

9 Korelacja

Jeśli badaniu statystycznemu podlega jednocześnie wiele cech, to jednym z podstawowych zagadnień staje się analiza zależności pomiędzy nimi.

Do wykrycia zależności pomocne są odpowiednie wykresy, współczynniki mierzące jej siłę oraz testy badające jej istotność.

Metody opisowe:

1. Tabelaryczna - tablice wielodzielcze (korelacyjne).
2. Graficzna - diagramy korelacyjne (wykresy rozrzutu).
3. Statystyki opisowe – współczynniki zależności i korelacji.

Współczynniki zależności i korelacji

Niech $\mathbf{x} = (x_1, \dots, x_n)'$ i $\mathbf{y} = (y_1, \dots, y_n)'$ będą obserwacjami zmiennej XX oraz zmiennej YY , odpowiednio.

DEFINICJA

Współczynnikiem korelacji liniowej rr -Pearsona nazywamy liczbę:

$$r = \frac{\sum_{k=1}^n (x_k - \bar{x})(y_k - \bar{y})}{\sqrt{\sum_{k=1}^n (x_k - \bar{x})^2 \sum_{k=1}^n (y_k - \bar{y})^2}}, \quad -1 \leq r \leq 1.$$

Dwuwymiarowy model normalny

Zakładamy, że łącznym rozkładem zmiennych XX i YY jest dwuwymiarowy rozkład normalny.

FAKT

Zmienne XX i YY są nieskorelowane wtedy i tylko wtedy, gdy $\rho = 0$, zatem w tym modelu pojęcia korelacji i niezależności są równoważne.

FAKT

Estymatorem największej wiarygodności współczynnika korelacji ρ jest statystyka:

$$\hat{\rho} = r = \frac{\sum_{k=1}^n (X_k - \bar{X})(Y_k - \bar{Y})}{\sqrt{\sum_{k=1}^n (X_k - \bar{X})^2 \sum_{k=1}^n (Y_k - \bar{Y})^2}}.$$

Test istotności dla współczynnika korelacji

Założenia: Dwuwymiarowy model normalny.

Hipoteza zerowa: Zmienne XX i YY nie są istotnie (skorelowane) zależne.

$$H_0: \rho = 0$$

Hipotezy alternatywne:

$$H_1: \rho \neq 0$$

$$H_1: \rho > 0$$

$$H_1: \rho < 0$$

Statystyka testowa:

$$t = \frac{r}{\sqrt{1-r^2}} \sqrt{n-2}.$$

Rozkład statystyki testowej:

$$t|_{H_0} \sim t(n-2).$$

p -wymiarowy model normalny

Gęstość:

$$f(\mathbf{x}) = (2\pi)^{-\frac{p}{2}} |\boldsymbol{\Sigma}|^{-\frac{p}{2}} \exp\left[-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})' \boldsymbol{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu})\right], \mathbf{x} \in \mathbf{R}^p,$$

gdzie

$$\boldsymbol{\mu} = \begin{bmatrix} \mu_1 \\ \mu_2 \\ \vdots \\ \mu_p \end{bmatrix} - \text{wektor wartości oczekiwanych},$$

$$\Sigma = \begin{pmatrix} \sigma_{11} & \sigma_{12} & \cdots & \sigma_{1p} \\ \sigma_{21} & \sigma_{22} & \cdots & \sigma_{2p} \\ \vdots & \vdots & & \vdots \\ \sigma_{p1} & \sigma_{p2} & \cdots & \sigma_{pp} \end{pmatrix} - \text{macierz wariancji-kowariancji.}$$

Funkcje związane z analizą zależności i korelacji cech:

cor - estymator macierzy korelacji, estymator współczynnika korelacji liniowej *rr*-Pearsona,

cor.test - test istotności współczynnika korelacji liniowej.

Redukcja wymiaru

Analiza składowych głównych jest techniką redukcji wymiaru. Jej celem jest znalezienie niewielkiej liczby składowych głównych, które wyjaśniają w maksymalnym stopniu całkowitą wariancję z próby pp zmiennych pierwotnych X_1, \dots, X_p , tj. wielkość

$\sum_{j=1}^p \text{Var}(X_j) = \text{tr}(\Sigma)$ $\sum_{j=1}^p \text{Var}(X_j) = \text{tr}(\Sigma)$, gdzie Σ jest macierzą kowariancji wektora $\mathbf{X} = (X_1, \dots, X_p)'$ $\mathbf{X} = (X_1, \dots, X_p)'$.

Składowe główne są unormowanymi kombinacjami liniowymi zmiennych pierwotnych:

$$Z_1 = \mathbf{a}_1' \mathbf{X},$$

$$Z_2 = \mathbf{a}_2' \mathbf{X},$$

$$\vdots$$

$$Z_p = \mathbf{a}_p' \mathbf{X}.$$

Przekształcone zmienne (składowe główne) są ortogonalne i nieskorelowane.

Uwaga: Ponieważ macierz Σ nie jest znana, posługujemy się jej oszacowaniem z próby, tj. macierzą \mathbf{SS} .

Algorytm składowych głównych

1. Wyznaczamy współczynniki $\mathbf{a}_1 = (a_{11}, \dots, a_{1p})'$ $\mathbf{a}_1 = (a_{11}, \dots, a_{1p})'$ pierwszej składowej głównej, tak aby
 - zmaksymalizować wariancję zmiennej Z_1 : $\mathbf{a}_1' \mathbf{S} \mathbf{a}_1$,
 - $\mathbf{a}_1' \mathbf{a}_1 = 1$.
2. Wyznaczamy współczynniki $\mathbf{a}_2 = (a_{21}, \dots, a_{2p})'$ $\mathbf{a}_2 = (a_{21}, \dots, a_{2p})'$ drugiej składowej głównej, tak aby
 - zmaksymalizować wariancję zmiennej Z_2 : $\mathbf{a}_2' \mathbf{S} \mathbf{a}_2$,
 - $\mathbf{a}_2' \mathbf{a}_2 = 1$,

- składowa Z_2 była nieskorelowana z Z_1 : $\mathbf{a}_2' \mathbf{a}_1 = 0$.

3. Powtarzamy krok 2 (dla następnych składowych głównych), aż do otrzymania współczynników wszystkich p składowych głównych.

Uwaga: Wektor \mathbf{a}_i jest wektorem charakterystycznym, odpowiadającym i -tej co do wielkości, wartości własnej λ_i macierzy \mathbf{S} .

Własności składowych głównych

Mamy

$$\sum_{j=1}^p \text{Var}(X_j) = \sum_{j=1}^p \text{Var}(Z_j) = \sum_{j=1}^p \lambda_j = \text{tr}(\mathbf{S}).$$

W analizie składowych głównych oczekujemy, że dla pewnego małego k , suma $\lambda_1 + \lambda_2 + \dots + \lambda_k$ będzie bliska $\text{tr}(\mathbf{S}) = \lambda_1 + \lambda_2 + \dots + \lambda_p$. Jeśli tak jest, to k pierwszych składowych głównych wyjaśnia dobrze zmienność wektora $\mathbf{X} = (X_1, X_2, \dots, X_p)'$ i pozostałe $p - k$ składowe główne wnoszą niewiele, ponieważ mają one małe wariancje z próby.

Wskaźnik

$$\frac{\lambda_1 + \dots + \lambda_k}{\lambda_1 + \dots + \lambda_p} 100\%$$

jest procentową miarą wyjaśniania zmienności wektora \mathbf{X} przez pierwszych k składowych głównych.

Dobór liczby składowych głównych

1. Jeśli dla pewnego k wskaźnik

$$\frac{\lambda_1 + \dots + \lambda_k}{\lambda_1 + \dots + \lambda_p} 100\% \geq \beta,$$

np. $\beta = 80\%$, to pozostałe $p - k$ składowe główne pomijamy.

2. Pomijamy te składowe główne, których wartości własne są mniejsze od średniej

$$\bar{\lambda} = \frac{1}{p} \sum_{j=1}^p \lambda_j.$$

Uwaga: W ustaleniu liczby użytecznych składowych głównych, pomocny jest **wykres osypiska**.

Interpretacja składowych głównych

1. Wartość modułu współczynnika a_{ji} , w j -tej składowej głównej, pokazuje wkład w jej budowę i -tej zmiennej pierwotnej (z uwzględnieniem udziału pozostałych zmiennych pierwotnych).
2. Wartość współczynnika korelacji pomiędzy i -tą zmienną pierwotną, a j -tą składową główną, pokazuje wkład i -tej zmiennej pierwotnej w budowę j -tej składowej głównej (bez uwzględnienia udziału pozostałych zmiennych pierwotnych).

Funkcje związane z analizą składowych głównych:

princomp - analiza składowych głównych, procedura główna.