

S 145 (M 23a)

Toets voor de hypothese $P_1 = P_2 = \dots = P_k$ met

behulp van een $2 \times k$ tabel

Stichting Mathematisch Centrum

Statistische Afdeling

1954

Toets voor de hypothese $p_1 = p_2 = \dots = p_k$ met behulp van een $2 \times k$ tabel ¹⁾

Wij beschouwen k reeksen R_1, R_2, \dots, R_k van onafhankelijke waarnemingen, waarbij iedere waarneming als resultaat het kenmerk A of het kenmerk \bar{A} (non A) kan geven. De kans op A is binnen ieder der reeksen constant en wel gelijk aan p_i voor de waarnemingen van reeks R_i ($i = 1, 2, \dots, k$).

Laat het aantal waarnemingen van reeks R_i ($i = 1, 2, \dots, k$) gelijk zijn aan n_i en laat hieronder het aantal met kenmerk A , m_i zijn. Gevraagd wordt dan de hypothese $H_0: p_1 = p_2 = \dots = p_k$ te toetsen op grond van deze gegevens.

De gegevens kunnen in een $2 \times k$ -tabel worden samengevat;

	R_1	R_2	...	R_i	...	R_k	totaal
A	m_1	m_2	...	m_i	...	m_k	m
\bar{A}	$n_1 - m_1$	$n_2 - m_2$...	$n_i - m_i$...	$n_k - m_k$	$n - m$
totaal	n_1	n_2		n_i		n_k	n

waarin dus $m_1 + m_2 + \dots + m_k = m$

en $n_1 + n_2 + \dots + n_k = n$.

De hypothese H_0 wordt getoetst met de grootheid

$$\chi_c^2 = \sum_i \frac{(m_i - \frac{m \cdot n_i}{n})^2}{\frac{m \cdot n_i}{n}} + \sum_i \frac{(n_i - m_i - \frac{(n-m)n_i}{n})^2}{\frac{(n-m)n_i}{n}}$$

$$= \sum_i \frac{(nm_i - mn_i)^2}{m(n-m)n_i} = \frac{n^2}{m(n-m)} \sum_i \frac{m_i^2}{n_i} - \frac{nm}{n-m}$$

Deze grootheid χ_c^2 ²⁾ heeft onder de hypothese H_0 bij benadering een χ^2 -verdeling met $k-1$ vrijheidsgraden (zie b.v. [1] p. 445 e.v.).

1) Dit memorandum is slechts bedoeld ter oriëntatie en streeft niet naar volledigheid of volledige exactheid.

2) Als wij grootheden als stochastische grootheden (dit zijn grootheden met een waarschijnlijkheidsverdeling) beschouwen geven wij dit door onderstreping aan. Niet onderstreepte letters geven waarden aan, die door de stochastische grootheden worden aangenomen.

Deze benadering is goed, indien $m \frac{n_i}{n} \geq 5$ voor iedere i (zie [2]). Indien H_0 onjuist is, dus als er bij verschillende reeksen verschillende kansen op A zijn, zal $\underline{\chi}_c^2$ gewoonlijk grotere waarden aannemen, dan wanneer H_0 juist is.

De kritieke zone bestaat uit die waarden van $\underline{\chi}_c^2$, waarvoor geldt $\underline{\chi}_c^2 \geq \chi_\alpha^2$. Hierin is χ_α^2 die waarde van $\underline{\chi}^2$, die voldoet aan

$$P[\underline{\chi}^2 \geq \chi_\alpha^2] = \alpha$$

met α als van te voren vastgelegde onbetrouwbaarheid.

De overschrijdingskans behorende bij een bepaalde gevonden waarde $\underline{\chi}_c^2$ van $\underline{\chi}^2$ is gedefinieerd als

$$P[\underline{\chi}_c^2 \geq \chi_c^2 | H_0]$$

waarin " $|H_0$ " aangeeft, dat deze kans berekend wordt op grond van H_0 . χ_α^2 en de overschrijdingskans kunnen in tabellen of nomogrammen worden opgezicht (zie [3]).

Opmerking. Indien niet voldaan is aan de voorwaarde $m \frac{n_i}{n} \geq 5$ voor iedere i , kan men een (meer bewerkelijke) exacte toets baseren op de voorwaardelijke waarschijnlijkheidsverdeling van de grootheden \underline{m}_i ($i = 1, \dots, k$), onder de voorwaarde, dat hun som de waarde m aanneemt:

$$P[\underline{m}_1 = m_1, \underline{m}_2 = m_2, \dots, \underline{m}_k = m_k | \underline{m}_1 + \underline{m}_2 + \dots + \underline{m}_k = m; H_0] = \\ = \frac{\binom{n_1}{m_1} \binom{n_2}{m_2} \dots \binom{n_k}{m_k}}{\binom{n}{m}}$$

De geldigheid van deze formule volgt direct uit de waarschijnlijkheidsverdelingen van de \underline{m}_i en van \underline{m} (onder H_0) en uit de definitie van een voorwaardelijke waarschijnlijkheid.

In dit geval definiëren wij de overschrijdingskans behorend bij een gevonden resultaat (m_1, m_2, \dots, m_k) met $m_1 + m_2 + \dots + m_k = m$ als de som van alle waarschijnlijkheden van bovengenoemde verdeling (met de gevonden waarde van m), die hoogstens gelijk zijn aan de waarschijnlijkheid van het gevonden resultaat.

Literatuur.

- [1] H.Cramér, Mathematical methods of statistics, Princeton University Press, 1946.
- [2] P.G.Hoel, On indices of dispersion, Ann. Math. Stat. 14 (1943), p. 155-163.

- [3] Tabellen en nomogrammen van de -verdeling.
M.G.Kendall, The advanced theory of statistics, I, 1947,
p. 444-446.
H.Cramér, Mathematical methods of statistics, Princeton
University Press, 1946, p. 559.
Statistica 1 (1946), p. 109.