

Analiza variance

Enakost povprečij $r \geq 2$ normalno porazdeljenih populacij s podobnimi variancami lahko preverimo s testom, ki se imenuje analiza variance. Pri tem testu ničelna hipoteza trdi, da so povprečja vseh populacij med seboj enaka, alternativna pa, da se vsaj eno povprečje razlikuje od drugih:

$$\begin{aligned} H_0 : m_1 &= m_2 = \cdots = m_r, \\ H_1 : m_i &\neq m_j, \quad \text{za vsaj en par } (i, j). \end{aligned} \quad (1)$$

Za preverjanje ničelne hipoteze potrebujemo iz vsake od r populacij en vzorec z n_i elementi, skupaj torej r vzorcev z $n = \sum_{i=1}^r n_i$ elementi.

Test temelji na primerjavi odstopanja izmerjene spremenljivke med vzorci z odstopanjem spremenljivke znotraj vzorcev. Odstopanje med vzorci izrazimo s povprečnim kvadratičnim odstopanjem med vzorci S_1^2 , odstopanje znotraj vzorcev pa s povprečnim kvadratičnim odstopanjem znotraj vzorcev S_2^2 . V kolikor so odstopanja med vzorci značilno večja kot odstopanja znotraj vzorcev, H_0 zavrnemo. Uporabimo testno statistiko:

$$F = \frac{S_1^2}{S_2^2}. \quad (2)$$

Testna statistika F je porazdeljena s porazdelitvijo F , ki ima $(r-1, n-r)$ prostostnih stopenj. Interval zavračanja je $(f_{r-1, n-r; \alpha}, \infty)$. Količine, potrebne za izračun testne statistike, vpisujemo v tabelo:

Odstopanje	Vsota kvadratičnih odstopanj	Število prostostnih stopenj	Povprečje kvadratičnih odstopanj	Statistika
Med skupinami	q_1	$r-1$	$S_1^2 = \frac{q_1}{r-1}$	$F = \frac{S_1^2}{S_2^2}$
Znotraj skupin	q_2	$n-r$	$S_2^2 = \frac{q_2}{n-r}$	
Celotno	q	$n-1$		

Pri izračunu si pomagamo z naslednjimi formulami:

$$q = \left(\sum_{i=1}^r \sum_{j=1}^{n_i} x_{ij}^2 \right) - nm^2, \quad m = \frac{1}{n} \sum_{i=1}^r \sum_{j=1}^{n_i} x_{ij}, \quad (3)$$

$$q_1 = \left(\sum_{i=1}^r n_i m_i^2 \right) - nm^2, \quad m_i = \frac{1}{n_i} \sum_{j=1}^{n_i} x_{ij}, \quad (4)$$

$$q_2 = q - q_1, \quad n = \sum_{i=1}^r n_i. \quad (5)$$

Linearna regresija

S korelacijskim koeficientom r opišemo primernost uporabe linearne regresije za ocenjevanje medsebojne odvisnosti naključnih spremenljivk X in Y :

$$r = \frac{\text{Cov}[X, Y]}{\sqrt{\text{Var}[X] \text{Var}[Y]}}. \quad (6)$$

Zaloga vrednosti za r je $[-1, 1]$. Vrednosti $|r| \approx 1$ kažejo na izrazito linearno odvisnost med X in Y , medtem ko vrednosti $|r| < 0.5$ kažejo, da je regresijska premica neprimerna za upodobitev povezanosti X in Y , saj sta spremenljivki bodisi neodvisni ali pa je njuna odvisnost nelinearna. Linearno regresijo je smiselno uporabljati, če je $r \geq 0.75$. Za vzorec meritev $\{(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)\}$ izračunamo korelacijski koeficient po enačbi:

$$r = \frac{n \sum_{i=1}^n x_i y_i - \sum_{i=1}^n x_i \sum_{i=1}^n y_i}{\sqrt{\left(n \sum_{i=1}^n x_i^2 - \left(\sum_{i=1}^n x_i \right)^2 \right) \left(n \sum_{i=1}^n y_i^2 - \left(\sum_{i=1}^n y_i \right)^2 \right)}}. \quad (7)$$

Pri linearni regresiji vzamemo za matematični model odvisnosti linearno funkcijo:

$$Y = aX + b, \quad (8)$$

cilj pa je določiti koeficiente a in b tako, da bo vsota kvadratov napak najmanjša. Rešitev je:

$$a = \frac{\text{Cov}[X, Y]}{\text{Var}[X]}, \quad b = \text{E}[Y] - a\text{E}[X]. \quad (9)$$

Cenilki \hat{a} in \hat{b} za koeficiente a in b določimo na podlagi vzorca meritev:

$$\hat{a} = \frac{n \sum_{i=1}^n x_i y_i - \sum_{i=1}^n x_i \sum_{i=1}^n y_i}{n \sum_{i=1}^n x_i^2 - \left(\sum_{i=1}^n x_i \right)^2}, \quad \hat{b} = \frac{1}{n} \sum_{i=1}^n y_i - \hat{a} \frac{1}{n} \sum_{i=1}^n x_i. \quad (10)$$

Linearni model lahko uporabimo tudi, kadar zveza med spremenljivkama ni linearna. V takih primerih si lahko pomagamo z logaritmiranjem zvezne, v kolikor jo s tem lineariziramo, ali pa uvedemo nove spremenljivke:

$$\begin{aligned} Y = b e^{aX} &\rightarrow \ln Y = \ln b + aX, \\ Y = aX^2 + b &\rightarrow Y = aZ + b, \quad \text{kjer je } X^2 = Z. \end{aligned} \quad (11)$$