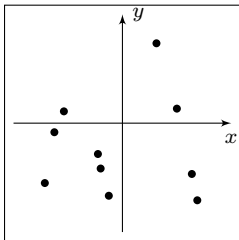
$$s_{XY} = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X})(Y_i - \bar{Y}).$$

- **Korelacja:**

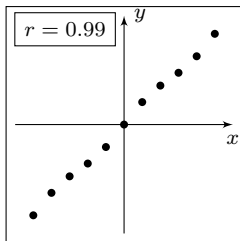
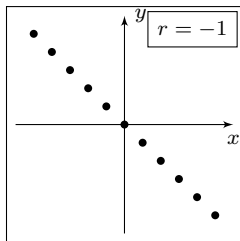
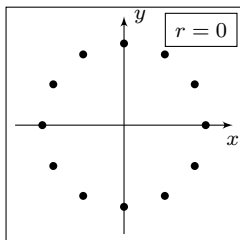
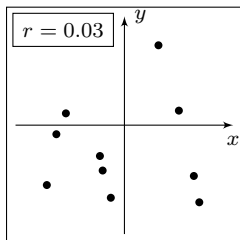
$$r = \frac{s_{XY}}{s_X s_Y} = \frac{\sum_{i=1}^n (X_i - \bar{X})(Y_i - \bar{Y})}{\sqrt{(\sum_{i=1}^n (X_i - \bar{X})^2)(\sum_{i=1}^n (Y_i - \bar{Y})^2)}}$$

Zachodzi $r \in [-1, 1]$, wartości skrajne $\{-1, 1\}$ przyjmowane są **wtedy i tylko wtedy** gdy Y_i są funkcją liniową X_i (lub odwrotnie)

Przykłady korelacji



Przykłady korelacji



Test na istotność korelacji

- **Układ hipotez:**

| | | | |
|---------|---------------|-----------------|-----------------|
| | | najczęściej | |
| $H_0 :$ | $\rho = 0$ | $(\rho \geq 0)$ | $(\rho \leq 0)$ |
| $H_1 :$ | $\rho \neq 0$ | $\rho < 0$ | $\rho > 0$ |

- **Statystyka testowa:**

$$T = \frac{r}{\sqrt{1-r^2}} \sqrt{n-2} \sim t(n-2)$$

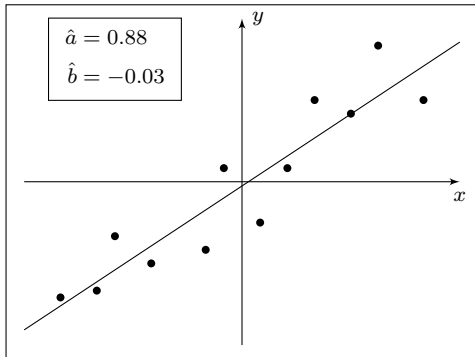
Wartość krytyczną (lub p -wartość) otrzymujemy z rozkładu t-Studenta z $n - 2$ stopniami swobody.

Regresja liniowa

Korelacja mierzy **siłę zależności liniowej**.

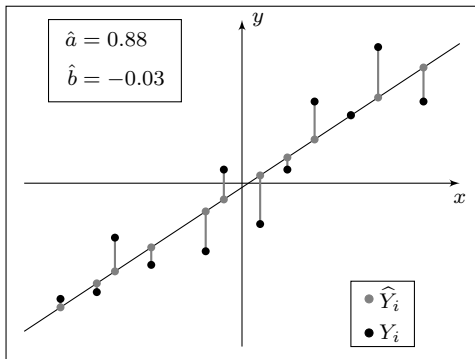
Regresja to wyznaczanie **współczynników zależności liniowej**:

Mając zbiór n punktów $(X_1, Y_1), \dots, (X_n, Y_n)$ wyznacz współczynniki \hat{a}, \hat{b} zależności liniowej $Y = \hat{a}X + \hat{b}$.



Metoda najmniejszych kwadratów

Dla każdego X_i błąd modelu liniowego to różnica między wartością odczytaną z prostej $\hat{Y}_i = \hat{a}X_i + \hat{b}$ a prawdziwą wartością Y_i :



Minimalizujemy **sumę kwadratów błędów**:

$$\hat{a}, \hat{b} \leftarrow \min_{a,b} \sum_{i=1}^n (Y_i - \hat{Y}_i)^2 = \sum_{i=1}^n (Y_i - aX_i - b)^2$$

Wyprowadzenie

$$L(a, b) = \sum_{i=1}^n (Y_i - aX_i - b)^2$$

Przyrównujemy pochodne cząstkowe do zera:

$$\frac{\partial L}{\partial b} = 0 \iff - \sum_{i=1}^n 2(Y_i - aX_i - b) = 0 \iff b = \bar{Y} - a\bar{X}$$

$$\frac{\partial L}{\partial a} = 0 \iff - \sum_{i=1}^n 2(Y_i - aX_i - b)X_i = 0$$

$$\iff \sum_{i=1}^n (Y_i - \bar{Y})X_i = a \sum_{i=1}^n (X_i - \bar{X})X_i$$

$$\iff \sum_{i=1}^n (Y_i - \bar{Y})(X_i - \bar{X}) = a \sum_{i=1}^n (X_i - \bar{X})(X_i - \bar{X})$$

$$\iff a = \frac{\sum_{i=1}^n (Y_i - \bar{Y})(X_i - \bar{X})}{\sum_{i=1}^n (X_i - \bar{X})^2} = \frac{s_{XY}}{s_X^2}.$$

Współczynniki regresji

$$\hat{a} = \frac{s_{XY}}{s_X^2} = r \frac{s_Y}{s_X}$$
$$\hat{b} = \bar{Y} - \hat{a}\bar{X}$$

- Linia regresji przechodzi przez punkt (\bar{X}, \bar{Y})
- Współczynnik kierunkowy \hat{a} ma **ten sam znak** co współczynnik korelacji r .