

"(k-1)-mean significance levels" van de asymptotische versie van de methode voor simultane uitspraken voor het k-steekproevenprobleem gebaseerd op de toets van Kruskal en Wallis

Citation for published version (APA):

Oude Voshaar, J. H. (1976). "(k-1)-mean significance levels" van de asymptotische versie van de methode voor simultane uitspraken voor het k-steekproevenprobleem gebaseerd op de toets van Kruskal en Wallis. (Memorandum COSOR; Vol. 7627). Technische Hogeschool Eindhoven.

Document status and date:

Gepubliceerd: 01/01/1976

Document Version:

Uitgevers PDF, ook bekend als Version of Record

Please check the document version of this publication:

- A submitted manuscript is the version of the article upon submission and before peer-review. There can be important differences between the submitted version and the official published version of record. People interested in the research are advised to contact the author for the final version of the publication, or visit the DOI to the publisher's website.
- The final author version and the galley proof are versions of the publication after peer review.
- The final published version features the final layout of the paper including the volume, issue and page numbers.

[Link to publication](#)

General rights

Copyright and moral rights for the publications made accessible in the public portal are retained by the authors and/or other copyright owners and it is a condition of accessing publications that users recognise and abide by the legal requirements associated with these rights.

- Users may download and print one copy of any publication from the public portal for the purpose of private study or research.
- You may not further distribute the material or use it for any profit-making activity or commercial gain
- You may freely distribute the URL identifying the publication in the public portal.

If the publication is distributed under the terms of Article 25fa of the Dutch Copyright Act, indicated by the "Taverne" license above, please follow below link for the End User Agreement:

www.tue.nl/taverne

Take down policy

If you believe that this document breaches copyright please contact us at:

openaccess@tue.nl

providing details and we will investigate your claim.

TECHNISCHE HOGESCHOOL EINDHOVEN

Onderafdeling der Wiskunde

SECTIE KANSREKENING, STATISTIEK EN OPERATIONS RESEARCH

Memorandum COSOR 76-27

"(k-1)-mean significance levels" van de
asymptotische versie van de methode voor
simultane uitspraken voor het k-steekproeven-
probleem gebaseerd op de toets van Kruskal en
Wallis

door

J.H. Oude Voshaar

Eindhoven, December 1976

Nederland

"(k-1)-mean significance levels" van de asymptotische versie van de methode voor simultane uitspraken voor het k-steekproevenprobleem gebaseerd op de toets van Kruskal en Wallis

door

J.H. Oude Voshaar

1. Inleiding en samenvatting

Laat $y_{11}, \dots, y_{1n_1}; \dots; y_{k1}, \dots, y_{kn_k}$ onafhankelijke stochastische grootheden zijn, waarbij y_{ij} een continue verdelingsfunctie F_i heeft ($k \geq 3$).

Een verdelingsvrije toets voor de nulhypothese $H_0: F_1 = F_2 = \dots = F_k$ is de toets van Kruskal en Wallis met als toetingsgrootheid:

$$\underline{H} := \frac{12}{N(N+1)} \sum_{i=1}^k n_i (\bar{R}_{i\cdot} - \bar{R}_{\cdot\cdot})^2 = \frac{12}{N(N+1)} \sum_{i=1}^k n_i \bar{R}_{i\cdot}^2 - 3N(N+1)$$

waarbij $N := \sum_{i=1}^k n_i$ en R_{ij} het rangnummer van y_{ij} onder alle waarnemingen is. Verder:

$$\bar{R}_{i\cdot} := \frac{1}{n_i} \sum_{j=1}^{n_i} R_{ij} \quad \text{en} \quad \bar{R}_{\cdot\cdot} := \frac{1}{N} \sum_{i,j} R_{ij} = \frac{N+1}{2}$$

Bij de onbetrouwbaarheidsdrempel α verwerpen we H_0 als $H > h_\alpha$. Als H_0 verworpen wordt weten we hiermee echter nog niet, voor welke paren (i,j) F_i en F_j verschillend zijn. Maar gebruik makend van het

lemma:

als $c > 0$, dan geldt:

$$\sum_{i=1}^k y_i^2 \leq c^2 \quad \text{d.e.s.d.} \quad \left| \sum_{i=1}^k a_i y_i \right| \leq c \left(\sum_{i=1}^k a_i^2 \right)^{\frac{1}{2}} \quad \text{voor alle } (a_1, \dots, a_k) \in \mathbb{R}^k$$

vinden we dat uit $P_{H_0} \{ \underline{H} \leq h_\alpha \} \geq 1-\alpha$ volgt:

$$(1.1) \quad P_{H_0} \left\{ \left| \bar{R}_{i\cdot} - \bar{R}_{j\cdot} \right| \leq \sqrt{h_\alpha \frac{N(N+1)}{12} \left(\frac{1}{n_i} + \frac{1}{n_j} \right)} \text{ voor alle } i \text{ en } j \right\} \geq 1-\alpha$$

We konkluderen dus dat $F_i \neq F_j$ als geldt:

$$|\bar{R}_{i.} - \bar{R}_{j.}| > \sqrt{h_\alpha \frac{N(N+1)}{12} \left(\frac{1}{n_i} + \frac{1}{n_j} \right)}$$

en deze methode heeft (onder H_0) een simultane onbetrouwbaarheid kleiner of gelijk aan α .

Opmerking: Deze methode is afkomstig van Nemenyi (1963).

Het nadeel van deze methode is echter, dat de simultane onbetrouwbaarheid, vooral voor grotere waarden van k , veel kleiner is dan α (zie tabel 4.1), omdat in (1.1) niet alle lineaire combinaties uit het lemma gebruikt worden.

In het geval van steekproeven van gelijke omvang, d.w.z. $n_1 = n_2 = \dots = n_k = n$, is, voor grotere n , de methode te verbeteren tot:

$$(1.2) \quad \text{konkludeer } F_i \neq F_j \text{ als } |\bar{R}_{i.} - \bar{R}_{j.}| > q_{k,\infty}^\alpha \sqrt{\frac{k(k+1)}{12}}$$

waarin $q_{k,\infty}^\alpha$ de rechterkritieke waarde is van de verdeling van de range van k standaard-normaal verdeelde grootheden (een grootheid met deze verdeling noteren we wel eens als $q_{k,\infty}$).

Deze methode (1.2) heeft voor $n \rightarrow \infty$ een simultane onbetrouwbaarheid die exakt gelijk is aan α (zie 4.2).

We zijn echter niet alleen in deze onbetrouwbaarheid onder H_0 geïnteresseerd: stel dat van de verdelingsfuncties F_1, \dots, F_k er p gelijk zijn en de overige $k-p$ verschillen. Dan gaat onze belangstelling nu vooral uit naar de kans dat uit die p identieke verdelingsfuncties toch sommige verschillend verklaard worden.

Miller (1966) noemt deze kans, gemaximaliseerd over alle mogelijke verdelingsfuncties F_1, \dots, F_k waarvan er p gelijk zijn, de "p-mean significance level α_p ".

Omdat $\bar{R}_{i.}$ diskreet verdeeld is, is het te moeilijk om α_p te bepalen voor de methode gebaseerd op formule (1.1).

Voor methode (1.2) is het wel mogelijk α_p te bepalen als $n \rightarrow \infty$. We zullen ons hier beperken tot het geval $p = k-1$, d.w.z. er zijn continue verdelingsfuncties F en G zodat geldt: $F_1 = \dots = F_{k-1} = F$ en $F_k = G$ en we vragen ons af of hierbij (voor $n \rightarrow \infty$) geldt:

$$\alpha_{k-1} \leq \alpha$$

waarin

$$\alpha_{k-1} := \sup_{F,G \text{ kontinu}} P \left\{ \max_{1 \leq i, j \leq k-1} |\bar{R}_{i.} - \bar{R}_{j.}| > q_{k,\infty}^\alpha \sqrt{\frac{k(kn+1)}{12}} \right\}$$

We berekenen hiertoe eerst (in hoofdstuk 2) de verwachting en variantie van $\bar{R}_{i.}$ en $\bar{R}_{j.}$, zodat we, gebruik makend van de in hoofdstuk 3 te bewijzen asymptotische normaliteit van de vektor $(\bar{R}_{1.}, \dots, \bar{R}_{k-1.})$, de asymptotische verdeling van $\max_{1 \leq i, j \leq k-1} |\bar{R}_{i.} - \bar{R}_{j.}|$ kennen (hoofdstuk 4).

Als we, behalve continuïteit, geen andere beperkingen aan F en G opleggen, dan zal blijken dat voor de gangbare waarden van α geldt: $\alpha_{k-1} > \alpha$ (zie hoofdstuk 5). Beperken we ons echter tot de klasse van Lehmann-alternatieven ($G = F^t$ voor zekere $t \in (0, \infty)$), dan geldt wel: $\alpha_{k-1} < \alpha$ (hoofdstuk 6).

Voor verschuivingsalternatieven zullen we in hoofdstuk 7 vinden:

Als F symmetrisch en unimodaal (= ééntoppig) is, dan geldt: $\alpha_{k-1} < \alpha$.

2. Verwachting, variantie en kovariantie van $\bar{R}_{1.}, \dots, \bar{R}_{k-1.}$

We veronderstellen dus dat $F_1 = F_2 = \dots = F_{k-1} = F$ en $F_k = G$ (F en G continu) en verder $n_1 = \dots = n_k = n$.

Dan zijn $\bar{R}_{1.}, \dots, \bar{R}_{k-1.}$ identiek verdeeld.

We definiëren de functie u door:

$$u(z) := \begin{cases} 0 & \text{als } z < 0 \\ 1 & \text{als } z \geq 0 \end{cases}$$

dan

$$\bar{R}_{1.} := \sum_{i=1}^n R_{1i} = \sum_{i=1}^n \sum_{j=2}^k \sum_{\ell=1}^n u(y_{1i} - y_{j\ell}) + \frac{1}{2} n(n+1)$$

Nu geldt dat:

$$E[u(y_{1i} - y_{j\ell})] = \begin{cases} \frac{1}{2} & \text{als } j = 2, \dots, k-1 \\ p & \text{als } j = k \end{cases}$$

$$(2.1) \text{ waarin } p := P(y_{k1} \leq y_{11}) = \int G(y) dF(y)$$

zodat:

$$E(\bar{R}_{1.}) = \frac{1}{2}n \sum_{j=2}^{k-1} n + pn^2 + \frac{1}{2}n(n+1)$$

en dus
$$\mathcal{E}(\bar{R}_1) = \frac{1}{2}(kn+1) + (p - \frac{1}{2})n$$

$$(\quad = \frac{1}{2}(N+1) + (p - \frac{1}{2})n)$$

Voor de berekening van $\text{var}(\bar{R}_1)$ vatten we de waarnemingen op als 3 steekproeven:

$$Y_{11}, \dots, Y_{1n}; Y_{21}, \dots, Y_{2, (k-2)n}; Y_{k1}, \dots, Y_{kn}$$

met verdelingsfuncties F, F en G.

Dan geldt nu:

$$\bar{R}_1 = \sum_{i=1}^n \sum_{j=1}^{(k-2)n} u(Y_{1i} - Y_{2j}) + \sum_{i=1}^n \sum_{\ell=1}^n (Y_{1i} - Y_{k\ell}) + \frac{1}{2}n(n+1)$$

Vanwege de onafhankelijkheid hebben we:

als $i \neq i'$ en $j \neq j'$ dan $\mathcal{E}[u(Y_{1i} - Y_{2j})u(Y_{1i'} - Y_{2j'})] = \frac{1}{4}$

als $i \neq i'$ en $\ell \neq \ell'$ dan $\mathcal{E}[u(Y_{1i} - Y_{k\ell})u(Y_{1i'} - Y_{k\ell'})] = p^2$

en als $i \neq i'$ dan $\mathcal{E}[u(Y_{1i} - Y_{2j})u(Y_{1i'} - Y_{k\ell})] = \frac{1}{2}p$

Verder geldt (omdat $u^2(z) = u(z)$ voor alle z):

$$\mathcal{E}[u^2(Y_{1i} - Y_{2j})] = \frac{1}{2}$$

$$\text{en } \mathcal{E}[u^2(Y_{1i} - Y_{k\ell})] = p$$

Als $j \neq j'$ dan

$$\mathcal{E}[u(Y_{1i} - Y_{2j})u(Y_{1i} - Y_{2j'})] = P(Y_{2j} \leq Y_{1i} \text{ en } Y_{2j'} \leq Y_{1i}) =$$

$$= \int F^2(y) dF(y) = \frac{1}{3} \quad (\text{omdat } F \text{ continu is})$$

Als $i \neq i'$ dan

$$\mathcal{E}[u(Y_{1i} - Y_{2j})u(Y_{1i'} - Y_{2j})] = \int (1-F(y))^2 dF(y) = \frac{1}{3}$$

als $l \neq l'$ dan

$$\mathcal{E} [u(y_{1i} - y_{kl}) u(y_{1i} - y_{kl}')] = \int G^2(y) dF(y) = q$$

als $i \neq i'$ dan

$$\begin{aligned} \mathcal{E} [u(y_{1i} - y_{kl}) u(y_{1i'} - y_{kl}')] &= \int (1-F(y))^2 dG(y) \\ &= \left[(1-F(y))^2 G(y) \right]_{y=-\infty}^{y=\infty} + 2 \int (1-F(y)) G(y) dF(y) = \\ &= 2 \int G(y) dF(y) - 2 \int F(y) G(y) dF(y) = 2(p-r) \end{aligned}$$

en tot slot geldt dat:

$$\mathcal{E} [u(y_{1i} - y_{2j}) u(y_{1i} - y_{kl}')] = \int F(y) G(y) dF(y) = r$$

Hierbij zijn q en r gedefinieerd door:

$$(2.2) \quad \begin{aligned} q &= \int G^2(y) dF(y) \\ r &= \int F(y) G(y) dF(y) \end{aligned}$$

We vinden dan:

$$\begin{aligned} \mathcal{E} (R_{1.}^2) &= \sum_{i=1}^n \sum_{j=1}^{(k-2)n} \sum_{i'=1}^n \sum_{j'=1}^{(k-2)n} \mathcal{E} [u(y_{1i} - y_{2j}) u(y_{1i'} - y_{2j}')] + \\ &+ \sum_{i=1}^n \sum_{l=1}^n \sum_{i'=1}^n \sum_{l'=1}^n \mathcal{E} [u(y_{1i} - y_{kl}) u(y_{1i'} - y_{kl}')] + \\ &+ 2 \sum_{i=1}^n \sum_{j=1}^{(k-2)n} \sum_{i'=1}^n \sum_{l=1}^n \mathcal{E} [u(y_{1i} - y_{2j}) u(y_{1i'} - y_{kl}')] + \\ &+ 2 \cdot \frac{1}{2} n(n+1) \mathcal{E} (R_{1.}) - \frac{1}{4} n^2 (n+1)^2 = \\ &= \frac{1}{4} n(n-1)(k-2)n \{ (k-2)n-1 \} + \frac{1}{3} n(k-2)n \{ (k-2)n-1 \} + \\ &+ \frac{1}{3} n(n-1)(k-2)n + \frac{1}{2} n(k-2)n + p^2 n(n-1)n(n-1) + \\ &+ qn^2(n-1) + 2(p-r)n(n-1)n + pn^2 + 2 \cdot \frac{1}{2} pn(n-1)(k-2)n^2 + \\ &+ 2 rn(k-2)n^2 + n(n+1) \left\{ \frac{1}{2} n(kn+1) + (p-\frac{1}{2})n^2 \right\} - \frac{1}{4} n^2 (n+1)^2 \end{aligned}$$

zodat:

$$\begin{aligned}
 (2.3) \quad \text{var}(\bar{R}_1) &= \frac{1}{n} \mathcal{E}(R_1^2) - (\mathcal{E}\bar{R}_1)^2 = \\
 &= \frac{1}{12} k^2 n + (2r - p - \frac{1}{4})kn + (4p - 2p^2 + q - 6r + \frac{1}{6})n + \\
 &+ \frac{1}{12} k - p + p^2 - q + 2r - \frac{1}{6}
 \end{aligned}$$

Teneinde de kovariantie van \bar{R}_1 en \bar{R}_2 te berekenen vatten we de waarnemingen op als 4 steekproeven:

$$Y_{11}, \dots, Y_{1n}; Y_{21}, \dots, Y_{2n}; Y_{31}, \dots, Y_{3(k-3)n}; Y_{k1}, \dots, Y_{kn}$$

met verdelingsfuncties F, F, F en G.

Dan geldt:

$$\begin{aligned}
 \bar{R}_1 &= \sum_{i=1}^n \sum_{j=1}^n u(Y_{1i} - Y_{2j}) + \sum_{i=1}^n \sum_{\ell=1}^{(k-3)n} u(Y_{1i} - Y_{3\ell}) \\
 &+ \sum_{i=1}^n \sum_{m=1}^n u(Y_{1i} - Y_{km}) + \frac{1}{2}n(n+1) \\
 \text{en} \quad \bar{R}_2 &= \sum_{j=1}^n \sum_{i=1}^n u(Y_{2j} - Y_{1i}) + \sum_{j=1}^n \sum_{\ell=1}^{(k-3)n} u(Y_{2j} - Y_{3\ell}) \\
 &+ \sum_{j=1}^n \sum_{m=1}^n u(Y_{2j} - Y_{km}) + \frac{1}{2}n(n+1)
 \end{aligned}$$

Door \bar{R}_1, \bar{R}_2 op dezelfde manier als bij de berekening van $\mathcal{E}(R_1^2)$ term voor term uit te vermenigvuldigen en de verwachting te nemen, vinden we (hierbij gebruik makend van de definities van p, q en r):

$$\begin{aligned}
 (2.4) \quad \text{cov}(\bar{R}_1, \bar{R}_2) &= \frac{1}{n} \mathcal{E}(R_1 R_2) - (\mathcal{E}\bar{R}_1)(\mathcal{E}\bar{R}_2) = \\
 &= -\frac{1}{12}kn + (3p - p^2 - 4r + \frac{1}{12})n - \frac{1}{12}
 \end{aligned}$$

Opmerking 1:

Onder H_0 (d.w.z. $F = G$) geldt:

$$p = \frac{1}{2} \text{ en } q = r = \frac{1}{3}$$

zodat dan

$$\text{var}(\bar{R}_{1.}) = \frac{1}{12}k(k-1)n + \frac{1}{12}(k-1)$$

$$\text{en cov}(\bar{R}_{1.}, \bar{R}_{2.}) = -\frac{1}{12}(kn+1)$$

hetgeen overeenkomt met de bekende uitdrukkingen hiervoor (zie bijvoorbeeld Miller, pag. 171).

Opmerking 2:

Voor n_1, \dots, n_k niet alle gelijk vindt men op dezelfde wijze:

$$\begin{aligned} \text{var } \bar{R}_{1.} &= \frac{1}{n_1} \left\{ \frac{1}{12}M^2 + (q-p^2)n_k^2 + (2r-p)Mn_k + \frac{1}{12}M + (2r-p-q+p^2)n_k \right\} \\ &+ \frac{1}{12}M + (2p-p^2-m)n_k \end{aligned}$$

$$\text{waarin } M := \sum_{i=2}^{k-1} n_i = N - n_1 - n_k$$

$$\text{en cov}(\bar{R}_{1.}, \bar{R}_{2.}) = -\frac{1}{12}(N+1) + (3p - p^2 - 4r + \frac{1}{12})n_k$$

3. Asymptotische normaliteit van $(\bar{R}_{1.}, \bar{R}_{2.}, \dots, \bar{R}_{k-1.})$

In hoofdstuk 4 zullen we nodig hebben, dat $(\bar{R}_{1.}, \dots, \bar{R}_{k-1.})$ asymptotisch multi-normaal verdeeld is voor $n \rightarrow \infty$. Dit is d.e.s.d. het geval als

elke lineaire combinatie $\sum_{i=1}^{k-1} c_i \bar{R}_{i.}$ asymptotisch normaal verdeeld is

en dit laatste bewijzen we met stelling 2.1 uit Hajek.(1968)..

De voorwaarden van deze stelling zijn hier vervuld d.e.s.d. als

$$\text{var} \left(\sum_{i=1}^{k-1} c_i \bar{R}_{i.} \right) \text{ naar oneindig gaat voor } n \rightarrow \infty.$$

Voordat we deze voorwaarde verifiëren formuleren we echter eerst de volgende lemma's:

lemma 3.1:

$$\text{cov}(\bar{R}_{1.}, \bar{R}_{2.}) \geq -\frac{1}{k-2} \text{var}(\bar{R}_{1.})$$

en het gelijkteken geldt d.e.s.d. als $\text{var} \sum_{i=1}^{k-1} \bar{R}_{i.} = 0$

Bewijs: Het gestelde volgt uit:

$$\text{var}\left(\sum_{i=1}^{k-1} \bar{R}_{i.}\right) = (k-1) \{ \text{var}(\bar{R}_{i.}) + (k-2) \text{cov}(\bar{R}_{i.}, \bar{R}_{2.}) \} \geq 0$$

Q.E.D.

lemma 3.2:

$$\text{cov}(\bar{R}_{1.}, \bar{R}_{2.}) < 0$$

Bewijs: Onder de voorwaarde $\bar{R}_{.k} = r_k$ (d.w.z. $\bar{R}_{k1} = r_{k1}, \dots, \bar{R}_{kn} = r_{kn}$) geldt dat $\sum_{i=1}^{k-1} \bar{R}_{i.} = \text{konstant}$, dus voor $\text{var}\left(\sum_{i=1}^{k-1} \bar{R}_{i.} \mid \bar{R}_{.k} = r_k\right) = 0$

lemma 3.1 geeft:

$$\text{cov}(\bar{R}_{1.}, \bar{R}_{2.} \mid \bar{R}_{.k} = r_k) = -\frac{1}{k-2} \text{var}(\bar{R}_{1.} \mid \bar{R}_{.k} = r_k) < 0$$

waaruit volgt: $\text{cov}(\bar{R}_{1.}, \bar{R}_{2.}) = \sum_{r_k} \text{cov}(\bar{R}_{1.}, \bar{R}_{2.} \mid \bar{R}_{.k} = r_k) \cdot P(\bar{R}_{.k} = r_k) < 0$

Q.E.D.

lemma 3.3

$$\sum_{i \neq j} c_i c_j \leq (k-2) \sum_{i=1}^{k-1} c_i^2 \text{ en het gelijkteken}$$

geldt d.e.s.d. als $c_1 = c_2 = \dots = c_{k-1}$

Bewijs: De bewering volgt uit de ongelijkheid van Cauchy-Schwarz:

$$(k-1) \sum_{i=1}^{k-1} c_i^2 \geq \left(\sum_{i=1}^{k-1} c_i\right)^2 = \sum_{i \neq j} c_i c_j + \sum_{i=1}^{k-1} c_i^2$$

Q.E.D.

Als we definiëren:

$$(3.1) \quad a_1 := \frac{1}{12}k^2 + (2r - p - \frac{1}{4})k + 4p - 2p^2 + q - 6r + \frac{1}{6}$$

$$a_2 := \frac{1}{12}k - p + p^2 - q + 2r - \frac{1}{6}$$

$$a_3 := -\frac{1}{12}k + 3p - p^2 - 4r + \frac{1}{12}$$

zodat: $\text{var} \bar{R}_{1.} = a_1 n + a_2$

$$\text{en cov}(\bar{R}_{1.}, \bar{R}_{2.}) = a_3 n - \frac{1}{12}$$

dan geldt:

$$\begin{aligned} \text{var} \sum_{i=1}^{k-1} (c_i \bar{R}_{-i}) &= \sum_{i=1}^{k-1} c_i^2 (a_1 n + a_2) + \sum_{i \neq j} c_i c_j (a_3 n - \frac{1}{12}) \\ &= (\sum c_i^2) (a_1 n + a_2) \left\{ 1 + \frac{\sum_{i \neq j} c_i c_j}{\sum c_i^2} \cdot \frac{a_3 n - \frac{1}{12}}{a_1 n + a_2} \right\} \end{aligned}$$

Uit de lemma's 3.1 en 3.2 volgt: $-a_1 \leq (k-2)a_3 \leq 0$ zodat

$\text{var}(\sum_{i=1}^{k-1} c_i \bar{R}_{-i})$ naar oneindig gaat voor $n \rightarrow \infty$ als de volgende twee

voorwaarden zijn vervuld:

$$(3.3) \quad \begin{aligned} &\underline{i} \quad a_1 > 0 \\ &\underline{ii} \quad \sum_i \sum_{j \neq i} c_i c_j < (k-2) \sum_{i=1}^{k-1} c_i^2 \quad \text{of} \quad -a_1 < (k-2) a_3 \end{aligned}$$

Dus als $-a_1 < (k-2)a_3$ (dan ook $a_1 > 0$), dan geldt dat de vektor $(\bar{R}_{-1}, \dots, \bar{R}_{-k-1}, \dots)$ asymptotisch normaal verdeeld is.

Nu geldt: $-a_1 < (k-2)a_3$ d.e.s.d. als:

$$(2p - 2r - p^2)(k-1) + q - p^2 > 0$$

waarin:

$$\begin{aligned} 2p - 2r - p^2 &= \int (1-F)^2 dG - \left(\int GdF \right)^2 = \\ &= \int (1-F)^2 dG - \left(\int (1-F)dG \right)^2 = \int \{F - \int FdG\}^2 dG \geq 0 \\ \text{en } q - p^2 &= \int G^2 dF - \left(\int GdF \right)^2 = \int \{G - \int GdF\}^2 dF \geq 0 \end{aligned}$$

$2p - 2r - p^2$ en $q - p^2$ zijn beide gelijk aan nul d.e.s.d. als

$$(3.4) \quad G \text{ gelijk is aan } 0 \text{ of } 1 \text{ op de drager van } F.$$

We hebben dus dat, behalve wanneer (3.4) geldt, de vektor $(\bar{R}_{-1}, \dots, \bar{R}_{-k-1}, \dots)$ asymptotisch normaal verdeeld is.

Als (3.4) wèl geldt, dan is a_1 gelijk aan $\frac{1}{12}k^2 - \frac{1}{4}k + \frac{1}{6}$ en dus positief (want $k \geq 3$), zodat volgens (3.3) en lemma 3.3. $\sum_{i=1}^{k-1} c_i \bar{R}_{-i}$ asymptotisch

normaal verdeeld is, als niet alle c_i 's gelijk zijn. Verder is dan

$\sum_{i=1}^{k-1} \bar{R}_{-i}$ konstant, zodat, ook als (3.4) geldt, $(\bar{R}_{-1}, \dots, \bar{R}_{-k-1}, \dots)$ asympto-

tisch normaal verdeeld is, maar met één dimensie lager, nl. dimensie $k-2$.

4. De asymptotische verdeling van $\max_{1 \leq i, j \leq k-1} |\bar{R}_i - \bar{R}_j|$

Teneinde de asymptotische verdeling van $\max_{1 \leq i, j \leq k-1} |\bar{R}_i - \bar{R}_j|$ te berekenen voor $n \rightarrow \infty$, beschouwen we de volgende $(k-1)$ -dimensionale stochastische vektor \underline{v}_n gedefinieerd door:

$$\underline{v}_n = (v_{n1}, \dots, v_{n,k-1}) := \left(\frac{\bar{R}_1 - \mathbb{E}\bar{R}_1}{\sqrt{n}}, \dots, \frac{\bar{R}_{k-1} - \mathbb{E}\bar{R}_{k-1}}{\sqrt{n}} \right)$$

Volgens hoofdstuk 3 is \underline{v}_n voor $n \rightarrow \infty$ asymptotisch multi-normaal verdeeld met kovariantiematrix:

$$\begin{pmatrix} a_1 & a_3 & \dots & a_3 \\ & a_1 & & a_3 \\ & & a_1 & a_3 \\ & & & a_1 \end{pmatrix}$$

Definiëren we $w_{ni} := v_{ni} - \gamma \bar{v}_n$

waarin $\gamma := 1 + \sqrt{\frac{a_1 - a_3}{a_1 + (k-2)a_3}}$ en $\bar{v}_n := \frac{1}{k-1} \sum_{j=1}^{k-1} v_{nj}$

dan is $(w_{n1}, \dots, w_{n,k-1})$ voor $n \rightarrow \infty$ asymptotisch multi-normaal verdeeld met als kovariantiematrix $(a_1 - a_3)I_{k-1}$ (waarin met I_{k-1} de identiteitsmatrix van afmeting $k-1$ bedoeld wordt), zodat de range van

$$\frac{w_{n1}}{\sqrt{a_1 - a_3}}, \dots, \frac{w_{n,k-1}}{\sqrt{a_1 - a_3}}$$

voor $n \rightarrow \infty$ $q_{k-1, \infty}$ verdeeld is.

Hieruit volgt dat ook de range van $\frac{\bar{R}_1}{\sqrt{n(a_1 - a_3)}}, \dots, \frac{\bar{R}_{k-1}}{\sqrt{n(a_1 - a_3)}}$ voor $n \rightarrow \infty$

ook $q_{k-1, \infty}$ verdeeld is zodat we vinden voor $n \rightarrow \infty$:

$$(4.1) \quad P\{|\bar{R}_i - \bar{R}_j| < q_{k-1, \infty}^{\alpha} \sqrt{n(a_1 - a_3)} \text{ voor alle } i, j \text{ met } 1 \leq i, j \leq k-1\} = 1 - \alpha$$

Op dezelfde wijze bewijst men dat onder $H_0 (F=G)$ voor $n \rightarrow \infty$ geldt:

$$(4.2) \quad P\{ |\bar{R}_{i\cdot} - \bar{R}_{j\cdot}| < q_{k,\infty}^\alpha \frac{1}{12} k^2 n \text{ voor alle } i, j \} = 1 - \alpha$$

(want onder H_0 is $\frac{\bar{R}_{1\cdot}}{\sqrt{n}}, \dots, \frac{\bar{R}_{k\cdot}}{\sqrt{n}}$ asymptotisch normaal verdeeld

met varianties $\frac{1}{12}k(k-1)$ en kovarianties $-\frac{1}{12}k$)

Onder het door ons beschouwde alternatief $F_1 = \dots = F_{k-1} = F$ en $F_k = G$ is de kans om uit uit $\{F_1, \dots, F_{k-1}\}$ sommige verschillend van elkaar te verklaren gelijk aan:

$$\begin{aligned} & P\{ \max_{1 \leq i, j \leq k-1} |\bar{R}_{i\cdot} - \bar{R}_{j\cdot}| > q_{k,\infty}^\alpha \sqrt{\frac{1}{12}k^2 n} \} = \\ & = P\{ q_{k-1,\infty} > q_{k,\infty}^\alpha \sqrt{\frac{\frac{1}{12}k^2}{a_1 - a_3}} \} \end{aligned}$$

zodat:

$$(4.3) \quad \alpha_{k-1} = \sup_{F, G} P\{ q_{k-1,\infty} > q_{k,\infty}^\alpha \sqrt{\frac{\frac{1}{12}k^2}{a_1 - a_3}} \}$$

F en G continu

Hierin hangen a_1 en a_3 van F en G af (zie(2.1),(2.2) en (3.1)).

Opmerking:

We kunnen nu ook de echte simultane onbetrouwbaarheid voor $n \rightarrow \infty$ berekenen van de methode gebaseerd op formule (1.1). Vanwege de asymptotische χ_{k-1}^2 - verdeling van de toetsingsgrootte van Kruskal-Wallis is die echte onbetrouwbaarheid (die afhangt van α en k) gelijk aan:

$$\begin{aligned} & \lim_{n \rightarrow \infty} P\{ \max_{1 \leq i, j \leq k} |\bar{R}_{i\cdot} - \bar{R}_{j\cdot}| > \sqrt{\frac{\chi^2(\alpha)}{k-1} \frac{k^2 n + k}{6}} \} \\ & = \lim_{n \rightarrow \infty} P\{ \frac{\max_{1 \leq i, j \leq k} |\bar{R}_{i\cdot} - \bar{R}_{j\cdot}|}{\sqrt{\frac{1}{12}k^2 n}} > \sqrt{\frac{\chi^2(\alpha)}{k-1} \frac{2(k^2 n + k)}{k^2 n}} \} \\ & = P\{ q_{k,\infty} > \sqrt{\frac{2\chi^2(\alpha)}{k-1}} \} \end{aligned}$$

Tabel 4.1: Simultane onbetrouwbaarheid van de methode gebaseerd op formule

(1.1) als $n \rightarrow \infty$

k=	3	4	5	6	7	8	9	10	12	15	20
$\alpha=0.10$.081	.060	.042	.029	.019	.012	.0079	.0050	.0019	.00044	.00003
0.05	.038	.027	.018	.011	.0071	.0044	.0027	.0016	.00057	.00011	.000008
0.01	.0068	.0042	.0025	.0014	.00082	.00046	.00025	.00014	.000042	.000007	.0000003

5. Berekening van α_{k-1} zonder restricties op F en G

Als we, behalve continuïteit, geen andere beperkingen opleggen aan de F en G uit formule (4.3) waarover het supremum genomen wordt, dan zal blijken dat voor de gangbare waarden van α geldt $\alpha_{k-1} > \alpha$. (zie tabel 5.1). Dit bewijzen we als volgt:

de kans in formule (4.3) is maximaal voor die F en G waarvoor $a_1 - a_3$ maximaal is.

Nu geldt:

$$(5.1) \quad a_1 - a_3 = \frac{1}{12}k^2 + (2r - p - \frac{1}{6})(k-1) + q - p^2 - \frac{1}{12}$$

waarin:

$$2r-p = \int_{\{x|F(x) \leq \frac{1}{2}\}} (2F-1)GdF + \int_{\{x|F(x) > \frac{1}{2}\}} (2F-1)GdF$$

zodat $2r-p$ maximaal is als er getallen x_0 en x_1 bestaan ($x_0 < x_1$), waarvoor geldt dat:

$$(5.2) \quad \begin{aligned} F(x_0) = F(x_1) = \frac{1}{2} \\ G(x) = 0 \text{ voor } x < x_0 \\ \text{en } G(x) = 1 \text{ voor } x > x_1 \end{aligned}$$

Nu is $q-p^2$ juist ook maximaal in deze situatie, want altijd geldt:

$$(5.3) \quad q-p^2 = \int G^2 dF - p^2 \leq \int GdF - p^2 = p(1-p) \leq \frac{1}{4}$$

en als (5.2) geldt, dan $q - p^2 = \frac{1}{4}$

Voor een F en G als in (5.2) geldt dus:

$$a_1 - a_3 = \frac{1}{12}(k^2 + k + 1)$$

zodat we vinden:

$$(5.4) \quad \alpha_{k-1} = P \left\{ q_{k-1, \infty} > q_{k, \infty}^{\alpha} \sqrt{\frac{k^2}{k^2 + k + 1}} \right\}$$

M.b.v. een tabel van de verdelingsfunctie van de range van normaal verdeelde grootheden, bijvoorbeeld Harter (1969), vinden we:

tabel 5.1: α_{k-1} zonder restricties op F en G voor $\alpha = 0.01$, $\alpha = 0.025$, $\alpha = 0.05$ en $\alpha = 0.10$.

	k = 3	4	5	6	7	8	9	10	12	15	20
$\alpha=0.01$.0153	.0181	.0182	.0178	.0172	.0167	.0162	.0158	.0151	.0143	.0134
0.025	.0303	.0361	.0386	.0385	.0379	.0372	.0365	.0358	.0347	.0334	.0318
0.05	.0512	.0643	.0682	.0690	.0688	.0682	.0674	.0667	.0652	.0633	.0612
0.10	.0877	.1123	.1208	.1240	.1250	.1250	.1245	.1238	.1224	.1202	.1172

6. α_{k-1} voor Lehmann-alternatieven

Nemen we in formule (4.3) het supremum alleen over die paren F,G waarvoor geldt: $G = F^t$ voor zekere $t \in (0, \infty)$, dan zal blijken dat wel geldt:

$$\alpha_{k-1} < \alpha. \text{ (tabel 6.1)}$$

Als: $G = F^t$ dan:

$$p = \int F^t dF = \frac{1}{t+1}$$

$$q = \int F^{2t} dF = \frac{1}{2t+1}$$

$$r = \int F^{t+1} dF = \frac{1}{t+2}$$

zodat invullen in (5.1) geeft:

$$a_1 - a_3 = \frac{1}{12}k^2 + (k-1)\left(\frac{2}{t+2} - \frac{1}{t+1} - \frac{1}{6}\right) + \frac{1}{2t+1} - \frac{1}{(t+1)^2} - \frac{1}{12}$$

In de volgende tabel zijn voor verschillende waarden van k de maximale waarden van $a_1 - a_3 - \frac{k^2}{12}$ als functie van t gegeven:

k=	3	4	5	6	7	8	9	10	12	15	20
$a_1 - a_3 - \frac{k^2}{12}$.01625	.02112	.02600	.03090	.03579	.04069	.04559	.05049	.06030	.07501	.09954

zodat we voor α_{k-1} vinden:

tabel 6.1: α_{k-1} voor Lehmann-alternatieven:

k=	3	4	5	6	7	8	9	10	12	15	20
$\alpha=0.01$.0040	.0057	.0067	.0073	.0078	.0081	.0083	.0085	.0088	.0090	.0093
0.025	.0100	.0144	.0168	.0184	.0194	.0202	.0208	.0213	.0220	.0226	.0233
0.05	.0204	.0291	.0339	.0370	.0391	.0406	.0418	.0427	.0440	.0454	.0466
0.10	.0408	.0597	.0690	.0750	.0790	.0820	.0842	.0860	.0886	.0911	.0935

7. α_{k-1} voor verschuivingsalternatieven

We beschouwen nu alternatieven waarvoor een $a \in \mathbb{R}$ bestaat, zodat $G(x)=F(x-u)$ voor alle $x \in \mathbb{R}$.

Verder nemen we eerst even aan dat $a > 0$ (verschuiving naar rechts).

Dan geldt: $G(x) \leq F(x)$ voor alle x

zodat:

$$\begin{aligned}
 2r - p &= \int G(2F-1) dF = \\
 &= \int_{\{x|F(x) \leq \frac{1}{2}\}} G(2F-1) dF + \int_{\{x|F(x) > \frac{1}{2}\}} G(2F-1) dF \leq \\
 (7.1) \quad &\leq 0 + \int_{\{x|F(x) > \frac{1}{2}\}} (2F^2 - F) dF = \left[\frac{2}{3}F^3 - \frac{1}{2}F^2 \right]_{F=\frac{1}{2}}^{F=1} = \frac{5}{24}
 \end{aligned}$$

verder geldt nog steeds (zie 5.3):

$$q - p^2 \leq \frac{1}{4}$$

We kunnen echter nog wel een scherpere bovengrens voor $q-p^2$ vinden (waarvan bovendien de geldigheid niet alleen beperkt is tot verschuivingsalternatieven):

Stelling 7: 1:

$$(7.2) \quad q - p^2 \leq 2r - p$$

Bewijs:
$$\begin{aligned}
 \int G^2 dF + \int F^2 dG &= P(Y_{11} > Y_{k1} \cap Y_{11} > Y_{k2}) + P(Y_{11} < Y_{k2} \cap Y_{12} < Y_{k2}) = \\
 &= P\{(Y_{11} > Y_{k1} \cap Y_{11} > Y_{k2}) \cup (Y_{11} < Y_{k2} \cap Y_{12} < Y_{k2})\} \leq \\
 &\leq P\{Y_{11} > Y_{k1} \cup Y_{12} < Y_{k2}\} =
 \end{aligned}$$

$$\begin{aligned}
 &= 1 - P\{Y_{11} < Y_{k1} \cap Y_{12} > Y_{k2}\} = \\
 (7.3) \quad &= 1 - p \cdot (1-p)
 \end{aligned}$$

Het laatste gelijkteken geldt vanwege de onafhankelijkheid van Y_{11} , Y_{12} , Y_{k1} en Y_{k2}

We vinden m.b.v. (7.3):

$$q - p^2 = \int G^2 dF - p^2 \leq - \int F^2 dG + 1 - p(1-p) - p^2 = 2r - p$$

waarbij het laatste gelijkteken volgt uit:

$$\int F^2 dG = 1 - 2 \int FG dF \quad (\text{partiële integratie})$$

Q.E.D.

Vullen we (7.1) en (7.2) in in (5.1) dan volgt:

$$(7.4) \quad a_1 - a_3 \leq \frac{1}{12}k^2 + \frac{1}{24}(k-1) + \frac{1}{8} = \frac{1}{12}(k^2 + \frac{1}{2}k + 1)$$

voor $G(x) = F(x-a)$ met $a > 0$.

Als $a < 0$ (verschuiving naar links) daar kijken we naar:

$$-Y_{11}, \dots, -Y_{1n}; \dots; -Y_{k1}, \dots; -Y_{kn}$$

Als F' de verdelingsfunctie is van $-Y_{11}$ en G' die van $-Y_{k1}$ dan geldt:

$$\begin{aligned}
 F'(x) &= 1-F(-x) \\
 \text{en } G'(x) &= F'(x+a)
 \end{aligned}$$

G' is dus een verschuiving van F' naar rechts (omdat $a < 0$). Verder geldt: als \underline{R}'_{ij} het rangnummer is van $-Y_{ij}$ onder alle $-Y_{11}, \dots, -Y_{kn}$ dan:

$$\underline{R}'_{ij} = N - \underline{R}_{ij} + 1$$

$$\begin{aligned}
 \text{zodat } a_1 - a_3 &= \lim_{n \rightarrow \infty} \frac{1}{2n} \text{var}(\bar{R}_1 - \bar{R}_2) = \lim_{n \rightarrow \infty} \frac{1}{2n} \text{var}(\bar{R}'_1 - \bar{R}'_2) \\
 &= a'_1 - a'_3
 \end{aligned}$$

Nu geldt dat $a'_1 - a'_3 \leq \frac{1}{12}(k^2 + \frac{1}{2}k+1)$, omdat G' een verschuiving naar rechts is van F' , zodat (7.4) ook geldt als $a < 0$.

(7.4) geeft een bovengrens voor $a_1 - a_3$ bij verschuivingsalternatieven. Dat deze bovengrens niet erg veel verlaagd kan worden, blijkt uit het volgende voorbeeld:

voorbeeld 7.1:

Neem F gedefinieerd door

$$F(x) = \begin{cases} x + \frac{1}{2} & \text{als } -\frac{1}{2} \leq x \leq 0 \\ \frac{x}{c} + \frac{1}{2} & \text{als } 0 \leq x \leq \frac{1}{2}c \end{cases}$$

en G gedefinieerd door $G(x) = F(x - \frac{1}{2})$

dan geldt (onderstel $c > \frac{1}{2}$):

$$\begin{aligned} 2r - p &= \int_0^{\frac{1}{2}c} G(2F-1)dF \rightarrow \int_0^{\frac{1}{2}c} (F(x) - \frac{1}{2}c)(2F(x) - 1) dF(x) = \\ &= \left[\frac{2}{3}F^3 - \frac{1}{2}(1+c)F^2 + \frac{1}{2}c F \right]_{F=\frac{1}{2}}^{F=\frac{1}{2}c} = \frac{5}{24} - \frac{c}{4} \end{aligned}$$

en $q - p^2 = \frac{29}{192} + O(\frac{1}{c})$

zonder geldt:

$$a_1 - a_3 = \frac{1}{12}k^2 + \frac{1}{24}k + \frac{5}{192} + O(\frac{1}{c})$$

Bij deze F en G geldt dus voor $c \rightarrow \infty$:

$$a_1 - a_3 \rightarrow \frac{1}{12}(k^2 + \frac{1}{2}k + \frac{5}{16})$$

Als we deze waarde van $a_1 - a_3$ en die uit (7.4) invullen in (4.3), dan vinden we de volgende onder- en bovengrenzen voor α_p :

Tabel 7.1

Onder- en bovengrenzen voor α_{k-1} bij verschuivingsalternatieven.

	$\alpha=0.01$	$\alpha=0.025$	$\alpha=0.05$	$\alpha=0.10$
k=3	.0079-.0100	.0175-.0213	.0325-.0381	.0612-.0695
4	.0101-.0119	.0230-.0263	.0431-.0483	.0816-.0894
5	.0109-.0123	.0253-.0279	.0478-.0519	.0909-.0972
6	.0113-.0124	.0263-.0284	.0501-.0535	.0958-.1009
7	.0114-.0123	.0268-.0285	.0514-.0541	.0987-.1030
8	.0114-.0121	.0271-.0285	.0521-.0544	.1005-.1040
9	.0114-.0120	.0273-.0284	.0526-.0544	.1019-.1046
10	.0114-.0119	.0273-.0283	.0529-.0544	.1025-.1049
12	.0114-.0117	.0273-.0281	.0531-.0543	.1034-.1053
15	.0112-.0115	.0272-.0277	.0532-.0540	.1039-.1052
20	.0111-.0112	.0270-.0273	.0530-.0534	.1041-.1049

We zien in tabel 7.1 dat voor verschuivingsalternatieven in het algemeen niet geldt dat $\alpha_{k-1} \leq \alpha$ (hoewel het niet veel scheelt). Als we echter twee extra eisen aan F opleggen (die in de praktijk wel eens vervuld zijn), dan blijkt dat altijd geldt dat $\alpha_{k-1} \leq \alpha$.

Stelling 7.2: Als G een verschuiving van F is, waarbij F symmetrisch en unimodaal is (en kontinu), dan geldt:

$$(7.5) \quad 2r - p \leq \frac{1}{6}$$

Bewijs:

Daar het probleem translatie-invariant is, is het geen beperking om F symmetrisch in $x = 0$ te nemen.

$$2r - p = \int_{-\infty}^0 F(x-a)(2F(x)-1)dF(x) + \int_0^{\infty} F(x-u)(2F(x)-1)dF(x)$$

Nu geldt: (omdat $F(-x) = 1 - F(x)$):

$$\begin{aligned} \int_{-\infty}^0 F(x-a)(2F(x)-1) &= \int_{+\infty}^0 F(-y-a) \{2F(-y)-1\} dF(-y) = \\ &= \int_0^{\infty} F(-y-a) \{1-2F(y)\} dF(y) = \int_0^{\infty} \{F(y+a) - 1\} \{2F(y)-1\} dF(y) \end{aligned}$$

zodat

$$(7.6) \quad 2r - p = \int_0^{\infty} \{ F(x+a) + F(x-a) - 1 \} \{2F(x)-1\} dF(x)$$

$$(7.7) \quad \text{Voor } x \geq 0 \text{ geldt: } F(x+a) + F(x-a) \leq 2F(x)$$

Als $x \geq |a|$ geldt dit vanwege de unimodaliteit van F, hetgeen betekent dat F konkaaf is voor $x \geq 0$.

Voor $0 \leq x < |a|$ heeft men de symmetrie en unimodaliteit beide nodig om in te zien dat (7.7) geldt.

Door (7.7) in te vullen in (7.6) vinden we:

$$2r - p \leq \int_0^{\infty} (2F(x)-1)^2 dF(x) = \left[\frac{1}{6}(2F(x) - 1)^3 \right]_{F(x)=\frac{1}{2}}^{F(x)=1} = \frac{1}{6}$$

Q.E.D.

Met behulp van stelling 7.1 vinden we dat onder de voorwaarden van stelling 7.2 ook geldt:

$$(7.8) \quad q - p^2 \leq \frac{1}{6}$$

Invullen van (7.5) en (7.8) in (5.1) geeft:

$$a_1 - a_3 \leq \frac{1}{12}k^2 + \frac{1}{12}$$

zodat geldt:

$$\alpha_{k-1} \leq P\left\{q_{k-1,\infty} > q_{k,\infty}^{\alpha} \sqrt{\frac{k^2}{k^2+1}}\right\}$$

tabel 7.2:

Bovengrens voor α_{k-1} voor verschuivingen van symmetrische, unimodale verdelingen.

k=	3	4	5	6	7	8	9	10	12	15	20
$\alpha=0.01$.0057	.0071	.0077	.0081	.0083	.0085	.0086	.0087	.0089	.0091	.0093
0.025	.0135	.0173	.0190	.0200	.0206	.0211	.0214	.0217	.0222	.0227	.0232
0.05	.0262	.0349	.0375	.0396	.0410	.0421	.0429	.0435	.0445	.0455	.0465
0.10	.0498	.0673	.0749	.0793	.0823	.0845	.0860	.0872	.0893	.0913	.0933

Opmerking: Als we de symmetrie-eis laten vallen, dan zitten we weer in de situatie van tabel 7.1, omdat de F uit voorbeeld 7.1 unimodaal is.

Literatuur:

Hajek J., Asymptotic normality of simple linear rank statistics under alternatives, The Annals of Mathematical Statistics, vol.39 (1968), pag. 325-346.

Harter H.L, Order statistics and their use in testing and estimation, vol I, Aerospace Research Laboratories, Government Printing Office, Washington (1969)

Miller R.G, Simultaneous statistical inference, McGraw-Hill, New York (1966)

Nemenyi P., Distribution-free multiple comparisons.
Unpublished doctoral thesis, Princeton University (1963)