

MATHEMATISCH CENTRUM

2e BOERHAAVESTRAAT 49

AMSTERDAM

STATISTISCHE AFDELING

Leiding: Prof. Dr D. van Dantzig

Chef van de Statistische Consultatie: Prof. Dr J. Hemelrijk

Rapport S 105 (M 41)

Uittreksel en uitwerking van het artikel

"A class of statistics with asymptotical-
ly normal distributions", door W.HOFFDING,

Ann. Math. Stat. 19 (1948), pp. 293-325.

door

D.J.Stoker

Een klasse van statistische grootheden met asymptotische normale verdeling met speciale verwijzing naar parameter vrije toetsingsmethoden voor onafhankelijkheid.

1. Definitie van een U -statistische grootheid.

Stel $x_v, v=1, \dots, n$ is een steekproef van n vectoren

$$(1) \quad x_v = (x_v^{(1)}, \dots, x_v^{(r)})$$

en $\Phi(x_1, \dots, x_m)$ een functie van $m (\leq n)$ vectorargumenten. De functie van de steekproef

$$(2) \quad U = U(x_1, \dots, x_n) = \frac{1}{n(n-1)\dots(n-m+1)} \sum'' \Phi(x_{\alpha_1}, \dots, x_{\alpha_m}),$$

waar \sum'' een sommatie over alle permutaties $(\alpha_1, \dots, \alpha_m)$ van gehele getallen aangeeft zodanig dat

$$(3) \quad 1 \leq \alpha_i \leq n, \quad \alpha_i \neq \alpha_j \text{ as } i \neq j, \quad (i, j = 1, \dots, m),$$

een U -statistische grootheid wordt genoemd. U is het gemiddelde van de waarden van Φ in de verzameling van gerangschikte deelverzamelingen van m elementen van de steekproef (1). U is symmetrisch in x_1, \dots, x_n (een functie $f = f(x_1, \dots, x_n)$ wordt symmetrisch genoemd indien deze invariant is ten aanzien van alle permutaties van zijn argumenten).

Een functie $\Phi(x_1, \dots, x_m)$ die aan (2) voldoet, wordt een kern van de statistische grootheid U genoemd.

Stel \mathcal{D} is een deelverzameling van alle verdelingsfuncties (v.f.'s) in de r -dimensionale ruimte waartoe $F = F(x) =$

$= F(x^{(1)}, \dots, x^{(r)})$ behoort. $\Theta(F)$ wordt een functioneel van F gedefinieerd op \mathcal{D} , genoemd indien aan iedere F die tot \mathcal{D} behoort, een grootheid $\Theta(F)$ wordt toegekend. Stel verder dat voor een functioneel $\Theta = \Theta(F)$ van $F(x)$ een steekproefgrootte

n bestaat waarvoor een zuivere schatting van Θ voor iedere v.f. F in \mathcal{D} bestaat. Dit wil zeggen, als x_1, \dots, x_n n onafhankelijke stochastische vectoren met dezelfde v.f. F zijn, dan bestaat er een functie $\Phi(x_1, \dots, x_n)$ van n vectorargumenten (1) zodanig dat

$$(4) \quad E \int_{R_n} \Phi(x_1, \dots, x_n) dF(x_1) \dots dF(x_n) = \Theta(F)$$

voor iedere F in \mathcal{D} . $\int_{R_n} \Phi(x_1, \dots, x_n) dF(x_1) \dots dF(x_n)$ is dus een zuivere schatting van Θ op \mathcal{D} . Een functioneel $\Theta(F)$ van de vorm (4) noemen wij regelmatig op \mathcal{D} .

Stel $m (\leq n)$ is de kleinste steekproefgrootte waarvoor een zuivere schatting $\Phi(x_1, \dots, x_m)$ van de regelmatige functioneel Θ op \mathcal{D} bestaat, d.w.z.

$$(5) \quad \Theta(F) = \int \dots \int \Phi(x_1, \dots, x_m) dF(x_1) \dots dF(x_m)$$

voor iedere F in \mathcal{D} , m wordt de graad van de regelmatige functioneel op \mathcal{D} genoemd.

Iedere functie $\Phi(x_1, \dots, x_m)$ die aan (5) voldoet, wordt de kern van de regelmatige functioneel $\Theta(F)$ genoemd.

$\Phi_0(x_1, \dots, x_m)$ is een symmetrische kern in x_1, \dots, x_m van $\Theta(F)$ indien

$$(6) \quad \Phi_0(x_1, \dots, x_m) = \frac{1}{m!} \sum'' \Phi(x_{\alpha_1}, \dots, x_{\alpha_m}),$$

waar de som \sum'' zich over alle permutaties $(\alpha_1, \dots, \alpha_m)$ van $(1, \dots, m)$ uitstrekt.

Voor $n=m$, reduceert de statistische grootheid U zich dus tot de symmetrische kern (6) van $\Theta(F)$.

Halmos [3] heeft zekere optimale eigenschappen van U -statistische grootheden als zuivere schattingen van regelmatige functionelen bewezen, n.l. dat, wanneer $\Theta(F)$ een regelmatige functioneel van de m^e graad op een verzameling \mathcal{D} , die alle geheel en al discontinue verdelingsfuncties bevat, is, U dan de enige zuivere schatting op \mathcal{D} is die symmetrisch in x_1, \dots, x_n is en U heeft de kleinste variantie van alle zuivere schattingen op \mathcal{D} , met andere woorden, U is de doeltreffendste schatting van $\Theta(F)$ op \mathcal{D} .

Uit (2) en (6) volgt nu dat wij een U -statistische grootheid in de volgende vorm kunnen schrijven:

$$(7) \quad U(x_1, \dots, x_n) = \binom{n}{m}^{-1} \sum' \Phi_0(x_{\alpha_1}, \dots, x_{\alpha_m})$$

daar de kern Φ_0 symmetrisch in zijn m vectorargumenten is en de som \sum' zich over alle indices α uitstrekt, zodanig dat

$$1 \leq \alpha_1 < \alpha_2 < \dots < \alpha_m \leq n.$$

2. De variantie van een U -statistische grootheid.

Stel x_1, \dots, x_n zijn n onafhankelijke stochastische vectoren met dezelfde v.f. $F(x) = F(x^{(1)}, \dots, x^{(n)})$ en stel

$$(8) \quad \underline{U} = U(x_1, \dots, x_n) = \binom{n}{m}^{-1} \sum' \Phi(x_{\alpha_1}, \dots, x_{\alpha_m})$$

daar $\Phi(x_1, \dots, x_m)$ symmetrisch in x_1, \dots, x_m is en \sum' dezelfde betekenis als in (7) heeft. Stel dat de functie Φ n niet bevat. Nu is $\mathcal{E} \underline{U} = \mathcal{E} \Phi(x_1, \dots, x_m) = \Theta = \Theta(F)$.

Stel

$$(9) \quad \Phi_c(x_1, \dots, x_c) = \mathcal{E} \Phi(x_1, \dots, x_c, x_{c+1}, \dots, x_m), \quad (c=1, \dots, m)$$

daar x_1, \dots, x_c willekeurig vaste vectoren zijn en de verwachting ten opzichte van de stochastische vectoren x_{c+1}, \dots, x_m genomen

wordt.

Dan is

$$(10) \quad \Phi_{c-1}(x_1, \dots, x_{c-1}) = \sum \Phi_c(x_1, \dots, x_{c-1}, x_c),$$

en

$$(11) \quad \sum \Phi_c(x_1, \dots, x_c) = \Theta \quad , \quad (c=1, \dots, m)$$

Definieer nu

$$(12) \quad \Psi(x_1, \dots, x_m) = \Phi(x_1, \dots, x_m) - \Theta$$

$$(13) \quad \Psi_c(x_1, \dots, x_c) = \Phi_c(x_1, \dots, x_c) - \Theta \quad , \quad (c=1, \dots, m).$$

Nu is ook

$$(14) \quad \Psi_{c-1}(x_1, \dots, x_{c-1}) = \sum \Psi_c(x_1, \dots, x_{c-1}, x_c),$$

$$(15) \quad \sum \Psi_c(x_1, \dots, x_c) = \sum \Psi(x_1, \dots, x_m) = 0 \quad , \quad (c=1, \dots, m).$$

Veronderstel dat de variantie van $\Psi_c(x_1, \dots, x_c)$ bestaat en stel

$$(16) \quad J_0 = 0 \quad , \quad J_c = \sum \Psi_c^2(x_1, \dots, x_c) \quad , \quad (c=1, \dots, m)$$

Nu is

$$J_c = \sum \Phi_c^2(x_1, \dots, x_c) - \Theta^2.$$

Als voor een zekere oorspronkelijke verdeling $F = F_0$ en een zeker geheel getal d $J_d(F_0) = 0$ is, betekent dit dat

$\Psi_d(x_1, \dots, x_d) = 0$ is met een waarschijnlijkheid 1. Uit (14) en (16) volgt dat $J_d = 0$ betekent dat $J_1 = \dots = J_{d-1} = 0$ is.

Als $J_1(F_0) = 0$ is, is de regelmatige functioneel $\Theta(F)$ stationnair voor $F = F_0$. Als

$$(17) \quad J_1(F_0) = \dots = J_d(F_0) = 0 \quad , \quad J_{d+1}(F_0) > 0 \quad , \quad (1 \leq d \leq m)$$

is $\Theta(F)$ stationnair van orde d voor $F = F_0$.

Als $(\alpha_1, \dots, \alpha_m)$ en $(\beta_1, \dots, \beta_m)$ twee verzamelingen zijn bestaande uit m verschillende gehele getallen $1 \leq \alpha_i, \beta_i \leq n$ en c het aantal gehele getallen is dat beide verzamelingen gemeen hebben, volgt uit de symmetrie van Ψ dat

$$(18) \quad E \{ \psi(x_{\alpha_1}, \dots, x_{\alpha_m}) \psi(x_{\beta_1}, \dots, x_{\beta_m}) \} = J_c$$

Als de variantie van \underline{U} bestaat, dan is deze gelijk aan

$$\begin{aligned} \sigma^2(\underline{U}) &= E(\underline{U}^2) - E^2(\underline{U}) \\ &= \binom{n}{m}^{-2} E \left\{ \sum' \Phi(x_{\alpha_1}, \dots, x_{\alpha_m}) \right\}^2 - \theta^2 \\ &= \binom{n}{m}^{-2} E \left\{ \sum' \psi(x_{\alpha_1}, \dots, x_{\alpha_m}) \right\}^2 \\ &= \binom{n}{m}^{-2} \sum_{c=0}^m \sum^{(c)} E \{ \psi(x_{\alpha_1}, \dots, x_{\alpha_m}) \psi(x_{\beta_1}, \dots, x_{\beta_m}) \} \end{aligned}$$

waar $\sum^{(c)}$ de sommatie over alle indices aanduidt, zodanig dat

$$1 \leq \alpha_1 < \alpha_2 < \dots < \alpha_m \leq n, \quad 1 \leq \beta_1 < \dots < \beta_m \leq n$$

en precies aan c vergelijkingen

$$\alpha_i = \beta_j$$

wordt voldaan. Uit (18) volgt dat elke term in $\sum^{(c)}$ gelijk is aan J_c . Het aantal termen in $\sum^{(c)}$ is gelijk aan

$$\binom{m}{c} \binom{n-m}{m-c} \binom{n}{m} = \frac{n(n-1)\dots(n-2m+c+1)}{c! (m-c)! (m-c)!}$$

(De m α 's kunnen op $\binom{n}{m}$ manieren uit de getallen $1, 2, \dots, n$ gekozen worden. c van de m α 's kunnen op $\binom{m}{c}$ manieren gekozen worden. Nu zijn c waarden van de β 's gekozen en de overblijvende $m-c$ β 's mogen niet gelijk zijn aan de m α 's. Dus kunnen de overblijvende β 's op $\binom{n-m}{m-c}$ manieren gekozen worden. Let goed op dat de α 's en β 's alle onderling verschillend zijn.)

Omdat $J_0 = 0$ is, is

$$(19) \quad \sigma^2(\underline{U}) = \binom{n}{m}^{-1} \sum_{c=1}^m \binom{m}{c} \binom{n-m}{m-c} J_c$$

Een overeenkomstige formulering kan gemakkelijk afgeleid worden voor het geval dat de verdelingen $F_\nu(x)$ van x_ν , $\nu=1, \dots, n$, verschillend zijn.

Hoefding [5] bewijst vervolgens de volgende twee stellingen aangaande de grootheden J_1, \dots, J_m en $\sigma^2(\underline{U})$.

Stelling A.

De grootheden J_1, \dots, J_m zoals gedefinieerd in (16) voldoen aan de ongelijkheden

$$(20) \quad 0 \leq \frac{J_c}{c} \leq \frac{J_d}{d} \quad \text{als} \quad 1 \leq c < d \leq m.$$

Stelling B.

De variantie $\sigma^2(\underline{U}_n)$ van een U -statistische grootheid $\underline{U}_n = U(x_1, \dots, x_n)$, waar x_1, \dots, x_n onafhankelijk en identiek verdeeld is, voldoet aan de ongelijkheden

$$(21) \quad \frac{m^2}{n} f_1 \leq \sigma^2(\underline{U}_n) \leq \frac{m}{n} f_m.$$

$n \sigma^2(\underline{U}_n)$ is een afnemende functie van n ,

$$(22) \quad (n+1) \sigma^2(\underline{U}_{n+1}) \leq n \sigma^2(\underline{U}_n),$$

die voor $n=m$ zijn bovenste grens $m f_m$ aanneemt en naar zijn onderste grens $m^2 f_1$ nadert als n toeneemt:

$$(23) \quad \sigma^2(\underline{U}_m) = f_m,$$

$$(24) \quad \lim_{n \rightarrow \infty} n \sigma^2(\underline{U}_n) = m^2 f_1.$$

Als $\sum \underline{U}_n = \Theta(F)$ stationnair is van de orde $\geq d-1$ voor de v.f. van \underline{x}_α , kan (21) vervangen worden door

$$(25) \quad \frac{m}{d} K_n(m, d) f_d \leq \sigma^2(\underline{U}_n) \leq K_n(m, d) f_m$$

waarbij

$$(26) \quad K_n(m, d) = \binom{n}{m}^{-1} \sum_{c=d}^m \binom{m-1}{c-1} \binom{n-m}{m-c}.$$

Opmerking: Uit (19) en (20) volgt dat een nodige en voldoende voorwaarde voor het bestaan van $\sigma^2(\underline{U})$ het bestaan van

$$(27) \quad f_m = E \left\{ \Phi(x_1, \dots, x_m) \right\}^2 - \theta^2$$

of van $E \left\{ \Phi^2(x_1, \dots, x_m) \right\}$ is.

Voor $f_1 > 0$ is $\sigma^2(\underline{U}) = O(n^{-1})$ (zie o.a. verg. (21)).

Als $\Theta(F)$ stationnair is van de d^e orde voor $F = F_0$, d.w.z. als aan (17) voldaan is, dan is $\sigma^2(\underline{U})$ van de orde $n^{-(d+1)}$.

Alleen wanneer, voor een zekere $F = F_0$, $\Theta(F)$ stationnair is van de m^e orde, waarbij m de graad van $\Theta(F)$ is, krijgen wij $\sigma^2(\underline{U}) = 0$ en \underline{U} is dan gelijk aan een constante (spr 0).

3. De covariantie van twee \underline{U} -statistische grootheden.

Beschouw een verzameling van g \underline{U} -statistische grootheden

$$\underline{U}^{(\gamma)} = \binom{n}{m(\gamma)}^{-1} \sum \Phi^{(\gamma)}(x_{\alpha_1}, \dots, x_{\alpha_{m(\gamma)}}), \quad (\gamma = 1, \dots, g),$$

waar $\underline{U}^{(\gamma)}$ een functie van dezelfde n onafhankelijke gelijk verdeelde stochastische vectoren $\underline{x}_1, \dots, \underline{x}_n$ is. Uit de veronderstelling volgt dat $\Phi^{(\gamma)}$ symmetrisch in zijn $m(\gamma)$ argumenten

$(\gamma = 1, \dots, g)$ is.

Stel

$$E \left\{ \underline{U}^{(\gamma)} \right\} = E \left\{ \Phi^{(\gamma)}(x_1, \dots, x_{m(\gamma)}) \right\} = \theta^{(\gamma)} \quad (\gamma = 1, \dots, g)$$

$$(28) \quad \psi^{(\gamma)}(x_1, \dots, x_{m(\gamma)}) = \Phi^{(\gamma)}(x_1, \dots, x_{m(\gamma)}) - \theta^{(\gamma)} \quad (\gamma=1, \dots, g)$$

$$(29) \quad \psi_c^{(\gamma)}(x_1, \dots, x_c) = \mathcal{E} \left\{ \psi^{(\gamma)}(x_1, \dots, x_c, \underline{x}_{c+1}, \dots, \underline{x}_m) \right\}$$

$$(c=1, \dots, m(\gamma); \gamma=1, \dots, g)$$

$$(30) \quad \zeta_c^{(\gamma, \delta)} = \mathcal{E} \left\{ \psi_c^{(\gamma)}(\underline{x}_1, \dots, \underline{x}_c) \psi_c^{(\delta)}(\underline{x}_1, \dots, \underline{x}_c) \right\} \quad (\gamma, \delta=1, \dots, g)$$

Voor $\gamma = \delta$ is

$$(31) \quad \zeta_c^{(\gamma)} = \zeta_c^{(\gamma, \gamma)} = \mathcal{E} \left\{ \psi_c^{(\gamma)}(\underline{x}_1, \dots, \underline{x}_c) \right\}^2$$

$$\text{Stel } \sigma(\underline{u}^{(\gamma)}, \underline{u}^{(\delta)}) = \mathcal{E} \left\{ (\underline{u}^{(\gamma)} - \theta^{(\gamma)}) (\underline{u}^{(\delta)} - \theta^{(\delta)}) \right\}$$

is de covariantie van $\underline{u}^{(\gamma)}$ en $\underline{u}^{(\delta)}$.

Soortgelijk als voor de variantie, wordt voor $m(\gamma) \leq m(\delta)$ gevonden

$$(32) \quad \sigma(\underline{u}^{(\gamma)}, \underline{u}^{(\delta)}) = \binom{n}{m(\gamma)}^{-1} \sum_{c=1}^{m(\gamma)} \binom{m(\delta)}{c} \binom{n-m(\delta)}{m(\gamma)-c} \zeta_c^{(\gamma, \delta)}$$

De rechterkant is symmetrisch in γ en δ .

Voor $\gamma = \delta$ is (32) de variantie van $\underline{u}^{(\gamma)}$ (zie (19)).

Uit (24) en (32) volgt dat

$$\lim_{n \rightarrow \infty} n \sigma^2(\underline{u}^{(\gamma)}) = m^2(\gamma) \zeta_1^{(\gamma)}$$

$$\lim_{n \rightarrow \infty} n \sigma(\underline{u}^{(\gamma)}, \underline{u}^{(\delta)}) = m(\gamma) m(\delta) \zeta_1^{(\gamma, \delta)}$$

Als dus $\zeta_1^{(\gamma)} \neq 0$ en $\zeta_1^{(\delta)} \neq 0$, nadert de productmomentcorrelatie $\rho(\underline{u}^{(\gamma)}, \underline{u}^{(\delta)})$ tussen $\underline{u}^{(\gamma)}$ en $\underline{u}^{(\delta)}$ naar de limiet

$$(33) \quad \lim_{n \rightarrow \infty} \rho(\underline{u}^{(\gamma)}, \underline{u}^{(\delta)}) = \frac{\zeta_1^{(\gamma, \delta)}}{\sqrt{\zeta_1^{(\gamma)} \zeta_1^{(\delta)}}}$$

4. Limietstellingen voor het geval van gelijk verdeelde \underline{x}_n 's.

Hoefding bewijst de volgende vijf stellingen in verband met de asymptotische verdeling van \underline{u} -statistische grootheden en zekere aanverwante functies.

Stelling 1.1:

Stel $\underline{x}_1, \dots, \underline{x}_n$ zijn n onafhankelijke, gelijk verdeelde stochastische vectoren

$$\underline{x}_\alpha = (x_\alpha^{(1)}, \dots, x_\alpha^{(n)}) \quad (\alpha = 1, \dots, n)$$

Stel verder

$$\Phi^{(\gamma)}(x_1, \dots, x_{m(\gamma)}) \quad (\gamma = 1, \dots, g)$$

zijn g reële-waarde functies die n niet bevatten en waarbij $\Phi^{(\gamma)}$ symmetrisch in zijn $m(\gamma)$ ($\leq n$) vectorargumenten $x_\alpha = (x_\alpha^{(1)}, \dots, x_\alpha^{(n)})$ ($\alpha = 1, \dots, m(\gamma)$; $\gamma = 1, \dots, g$) is. Definieer nu

$$\underline{U}^{(\gamma)} = \binom{n}{m(\gamma)}^{-1} \sum \Phi^{(\gamma)}(x_{\alpha_1}, \dots, x_{\alpha_{m(\gamma)}}), \quad (\gamma = 1, \dots, g)$$

waarbij de sommatie zich over alle indices uitstrekt zodanig dat $1 \leq \alpha_1 < \dots < \alpha_{m(\gamma)} \leq n$. Indien de verwachtingen

$$\theta^{(\gamma)} = \mathbb{E} \left\{ \Phi^{(\gamma)}(x_1, \dots, x_{m(\gamma)}) \right\} \quad (\gamma = 1, \dots, g)$$

en

$$\mathbb{E} \left\{ \Phi^{(\gamma)}(x_1, \dots, x_{m(\gamma)}) \right\}^2 \quad (\gamma = 1, \dots, g)$$

bestaan, dan nadert de simultane verdelingsfunctie van

$$\sqrt{n}(\underline{U}^{(1)} - \theta^{(1)}), \dots, \sqrt{n}(\underline{U}^{(g)} - \theta^{(g)})$$

voor $n \rightarrow \infty$, tot de g -dimensionale normale verdelingsfunctie met gemiddelde 0 en $\left[m(\gamma)m(\delta) \gamma^{(\gamma, \delta)} \right]$, waarbij $\gamma^{(\gamma, \delta)}$ door (30) gedefinieerd is, als covariantiematrix. De limietverdeling is niet singulier als de determinant $|\gamma^{(\gamma, \delta)}| > 0$.

Opmerkingen: Voor $g=1$ volgt uit de stelling dat de verdeling van een \underline{U} -statistische grootte onder zekere voorwaarden tot de normale vorm nadert. Voor $m=1$ is \underline{U} de som van n onafhankelijke stochastische veranderlijken en is stelling 1.1 dus de Centrale Limietstelling voor zulke sommen. Voor $m>1$ is \underline{U} de som van stochastische veranderlijken die in het algemeen niet onafhankelijk zijn.

De stelling toont aan dat in het geval van onderling onafhankelijke en gelijk verdeelde x_α 's, het bestaan van $\mathbb{E} \left\{ \Phi^2(x_1, \dots, x_m) \right\}$ voldoende voor de asymptotische normaliteit van \underline{U} is. Aan de functie Φ worden geen regelmatigheidsvoorwaarden opgelegd.

Volgens stelling B overschrijdt $\sigma^2(\underline{U})$ zijn asymptotische waarde $\frac{m^2}{n} \gamma$, voor iedere eindige n . Als wij dus stelling 1.1 gebruiken als een benadering van de verdeling van \underline{U} voor n groot maar eindig, dan onderschatten wij de variantie van \underline{U} . In vele toepassingen is dit ongewenst en moeten wij liever de volgende stelling gebruiken:

Stelling 1.2:

Onder de voorwaarden van stelling 1.1 en als

$$\xi_{\gamma}^{(\gamma)} > 0, \quad (\gamma = 1, \dots, g)$$

nadert de simultane v.f. van

$$\frac{\underline{u}^{(1)} - \theta^{(1)}}{\sigma(\underline{u}^{(1)})}, \dots, \frac{\underline{u}^{(g)} - \theta^{(g)}}{\sigma(\underline{u}^{(g)})}$$

voor $n \rightarrow \infty$ tot de g -dimensionale normale v.f. met gemiddelde 0 en $[\xi_{\gamma, \delta}^{(\gamma, \delta)}]$ als covariantiematrix, waarbij

$$\xi_{\gamma, \delta}^{(\gamma, \delta)} = \lim_{n \rightarrow \infty} \frac{\sigma(\underline{u}^{(\gamma)}, \underline{u}^{(\delta)})}{\sigma(\underline{u}^{(\gamma)}) \sigma(\underline{u}^{(\delta)})} = \frac{\xi_{\gamma, \delta}^{(\gamma, \delta)}}{\sqrt{\xi_{\gamma, \gamma}^{(\gamma, \gamma)} \xi_{\delta, \delta}^{(\delta, \delta)}}},$$

$(\gamma, \delta = 1, \dots, g)$

Opmerking: De stellingen (1.1) en (1.2) betreffen dus de asymptotische verdeling van $\underline{u}^{(1)}, \dots, \underline{u}^{(g)}$ die zuivere schattingen zijn van $\theta^{(1)}, \dots, \theta^{(g)}$. De zuiverheid van een statistische grootheid speelt natuurlijk geen rol bij het asymptotisch gedrag. Stelling (1.1) kan uitgebreid worden tot een grotere klasse van statistische grootheden, vandaar

Stelling 1.3:

Stel $\underline{u}^{(\gamma)'} = \underline{u}^{(\gamma)} + \frac{\underline{b}_n^{(\gamma)'}}{\sqrt{n}}, \quad (\gamma = 1, \dots, g)$

waarbij $\underline{u}^{(\gamma)}$ in stelling (1.1) gedefinieerd is en $\underline{b}_n^{(\gamma)'}$ een aselechte veranderlijke is. Indien aan de voorwaarden van stelling (1.1) is voldaan en $\lim_{n \rightarrow \infty} E \{ \underline{b}_n^{(\gamma)'} \}^2 = 0, \quad (\gamma = 1, \dots, g)$ is,

dan nadert de simultane verdeling van

$$\sqrt{n} (\underline{u}^{(1)'} - \theta^{(1)}), \dots, \sqrt{n} (\underline{u}^{(g)'} - \theta^{(g)})$$

tot de normale verdeling met gemiddelde 0 en $\{m(\gamma) m(\delta) \xi_{\gamma, \delta}^{(\gamma, \delta)}\}$ als covariantiematrix.

Stelling 1.4:

Stel $\underline{x}_1, \dots, \underline{x}_n$ is een aselechte steekproef uit een α -dimensionale verzameling met $F(x) = F(x^{(1)}, \dots, x^{(\alpha)})$ als v.f. en stel

$$\theta^{(\gamma)}(F) = \int \dots \int \Phi^{(\gamma)}(x_1, \dots, x_{m(\gamma)}) dF(x_1) \dots dF(x_{m(\gamma)})$$

$(\gamma = 1, \dots, g)$

zijn regelmatige functionels van F , waarbij $\Phi^{(\gamma)}(x_1, \dots, x_{m(\gamma)})$ symmetrisch in de vectoren $x_1, \dots, x_{m(\gamma)}$ is en n niet bevat.

Als nu $S(x)$ de v.f. van de aselechte steekproef is, en indien de variantie van

$$\theta^{(g)}(\underline{s}) = \frac{1}{n^m} \sum_{\alpha_1=1}^n \dots \sum_{\alpha_m(\gamma)=1}^n \Phi^{(g)}(\underline{x}_{\alpha_1}, \dots, \underline{x}_{\alpha_m(\gamma)})$$

bestaat, dan nadert de simultane v.f. van

$$\sqrt{n} \{ \theta^{(1)}(\underline{s}) - \theta^{(1)}(F) \}, \dots, \sqrt{n} \{ \theta^{(g)}(\underline{s}) - \theta^{(g)}(F) \}$$

tot de g -dimensionale normale v.f. met gemiddelde 0 en $\{ m(\gamma) m(\delta) \} \gamma^{(g, \delta)}$ als covariantiematrix.

Stelling 1.5:

Stel $(\underline{u}') = (\underline{u}^{(1)'}, \dots, \underline{u}^{(g)'})$ is een stochastische vector, waarbij $\underline{u}^{(g)'}$ in stelling (1.3) gedefinieerd is en veronderstel dat aan de voorwaarden van stelling (1.3) voldaan is. Indien de functie $h(y) = h(y^{(1)}, \dots, y^{(g)})$, niet bevat en evenals de partiële afgeleiden van de 2^e orde in de omgeving van het punt $(y) = (\theta) = (\theta^{(1)}, \dots, \theta^{(g)})$ continu is, dan nadert de verdeling van de stochastische veranderlijke $\sqrt{n} \{ h(\underline{u}') - h(\theta) \}$ tot de normale verdeling met gemiddelde 0 en

$$\sum_{\gamma=1}^g \sum_{\delta=1}^g m(\gamma) m(\delta) \left(\frac{\partial h(y)}{\partial y^{(\gamma)}} \right)_{y=\theta} \left(\frac{\partial h(y)}{\partial y^{(\delta)}} \right)_{y=\theta} \} \gamma^{(g, \delta)}$$

als variantie.

Opmerking: Stelling (1.5) is dus een uitbreiding van de vorige stellingen en betreft de asymptotische verdeling van een functie van \underline{u} - of \underline{u}' -statistische grootheden.

Aangezien ieder moment waaruit de variantie bestaat in de vorm $\underline{u}' = \theta(\underline{s})$ geschreven kan worden (zie par. 5 A en stelling (1.4)), is stelling (1.5) een uitbreiding van de stelling over een functie van momenten.

De hierbovengenoemde limietstellingen kunnen uitgebreid worden tot het geval waarbij de \underline{x}_{α} 's verschillende verdelingen hebben. Hoefding [5] zet uiteen hoe de uitbreidingen kunnen geschieden.

5. Speciale gevallen van \underline{u} -statistische grootheden.

A. Momenten en functies van momenten.

Stel $S = S(x)$ is de verdelingsfunctie van de steekproef (1). Substitutie van S voor F in (5) geeft

$$(A.1) \quad \theta(s) = \frac{1}{n^m} \sum_{\alpha_1=1}^n \dots \sum_{\alpha_m=1}^n \Phi(x_{\alpha_1}, \dots, x_{\alpha_m})$$

Als $m=1$ is $\theta(s) = \underline{u}$. Als $m=2$, dan is

$$\begin{aligned}
\theta(S) &= \frac{1}{n^2} \sum_{\alpha_1=1}^n \sum_{\alpha_2=1}^n \Phi(x_{\alpha_1}, x_{\alpha_2}) \\
&= \frac{1}{n^2} \sum_{\substack{\alpha_1=1 \\ \alpha_1 \neq \alpha_2}}^n \sum_{\alpha_2=1}^n \Phi(x_{\alpha_1}, x_{\alpha_2}) + \frac{1}{n} \left[\frac{1}{n} \sum_{\alpha=1}^n \Phi(x_{\alpha}, x_{\alpha}) \right] \\
&= \frac{n-1}{n} \mathcal{U} + \frac{1}{n} \left\{ \frac{1}{n} \sum_{\alpha=1}^n \Phi(x_{\alpha}, x_{\alpha}) \right\},
\end{aligned}$$

en $\theta(S)$ is een lineaire functie van \mathcal{U} -statistische grootheden met coëfficiënten die van n afhangen. Dit geldt voor iedere m . In het algemeen is $\theta(S)$ echter geen zuivere schatting van $\theta(F)$. Indien de verwachting van $\theta(S)$ echter bestaat voor iedere F in \mathcal{D} , dan is

$$E\{\theta(S)\} = \theta(F) + O(n^{-1}) \quad (\text{zie ook (4)})$$

en dus kan de schatting $\theta(S)$ van $\theta(F)$ asymptotisch zuiver over \mathcal{D} genoemd worden.

De steekproefmomenten hebben de vorm (A.1). Hieruit volgt dus dat de zuivere schattingen van de momenten \mathcal{U} -statistische grootheden zijn.

Uit de stellingen (1.1), (1.2) en (1.4) volgt dat de steekproefmomenten asymptotisch normaal verdeeld zijn en uit stelling (1.5) volgt hetzelfde voor een functie van die momenten indien aan de vereiste voorwaarden voldaan is. Dit resultaat is bekend (zie Cramér: Math. Methods of Statistics).

De steekproefmomenten hebben de vorm (A.1). Hun kern Φ kan als volgt verkregen worden:

Momenten t.o.v. de oorsprong:

$$\mu'_{\nu_1, \dots, \nu_r} = \int \dots \int (x^{(1)})^{\nu_1} \dots (x^{(r)})^{\nu_r} dF(x^{(1)}, \dots, x^{(r)})$$

Momenten t.o.v. het gemiddelde: Een moment t.o.v. het gemiddelde is een veelterm in momenten μ' t.o.v. de oorsprong. Voor $\nu=1$ (d.w.z. in het geval van een ééndimensionale verdelingsfunctie) is

$$\sigma^2 = \int (x - \mu)^2 dF(x)$$

waarbij μ het gemiddelde van de oorspronkelijke verzameling is. Geef met $\hat{\sigma}^2$ de geschatte waarde uit de steekproef van σ^2 aan. Dan is

$$\hat{\sigma}^2 = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2 \quad \text{waarbij} \quad \bar{x} = \frac{1}{n} \sum_{i=1}^n x_i$$

Nu is $\mathcal{E} \left\{ \frac{n}{n-1} \hat{\sigma}^2 \right\} = \sigma^2$ zoals bekend is.

Het is gemakkelijk om \mathcal{U} -statistische grootheden te verkrijgen, die parameter vrije zuivere schattingen van de momenten t.o.v. het gemiddelde van enige orde zijn (zie Halmos [3])¹). $\hat{\Phi}(x_1, \dots, x_n)$ wordt een parameter vrije zuivere schatting van $\theta(F)$ genoemd indien $\mathcal{E}(\hat{\Phi}(x_1, \dots, x_n)) = \theta(F)$ wanneer $\hat{\Phi}(x_1, \dots, x_n)$ ook al bestaan, waarbij x_1, \dots, x_n n onafhankelijke stochastische vectoren zijn.

Ook de \mathcal{K} -statistische grootheden van Fisher zijn \mathcal{U} -statistische grootheden. Dit volgt ook uit hun definitie als zuivere schattingen van de cumulanten, symmetrisch in de steekproefwaarden (Kendall I, blz. 256). Zij zijn dus ook asymptotisch normaal verdeeld.

B. Gemiddeld verschil en concentratie-coëfficiënt.

Als y_1, \dots, y_n n onafhankelijke, reële, aselechte veranderlijken zijn, dan wordt het gemiddelde verschil van Gini (zonder herhaling) gedefinieerd door

$$\underline{d} = \frac{1}{n(n-1)} \sum_{\alpha \neq \beta} |y_\alpha - y_\beta|$$

Indien de y_α 's allemaal dezelfde verdeling F hebben, is het gemiddelde van \underline{d}

$$\delta = \mathcal{E}(\underline{d}) = \iint |y_1 - y_2| dF(y_1) dF(y_2)$$

1) Als de steekproefmomenten van de m^e orde gedefinieerd worden door

$$g_m(x_1, \dots, x_n) = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^m$$

waar $\bar{x} = \frac{1}{n} \sum_i x_i$, dan bewijst Halmos dat, als $f_m(x_1, \dots, x_n)$ zuivere schattingen van de momenten t.o.v. het gemiddelde van de oorspronkelijke verzameling zijn, deze gegeven worden door

$$f_1(x_1, \dots, x_n) = g_1(x_1, \dots, x_n)$$

$$f_2(x_1, \dots, x_n) = \frac{n}{n-1} g_2(x_1, \dots, x_n)$$

$$f_3(x_1, \dots, x_n) = \frac{n^2}{(n-1)(n-2)} g_3(x_1, \dots, x_n)$$

en voor $m > 3$ kan f_m nog in termen van de g_m 's uitgedrukt worden, maar niet meer als een constant veelvoud van de g_m 's. In het algemeen is f_m een lineaire combinatie van g_1, \dots, g_m met rationale coëfficiënten, waarvan de noemers $(n-1)(n-2) \dots (n-m+1)$ zijn.

en de variantie van \underline{d} is

$$\sigma^2(\underline{d}) = \binom{n}{m}^{-1} \sum_{l=1}^m \binom{m}{l} \binom{n-m}{m-l} \zeta_l \quad \text{waarbij } m=2 \text{ is,}$$

$$= \binom{n}{2}^{-1} \{ 2(n-2) \zeta_1(\delta) + \zeta_2(\delta) \}$$

waar

$$(B.1) \quad \zeta_1(\delta) = \mathcal{E} \{ \psi_1^2(\underline{y}, \delta) \}$$

$$= \mathcal{E} \{ \phi_1^2(\underline{y}, \delta) \} - \delta^2$$

$$= \int \left\{ \int |y_1 - y_2| dF(y_2) \right\}^2 dF(y_1) - \delta^2$$

en

$$(B.2) \quad \zeta_2(\delta) = \int \int (|y_1 - y_2|)^2 dF(y_1) dF(y_2) - \delta^2$$

$$= \int \int (y_1 - y_2)^2 dF(y_1) dF(y_2) - \delta^2 = 2\sigma^2(\underline{y}) - \delta^2$$

De notatie $\zeta_1(\delta)$, $\zeta_2(\delta)$ dient om het verband tussen de functionels van F en de functionel $\delta(F)$ aan te duiden; δ is alleen het symbool van de functionel en niet een bepaalde waarde daarvan. Soortgelijk is $\phi_1(y_1, y_2, \delta) = |y_1 - y_2|$.

Nair [10] heeft $\sigma^2(\underline{d})$ voor verscheidene speciale verdelingen berekend.

Volgens stelling (1.1) is $\sqrt{n}(\underline{d} - \delta)$ asymptotisch normaal indien $\zeta_2(\delta)$ bestaat.

De concentratiecoëfficiënt van Gini wordt voor niet-negatieve waarden van y_1, \dots, y_n gedefinieerd door

$$G = \frac{\underline{d}}{2\bar{y}}$$

waarbij $\bar{y} = \frac{\sum y_\alpha}{n}$. G is hier een functie van twee \underline{U} -statistische grootheden. Indien de y_α 's gelijk verdeeld zijn, indien $\mathcal{E} \{ y^2 \}$ bestaat en indien $\mu = \mathcal{E} \{ y \} > 0$, dan volgt uit stelling (1.5) dat $\sqrt{n}(G - \frac{\delta}{2\mu})$ asymptotisch normaal is met gemiddelde 0 en variantie

$$\sum_{\gamma=\delta, \mu} \sum_{\lambda=\delta, \mu} m(\gamma)m(\lambda) \left(\frac{\partial(\frac{\delta}{2\bar{y}})}{\partial \gamma} \right)_{\substack{\underline{d}=\delta \\ \bar{y}=\mu}} \left(\frac{\partial(\frac{\delta}{2\bar{y}})}{\partial \lambda} \right)_{\substack{\underline{d}=\delta \\ \bar{y}=\mu}} \zeta_1(\gamma, \lambda)$$

waarbij $m(\delta) = 2$ en $m(\mu) = 1$,

$$= \frac{1}{\mu^2} \zeta_1(\delta) - \frac{\delta}{\mu^3} \zeta_1(\delta, \mu) + \frac{\delta^2}{4\mu^4} \zeta_1(\mu)$$

waarbij $\int_1(\delta)$ gegeven wordt door (B.1) en

$$\int_1(\mu) = \int y^2 dF(y) - \mu^2 = \sigma^2(y)$$

$$\int_1(\mu, \delta) = \iint y_1 |y_1 - y_2| dF(y_1) dF(y_2) - \mu \delta.$$

C. Functies van rangen en tekens van verschillen tussen veranderlijken.

Definieer $S(u)$ door

$$(C.1) \quad S(u) = \begin{cases} -1 & < \\ 0 & \text{als } u = 0 \\ 1 & > \end{cases}$$

en stel

$$(C.2) \quad C(u) = \frac{1}{2} \{1 + S(u)\} = \begin{cases} 0 & < \\ \frac{1}{2} & \text{als } u = 0 \\ 1 & > \end{cases}$$

Als

$$x_\alpha = (x_\alpha^{(1)}, \dots, x_\alpha^{(n)}) \quad (\alpha = 1, \dots, n)$$

een steekproef van n vectoren met r componenten elk is, definiëren wij de rang $R_\alpha^{(i)}$ van $x_\alpha^{(i)}$ door

$$(C.3) \quad \begin{aligned} R_\alpha^{(i)} &= \frac{1}{2} + \sum_{\beta=1}^n C(x_\alpha^{(i)} - x_\beta^{(i)}) \quad (i = 1, \dots, r) \\ &= \frac{n+1}{2} + \frac{1}{2} \sum_{\beta=1}^n S(x_\alpha^{(i)} - x_\beta^{(i)}) \end{aligned}$$

Indien de getallen $x_1^{(i)}, x_2^{(i)}, \dots, x_n^{(i)}$ allemaal verschillend zijn, heeft het kleinste de rang 1, of op één na het kleinste de rang 2, enz. Als sommige gelijk zijn, wordt de gemiddelde rang door (C.3) bepaald.

Iedere functie van de rang is een functie van uitdrukkingen $C(x_\alpha^{(i)} - x_\beta^{(i)})$ of $S(x_\alpha^{(i)} - x_\beta^{(i)})$.

Omgekeerd, aangezien $S(x_\alpha^{(i)} - x_\beta^{(i)}) = S(R_\alpha^{(i)} - R_\beta^{(i)})$, is iedere functie van uitdrukkingen $S(x_\alpha^{(i)} - x_\beta^{(i)})$ of $C(x_\alpha^{(i)} - x_\beta^{(i)})$ een functie van de rang.

Beschouw een regelmatige functionel $\theta(F)$ waarvan de kern $\phi(x_1, \dots, x_m)$ alleen afhangt van de tekens van de verschillen tussen de veranderlijken

$$(C.4) \quad S(x_\alpha^{(i)} - x_\beta^{(i)}) \quad (\alpha, \beta = 1, \dots, m; i = 1, \dots, r)$$

De overeenkomstige \mathcal{L} -statistische grootheid is een functie van de rang van de steekproefveranderlijken.

De functie Φ kan alleen een eindig aantal waarden c_1, \dots, c_N aannemen. Als $\pi_i = P\{\Phi = c_i\}$, ($i = 1, \dots, N$), krijgen wij

$$\theta = c_1 \pi_1 + \dots + c_N \pi_N, \quad \sum_{i=1}^N \pi_i = 1.$$

π_i is een regelmatig functioneel waarvan de kern $\Phi_i(x_1, \dots, x_m)$ gelijk is aan 1 of 0 al naar gelang $\Phi = c_i$ of $\neq c_i$. Dus is

$$\Phi = c_1 \Phi_1 + \dots + c_N \Phi_N.$$

Wil $\theta(F)$ bestaan, dan moet c_i eindig zijn en dus is Φ begrensd. Dus bestaat $\mathcal{E}\{\Phi^2\}$ en als x_1, x_2, \dots, x_n gelijk verdeeld zijn, nadert de v.f. van $\sqrt{n}(\mathcal{L} - \theta)$ volgens stelling (1.1) tot een normale v.f. die niet-singulier is voor $\gamma_i > 0$.

Een voorbeeld van een zodanig functioneel is

D. De correlatie tussen tekens van verschillen.

Beschouw de tweedimensionale steekproef

$$(D.1) \quad (x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)$$

$$\left[\text{of } (x_1^{(1)}, x_1^{(2)}), (x_2^{(1)}, x_2^{(2)}), \dots, (x_n^{(1)}, x_n^{(2)}) \right].$$

Met elke twee waarnemingen van deze steekproef komen er een paar tekens van verschillen tussen de desbetreffende veranderlijken overeen:

$$(D.2) \quad S(x_\alpha - x_\beta), S(y_\alpha - y_\beta), \quad (\alpha \neq \beta; \alpha, \beta = 1, \dots, n).$$

(D.2) is een verzameling van $n(n-1)$ paren van tekens van verschillen.

Omdat

$$\sum_{\alpha \neq \beta} S(x_\alpha - x_\beta) = 0 = \sum_{\alpha \neq \beta} S(y_\alpha - y_\beta)$$

is de covariantie t van de tekens van de verschillen (D.2)

$$(D.3) \quad t = \frac{1}{n(n-1)} \sum_{\alpha \neq \beta} S(x_\alpha - x_\beta) S(y_\alpha - y_\beta)$$

of $t = \frac{S}{n(n-1)}$, waarbij S de S van Kendall [7] is.

We noemen t de verschiltekencovariantie van de steekproef (D.1).

Als alle x 's en alle y 's verschillend zijn, dan is

$$\sum_{\alpha \neq \beta} S^2(x_\alpha - x_\beta) = n(n-1) = \sum_{\alpha \neq \beta} S^2(y_\alpha - y_\beta)$$

en is t de productmomentcorrelatie van de tekens van de verschillen. t is een lineaire functie van het (minimum) aantal

inversies in de permutatie van de rang van x en y (Kendall [7], blz. 8).

\underline{t} is een \mathcal{U} -statistische grootheid. \underline{t} is een functie van een aselechte steekproef uit een tweedimensionale verzameling en is dus een zuivere schatting van de regelmatige functioneel van de tweede graad

$$(D.4) \quad \tau = \iiint s(x_1 - x_2) s(y_1 - y_2) dF(x_1, y_1) dF(x_2, y_2)$$

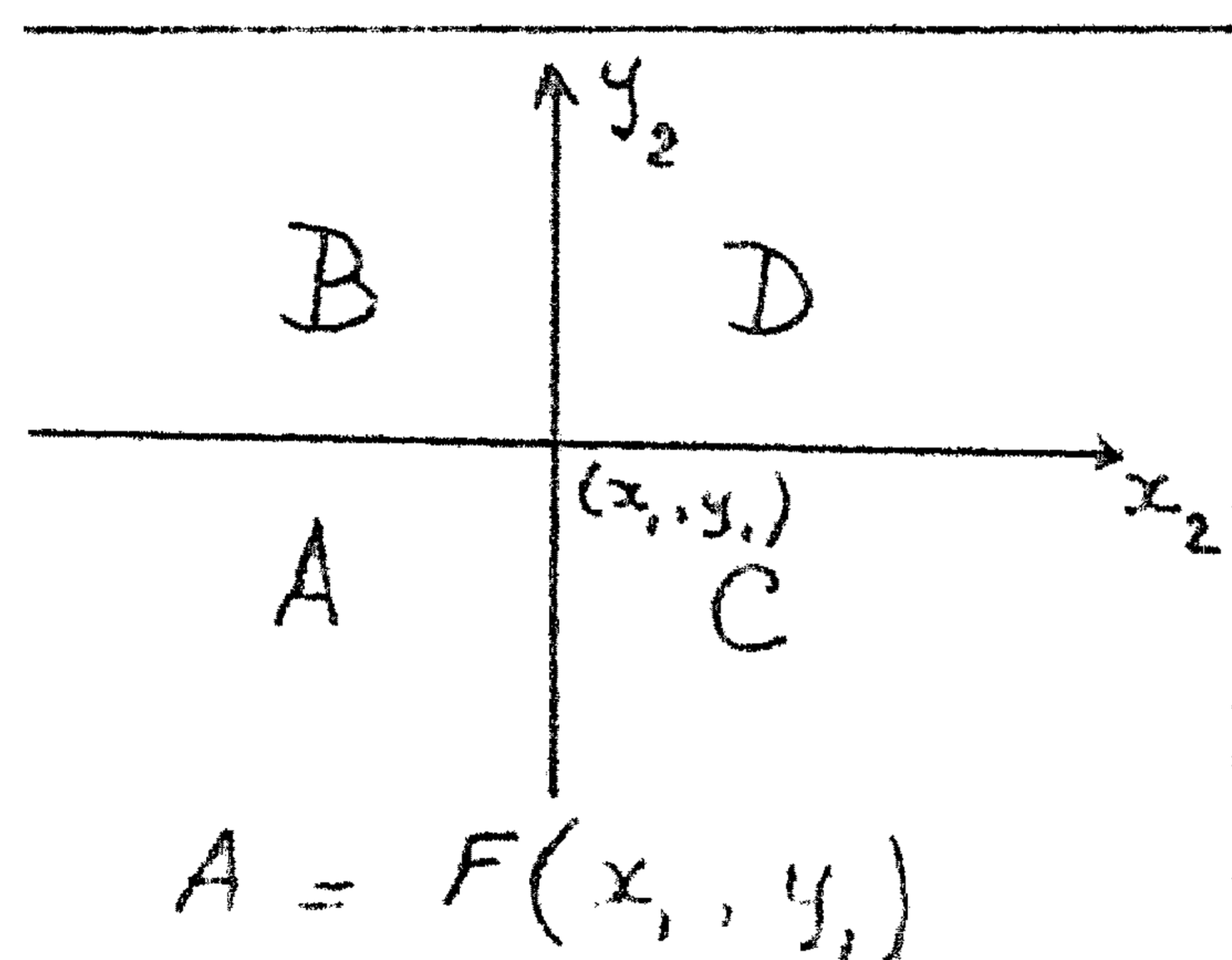
τ is de covariantie van de tekens van de verschillen van de overeenkomstige componenten van (x_1, y_1) en (x_2, y_2) in de verzameling van paren van onafhankelijke vectoren $(x_1, y_1), (x_2, y_2)$ met gelijke v.f. $F(x, y)$. Als $F(x, y)$ continu is, is τ de productmomentcorrelatie van de tekens van de verschillen.

Stel nu

$$(D.5) \quad \bar{F}(x, y) = \frac{1}{4} \{F(x-0, y-0) + F(x-0, y+0) + F(x+0, y-0) + F(x+0, y+0)\},$$

dan is

$$(D.6) \quad \Phi_1(x, y | \tau) = \mathcal{E} \{ \Phi_1(x_1, y_1; x_2, y_2) \}$$



of $F(x_1, -0, y_1, -0)$ in het discontinue geval.

$$A + B = F(x_1, \infty)$$

$$(A + C) = F(\infty, y_1)$$

$$A + B + C + D = 1.$$

$$= \iint s(x_1 - x_2) s(y_1 - y_2) dF(x_2, y_2)$$

$$= \int_{\substack{x_2 < x_1 \\ y_2 < y_1}} dF - \int_{\substack{x_2 > x_1 \\ y_2 < y_1}} dF - \int_{\substack{x_2 < x_1 \\ y_2 > y_1}} dF + \int_{\substack{x_2 > x_1 \\ y_2 > y_1}} dF$$

$$= F(x_1 - 0, y_1 - 0) - [F(\infty, y_1 - 0) - F(x_1 + 0, y_1 - 0)] +$$

$$- [F(x_1 - 0, \infty) - F(x_1 - 0, y_1 + 0)] +$$

$$+ [1 - F(\infty, y_1 + 0) - F(x_1 + 0, \infty) + F(x_1 + 0, y_1 + 0)]$$

$$= 1 - 2\bar{F}(x_1, \infty) - 2\bar{F}(\infty, y_1) + 4\bar{F}(x_1, y_1)$$

(uit D.5)

Verder is
$$\tau = \mathcal{E} \{ \Phi_1(x_1, y_1 | \tau) \}$$

Nu is de variantie van \underline{t}

$$(D.7) \quad \sigma^2(\underline{t}) = \binom{n}{2}^{-1} \sum_{c=1}^2 \binom{2}{c} \binom{n-2}{2-c} \zeta_c(\tau) = \frac{2}{n(n-1)} [2(n-2)\zeta_1(\tau) + \zeta_2(\tau)]$$

waarbij

$$(D.8) \quad \zeta_c(\tau) = \mathcal{E} \{ \Phi_1^2(x_1, y_1 | \tau) \} - \tau^2$$

$$(D.9) \quad \gamma_2(\tau) = \mathcal{E} \left\{ S^2(x_1 - x_2) S^2(y_1 - y_2) \right\} - \tau^2$$

Als $F(x, y)$ continu is, dan is

$$\gamma_2(\tau) = 1 - \tau^2$$

en $\overline{F}(x, y)$ in (D.6) kan vervangen worden door $F(x, y)$.

Als \underline{x} en \underline{y} onafhankelijk zijn en een continue verdelingsfunctie hebben, dan is ($\tau=0$)

$$\begin{aligned} \gamma_1(\tau) &= \gamma_1 = \mathcal{E} \left\{ \Phi_1^2(\underline{x}, \underline{y}) \right\} \\ &= \iiint \left\{ \iint S(x_1 - x_2) S(y_1 - y_2) dF(x_2) dF(y_2) \right\}^2 dF(x_1) dF(y_1) \\ &= \iiint [2F(x_1) - 1]^2 [2F(y_1) - 1]^2 dF(x_1) dF(y_1) \end{aligned}$$

zoals volgt uit (D.6) voor het onafhankelijke en continue geval

$$\text{Nu is} \quad \int (2u - 1)^2 du = \left[\frac{2}{3} u^3 - \frac{2}{2} u^2 + u \right]_0^1 = \frac{2}{3} - 2 + 1 = \frac{1}{3}$$

$$\therefore \gamma_1 = \frac{1}{9}$$

Ook is

$$\gamma_2(\tau) = \gamma_2 = 1$$

Dus

$$(D.10) \quad \sigma^2(\underline{t}) = \frac{2}{n(n-1)} \left[\frac{2(n-2)}{9} + 1 \right] = \frac{2(2n+5)}{9n(n-1)}$$

Dit stemt overeen met Kendall [7], waar $\sigma^2(\underline{s}) = \frac{1}{15} n(n-2)(2n+5)$ en $\underline{t} = \frac{2\underline{S}}{n(n-1)}$.

In het geval (D.10) is de verdeling van \underline{t} onafhankelijk van de ééndimensionale verdelingen van \underline{x} en \underline{y} . Dit is echter niet meer het geval als de onafhankelijke veranderlijken discontinu zijn. Dan hangt $\sigma^2(\underline{t})$ af van $\mathcal{P}\{\underline{x}_1 = \underline{x}_2\}$, $\mathcal{P}\{\underline{y}_1 = \underline{y}_2\}$ en $\mathcal{P}\{\underline{x}_1 = \underline{x}_2 = \underline{x}_3\}$, $\mathcal{P}\{\underline{y}_1 = \underline{y}_2 = \underline{y}_3\}$.

Uit stelling (1.1) volgt dat de v.f. van $\sqrt{n}(\underline{t} - \tau)$ naar de normale vorm nadert. Vergelijk ook Kendall [11] voor het speciale geval dat alle permutaties van de rang van \underline{x} en \underline{y} even waarschijnlijk zijn, hetgeen overeenkomt met de onafhankelijkheid van de continue (stochastische) aselechte veranderlijken \underline{x} en \underline{y} .

In het algemene geval is de asymptotische normaliteit van \underline{t} bewezen door Daniels en Kendall [2] en door Hoefding [4].

De functionel $\tau(F)$ is stationnair (en dus is de normale limietverdeling van $\sqrt{n}(\underline{t} - \tau)$ dan singulier) als $\gamma_1 < 0$ is. In het continue geval betekent dit dat $\Phi_1(\underline{x}, \underline{y} | \tau) = \tau$ of dat aan

$$(D.11) \quad 4F(\underline{x}, \underline{y}) = 2F(\underline{x}, \infty) + 2F(\infty, \underline{y}) - 1 + \tau.$$

wordt voldaan met een waarschijnlijkheid 1. Dit is b.v. het geval wanneer y een toenemende functie van x is. Dan is $\underline{t} = \bar{t} = 1$ met waarschijnlijkheid 1 en $\sigma^2(\underline{t}) = 0$.

Een voorbeeld van een geval waar aan (D.11) voldaan is en $\sigma^2(\underline{t}) > 0$ is, is het volgende: x is homogeen verdeeld over het interval $(0,1)$ en

$$(D.12) \quad \underline{y} = \underline{x} + \frac{1}{2} \quad \text{voor } 0 \leq \underline{x} < \frac{1}{2}, \quad \underline{y} = \underline{x} - \frac{1}{2} \quad \text{voor } \frac{1}{2} \leq \underline{x} \leq 1.$$

Hier is $\bar{t} = 0$, $\xi_2 = 1$ en $\sigma^2(\underline{t}) = \frac{2}{n(n-1)}$ uit (D.7) met $\xi_1 = 0$.

E De toets van Mann tegen "verloop" ("trend").

Stel y_1, \dots, y_n zijn n onafhankelijke, reële, stochastische veranderlijken, waarbij y_α de continue v.f. $F_\alpha(y)$, ($\alpha = 1, \dots, n$) heeft. Wij willen de hypothese H_0 toetsen dat de y_α 's ($\alpha = 1, \dots, n$) alle dezelfde verdeling hebben, dus

$$H_0 : F_1(y) = \dots = F_n(y)$$

tegen de alternatieve hypothese H_1 van een "stijgend verloop", dus

$$H_1 : F_1(y) < F_2(y) < \dots < F_n(y).$$

Mann [9] heeft een toets van H_0 tegenover H_1 voorgesteld die gebaseerd is op het aantal T ongelijkheden $y_\alpha < y_\beta$, waarbij $\alpha < \beta$ is. Nu is

$$\begin{aligned} 2T - \frac{n(n-1)}{2} &= \sum_{\alpha < \beta} s(y_\beta - y_\alpha) && \text{(Kendall [7], verg. (17))} \\ &= \sum_{\alpha < \beta} s(\alpha - \beta) s(y_\alpha - y_\beta). \end{aligned}$$

De \underline{t} -statistische grootheid

$$\underline{t} = \frac{4T}{n(n-1)} - 1 = \frac{2}{n(n-1)} \sum_{\alpha < \beta} s(\alpha - \beta) s(y_\alpha - y_\beta)$$

is dezelfde als (D.3) voor het speciale geval waar één component geen stochastische veranderlijke is.

Stel

$$\begin{aligned} \bar{T}_{\alpha\beta} &= s(\alpha - \beta) \iint s(y_1 - y_2) dF_\alpha(y_1) dF_\beta(y_2) \\ &= s(\alpha - \beta) \left\{ 2 \int F_\beta(y) dF_\alpha(y) - 1 \right\}. \end{aligned}$$

Indien H_0 geldt, is $\bar{T}_{\alpha\beta} = s(\alpha - \beta) \left\{ 2 \int F(y) dF(y) - 1 \right\} = 0$

Indien H_1 geldt, is $\bar{T}_{\alpha\beta} < 0$, want $s(\alpha - \beta) = -1$ en $\left\{ \dots \right\} > 0$.

Aangezien $\xi(\underline{t}) = \bar{T}_n = \frac{2}{n(n-1)} \sum_{\alpha < \beta} \bar{T}_{\alpha\beta}$

volgt dat

$$(E.1) \quad \xi(\underline{t}) = 0 \quad \text{onder } H_0 \quad \text{en} \quad \xi(\underline{t}) < 0 \quad \text{onder } H_1.$$

Veronderstel $P(\underline{t} < a_n | H_0) \leq \varepsilon$, waarbij a_n het grootste getal is dat hieraan voldoet, dan wordt ε de onbetrouwbaarheidsdrempel genoemd.

Dit betekent dus dat de waarschijnlijkheid dat een aselechte steekproef een waarde van \underline{t} zal geven kleiner dan a_n , indien H_0 juist is, ten hoogste gelijk is aan ε , de onbetrouwbaarheidsdrempel. ε is dus de kans op het ten onrechte verwerpen van de juiste hypothese (fout van de eerste soort). Maar wij kunnen ook een fout van de tweede soort maken, d.w.z. H_0 wordt niet verworpen, terwijl een alternatieve hypothese $H_1 \neq H_0$ geldt. Dit wordt gegeven door

$$(E.2) \quad P[\underline{t} \geq a_n | H_1] = \beta \quad \text{waarbij } H_1 \neq H_0 \text{ geldt.}$$

Het is uiteraard ook gewenst dat β zo klein mogelijk zal zijn. Het onderscheidingsvermogen van de toets van de hypothese H_0 tegen de alternatieve hypothese H_1 wordt gedefinieerd door

$$(E.3) \quad P[\underline{t} < a_n | H_1] = 1 - \beta,$$

dat is de waarschijnlijkheid dat H_0 terecht verworpen zal worden wanneer $H_1 \neq H_0$ geldt.

Uit de definitie van a_n volgt dat $a_n \rightarrow 0$ voor $n \rightarrow \infty$ en omdat $\sigma^2(\underline{t}) = O(\frac{1}{n})$ (want \underline{t} is een \mathcal{U} -statistische grootheid), volgt uit de ongelijkheid van Tchebycheff (Cramér [1], blz. 182) dat $P[\underline{t} < a_n | H_1] \rightarrow 1$, d.w.z. dat de toets bruikbaar is en dientengevolge asymptotisch zuiver is.

$$\begin{aligned} \text{(Let wel dat onder } H_0: \tau_n = \mathcal{E}(\underline{t}) = 0 \\ \text{en onder } H_1: \tau_n = \mathcal{E}(\underline{t}) < 0 \text{ .)} \end{aligned}$$

Volgens de ongelijkheid van Tchebycheff is

$$P\{|\underline{t} - \tau_n| \geq k\} \leq \frac{1}{k^2} \quad \text{waarbij } k > 0$$

$$\text{of } P\{|\underline{t} - \tau_n| \geq k\} \leq \frac{\sigma^2}{k^2} = \frac{\sigma^2(\underline{t})}{k^2}$$

$$\text{d.i. } P\{(\underline{t} - \tau_n)^2 \geq k^2\} \leq \frac{\sigma^2(\underline{t})}{k^2}$$

$$\therefore \beta = P\{\underline{t} \geq a_n | H_1\} = P\{\underline{t} - \tau_n \geq a_n - \tau_n | H_1\} \\ \leq \frac{\sigma^2(\underline{t})}{(a_n - \tau_n)^2} = \frac{O(\frac{1}{n})}{(a_n - \tau_n)^2}$$

Hierin is $\tau_n \neq 0$ en $a_n \rightarrow 0$ voor $n \rightarrow \infty$, $\therefore \beta \rightarrow 0$, waaruit volgt dat $P\{\underline{t} < a_n | H_1\} \rightarrow 1$.)

Mann heeft ook aangetoond dat de toets asymptotisch zuiver is en heeft bovendien voldoende voorwaarden gegeven waaronder de toets zuiver is voor eindige n .

Uit de algemene theorie van de \mathcal{U} -statistische grootheden voor het geval waar de x_α 's verschillende v.f.'s hebben, volgt dat de verdeling van $\frac{\underline{t} - \tau_n}{\sigma(\underline{t})}$ asymptotisch normaal is

$$\left(\begin{array}{l} \tau_n < 0 \text{ onder } H_1 \\ \text{en } |\tau_n| > |a_n| \\ \text{voor voldoende} \\ \text{grote } n \\ \therefore a_n - \tau_n > 0 \end{array} \right)$$

(Hoofding [5], stellingen 3.1 en 3.2).

F. Parameter vrije toets voor onafhankelijkheid.

Veronderstel dat de stochastische veranderlijken x, y een continue simultane v.f. $F(x, y)$ hebben. Wij willen de hypothese H_0 toetsen dat x en y onafhankelijk zijn, d.w.z. dat

$$F(x, y) = F(x, \infty) \cdot F(\infty, y).$$

Omdat de verdeling van iedere statistische grootte, die alleen de rang van de veranderlijken bevat onder de hypothese H_0 niet van de v.f. van de verzameling afhangt, is gesuggereerd verscheidene "order statistics", o.a. ook de verschiltekencorrelatie \underline{t} , te gebruiken om voor onafhankelijkheid te toetsen

Uit de voorgaande resultaten kunnen wij het asymptotisch onderscheidingsvermogen van de toets voor onafhankelijkheid, dat op \underline{t} gebaseerd is, verkrijgen. Indien H_0 geldt, is $E(\underline{t}) = \bar{t} = 0$ (vgl. par. D) en de kritieke zone van grootte ε van de \underline{t} -toets kan gedefinieerd worden door $|\underline{t}| > c_n$, waarbij c_n het kleinste getal is dat voldoet aan de ongelijkheid

$$(F.1) \quad P\{|\underline{t}| > c_n | H_0\} \leq \varepsilon$$

ε wordt de onbetrouwbaarheidsdrempel genoemd. Volgens stelling (1.2) en (D.10) kunnen wij schrijven $c_n = \frac{2\lambda_n}{3\sqrt{n}}$, waarbij λ_n tot een positieve constante λ nadert, die van ε afhangt.

Aangezien $\sigma^2(\underline{t}) = O(\frac{1}{n})$, nadert het onderscheidingsvermogen

$$P\left\{|\underline{t}| \geq \frac{2\lambda_n}{3\sqrt{n}} | H\right\}$$

tot 1 als $n \rightarrow \infty$ in de alternatieve hypothese H met $\tau(F) \neq 0$.

Als $\tau = 0$, is $\lim_{n \rightarrow \infty} P\left\{|\underline{t}| \geq \frac{2\lambda_n}{3\sqrt{n}} | H\right\} < 1$ (uit ongelijkheid van Tchebycheff)

Als $\tau = 0$ en $\int (\bar{t}) < \frac{1}{9}$, krijgen wij zelfs

$$\lim_{n \rightarrow \infty} P\left\{|\underline{t}| \geq \frac{2\lambda_n}{3\sqrt{n}} | H\right\} < \varepsilon$$

en ten opzichte van deze alternatieven is de toets asymptotisch onzuiver. In het geval van de verdeling (D.12) geldt b.v. dat

$$P\left\{|\underline{t}| \geq \frac{2\lambda_n}{3\sqrt{n}} | H\right\} \rightarrow 0.$$

In dit geval is er een functionele betrekking tussen de veranderlijken, en de verdeling verschilt dus enorm van het geval van onafhankelijkheid.

De vraag is nu of er parameter vrije toetsingsmethoden voor onafhankelijkheid bestaan die zuiver zijn of asymptotisch zuiver zijn.

Verdere analyse van parameter vrije toetsingsmethoden in het geval van onafhankelijkheid.

In een parameter vrije toetsingsmethode van een statistische hypothese worden geen onderstellingen omtrent de functionele vorm van de verdeling van de verzameling gemaakt. Een algemene theorie van parameter vrije toetsingsmethoden is nog niet ontwikkeld en het ziet er naar uit of er geen bevredigende definitie van een "beste" parameter vrije toetsingsmethode verkregen kan worden. Wenselijke eigenschappen van een "goede" parameter vrije toetsingsmethode zijn zuiverheid en bruikbaarheid. Een toets van een hypothese H_0 wordt bruikbaar ten opzichte van een bepaalde klasse toelaatbare hypothesen genoemd, indien de waarschijnlijkheid van het niet verwerpen van H_0 wanneer een hypothese $\neq H_0$ van die klasse geldt, tot nul nadert met toenemende grootte van de steekproef.

Een aselechte steekproef van grootte n is gegeven en daaruit willen wij bepalen of de twee stochastische veranderlijken \underline{x} en \underline{y} onafhankelijk zijn. Wij nemen aan dat de v.f. $F(x, y)$ van $(\underline{x}, \underline{y})$ continu is. Stel Ω' is de klasse van continue v.f.'s $F(x, y)$ en Ω'' is de klasse van v.f.'s die continue simultane en marginale waarschijnlijkheden hebben,

$$f(x, y) = \frac{\partial^2 F(x, y)}{\partial x \partial y} \quad , \quad f_1(x) = \int f(x, y) dy \quad , \quad f_2(y) = \int f(x, y) dx$$

$H_0: F(x, y) = F(x, \infty) \cdot F(\infty, y)$ is de hypothese die getoetst moet worden.

$F(x, y)$ is dus een tweedimensionale v.f. en wij schrijven

$$D(x, y) = F(x, y) - F(x, \infty) \cdot F(\infty, y)$$

en

(F.2) $\Delta = \Delta(F) = \int D^2(x, y) dF(x, y)$, waarbij de integraal zich over R_2 (zie Cramér [1]) uitstrekt.

De stochastische veranderlijken $\underline{x}, \underline{y}$ met de v.f. $F(x, y)$ is dan en slechts dan onafhankelijk als $D(x, y) \equiv 0$.

Stelling (F.1):

Als $F(x, y)$ tot Ω'' behoort, is $\Delta(F) = 0$ dan en slechts dan als $D(x, y) \equiv 0$.

Bewijs:

Als $D(x, y) \equiv 0$, dan is $\Delta(F) = 0$.

Omgekeerd, veronderstel $D(x, y) \not\equiv 0$. Aangezien $F(x, y)$ in Ω'' is, is de functie $f(x, y) - f_1(x) f_2(y)$ continu. Nu is

$$D(x, y) = \int_{-\infty}^x \int_{-\infty}^y d(u, v) du dv \quad .$$

Voor $D(x, y) \neq 0$ is $d(x, y) \neq 0$ en aangezien

$$\begin{aligned} \iint d(x, y) dx dy &= \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} \{f(x, y) - f_1(x) f_2(y)\} dx dy \\ &= \iint f(x, y) dx dy - \int f_1(x) dx \int f_2(y) dy \\ &= 1 - 1 = 0 \end{aligned}$$

bestaat er een rechthoek Q in \mathcal{R}_2 zodanig dat $d(x, y) > 0$ indien (x, y) in Q is. Dientengevolge is $D(x, y)$ bijna overal in $Q \neq 0$ en $f(x, y) > 0$ in Q . Dus is

$$\Delta(F) \geq \iint_Q D^2(x, y) f(x, y) dx dy > 0$$

$\neq 0$, hetgeen de stelling bewijst.

Als $F(x, y)$ discontinu is, kan $\Delta(F) = 0$ zijn terwijl $D(x, y) \neq 0$. Dit is b.v. het geval voor de verdeling

$$P\{\underline{x} = 0, \underline{y} = 1\} = P\{\underline{x} = 1, \underline{y} = 0\} = \frac{1}{2}.$$

Het is nog onbekend of $D(x, y) \equiv 0$ is als $\Delta(F) = 0$ is voor continue of absoluut continue $F(x, y)$.

Uit stelling (F.1) zien wij dus dat alleen wanneer de stochastische veranderlijken $\underline{x}, \underline{y}$ onafhankelijk zijn $\Delta(F) = 0$ is.

Veronderstel nu

$$(F.3) \quad C(u) = \begin{cases} 1 & \text{voor } u \geq 0 \\ 0 & \text{voor } u < 0 \end{cases}$$

$$\Psi(x_1, x_2, x_3) = C(x_1 - x_2) - C(x_1 - x_3)$$

$$\Phi(x_1, y_1; \dots; x_5, y_5) = \frac{1}{4} \Psi(x_1, x_2, x_3) \Psi(x_1, x_4, x_5) \Psi(y_1, y_2, y_3) \Psi(y_1, y_4, y_5)$$

Wij kunnen schrijven

$$(F.4) \quad \Delta = \int \dots \int \Phi(x_1, y_1; \dots; x_5, y_5) dF(x_1, y_1) \dots dF(x_5, y_5).$$

(I) De statistische grootheid D . Stel $(\underline{x}_1, \underline{y}_1), \dots, (\underline{x}_n, \underline{y}_n)$ is een aselechte steekproef uit de verzameling met de v.f. $F(x, y)$. $n \geq 5$ en stel

$$(F.5) \quad \underline{D} = \underline{D}_n = \frac{1}{n(n-1) \dots (n-4)} \sum'' \Phi(\underline{x}_{\alpha_1}, \underline{y}_{\alpha_1}; \dots; \underline{x}_{\alpha_5}, \underline{y}_{\alpha_5})$$

waarbij \sum'' de sommatie over alle α aanduidt zodanig dat

$$\alpha_i = 1, \dots, n; \alpha_i \neq \alpha_j \text{ als } i \neq j \text{ (} i, j = 1, \dots, 5 \text{)}.$$

Omdat het aantal termen in \sum'' $n(n-1) \dots (n-4)$ is, volgt uit (F.4) dat

$$(F.6) \quad \mathcal{E} \underline{D} = \mathcal{E}(\underline{D}) = \Delta$$

In het geval van onafhankelijkheid is $E \underline{D} = 0$ (want $\Delta = 0$ voor onafhankelijkheid zoals volgt uit (F.3) en (F.4)) en dien-
tengevolge kan \underline{D} zowel negatieve als positieve waarden aannemen.
Het zal bewezen worden dat $-\frac{1}{50} \leq \underline{D}_n \leq \frac{1}{50}$, waarbij de bovenste
grens $\frac{1}{50}$ voor iedere n bereikt kan worden, terwijl het mini-
mum van \underline{D}_n , $-\frac{1}{50}$, blijkbaar toeneemt met n .

De stochastische veranderlijke \underline{D} zoals gedefinieerd in
(F.5) behoort tot de klasse van \mathcal{U} -statistische grootheden.
Hieruit volgt dus onmiddellijk dat:

$$\begin{aligned} \Phi(x_1, y_1; \dots; x_5, y_5) &= D_5 = \frac{1}{5!} \sum'' \Phi(x_{\alpha_1}, y_{\alpha_1}; \dots; x_{\alpha_5}, y_{\alpha_5}) \\ \Phi(x_1, y_1; \dots; x_k, y_k) &= \int \dots \int \Phi(x_1, y_1; \dots; x_k, y_k; x_{k+1}, y_{k+1}; \dots; x_5, y_5) \\ &\quad dF(x_{k+1}, y_{k+1}) \dots dF(x_5, y_5), \quad (k=1, \dots, 5) \end{aligned}$$

$$\zeta_k = \int \dots \int \left\{ \Phi(x_1, y_1; \dots; x_k, y_k) - \Delta \right\}^2 dF(x_1, y_1) \dots dF(x_k, y_k)$$

Dan is de variantie van \underline{D}_n

$$(F.7) \quad \text{Var } \underline{D}_n = \binom{n}{5}^{-1} \sum_{k=1}^5 \binom{5}{k} \binom{n-5}{5-k} \zeta_k$$

Verder is $25 \zeta_1 \leq n \text{Var } \underline{D}_n \leq 5 \zeta_5$
 $n \text{Var } \underline{D}_n$ is een afnemende functie van n en

$$(F.8) \quad \lim_{n \rightarrow \infty} n \text{Var } \underline{D}_n = 25 \zeta_1.$$

Volgens stelling (1.1) is de stochastische veranderlijke
 $\sqrt{n}(\underline{D}_n - \Delta)$ asymptotisch normaal met gemiddelde 0 en variantie
 $25 \zeta_1$.

Het zal ook bewezen worden dat in het geval van onafhan-
kelijkheid $\zeta_1 = 0$ is zodat $\sqrt{n} \underline{D}_n$ een ontaarde normale limiet-
verdeling heeft. ($\Delta = 0$ voor het geval van onafhankelijkheid.)
Hoefding heeft ook aangetoond dat $n \underline{D}_n$ een niet-normale limiet-
verdeling heeft, maar aangezien het bewijs daarvan tamelijk ge-
compliceerd is, wordt dit hier niet gegeven.

(II) De berekening van \underline{D} . De rang van \underline{x}_α wordt gegeven door

$$\sum_{\beta=1}^n C(\underline{x}_\alpha - \underline{x}_\beta).$$

Schrijf nu

$$\underline{a}_\alpha = \sum_{\beta=1}^n C(\underline{x}_\alpha - \underline{x}_\beta) - 1$$

$$\underline{b}_\alpha = \sum_{\beta=1}^n C(\underline{y}_\alpha - \underline{y}_\beta) - 1$$

en
$$\underline{c}_\alpha = \sum_{\beta=1}^n C(\underline{x}_\alpha - \underline{x}_\beta) \cdot C(\underline{y}_\alpha - \underline{y}_\beta) - 1.$$

Hier is $a_\alpha + 1$ en $b_\alpha + 1$ dus de rang van resp. x_α en y_α . c_α is het aantal elementen uit de steekproef waarvoor beide $x_\beta < x_\alpha$ en $y_\beta < y_\alpha$. (Omdat $F(x, y)$ continu is, kunnen wij aannemen dat $x_\alpha \neq x_\beta$ en $y_\alpha \neq y_\beta$ als $\alpha \neq \beta$.) Schrijf verder

$$(F.9) \quad \begin{aligned} A &= \sum_{\alpha=1}^n a_\alpha (a_\alpha - 1) b_\alpha (b_\alpha - 1) \\ B &= \sum_{\alpha=1}^n (a_\alpha - 1)(b_\alpha - 1) c_\alpha \\ C &= \sum_{\alpha=1}^n c_\alpha (c_\alpha - 1) \end{aligned}$$

dan is $\sum'' \Phi(x_{\alpha_1}, y_{\alpha_1}; \dots; x_{\alpha_s}, y_{\alpha_s}) =$

$$= \sum'' \left[\left\{ c(x_{\alpha_1} - x_{\alpha_2}) - c(x_{\alpha_1} - x_{\alpha_3}) \right\} \left\{ c(x_{\alpha_1} - x_{\alpha_4}) - c(x_{\alpha_1} - x_{\alpha_5}) \right\} \left\{ c(y_{\alpha_1} - y_{\alpha_2}) + \right. \right. \\ \left. \left. - c(y_{\alpha_1} - y_{\alpha_3}) \right\} \left\{ c(y_{\alpha_1} - y_{\alpha_4}) - c(y_{\alpha_1} - y_{\alpha_5}) \right\} \right]$$

uit (F.3) en uit (F.5) en geeft na rangschikking in verschillende gevallen

$$(F.10) \quad = A - 2(n-2)B + (n-2)(n-3)C.$$

$$\therefore \underline{D} = \frac{A - 2(n-2)B + (n-2)(n-3)C}{n(n-1)(n-2)(n-3)(n-4)}$$

Om dus D te berekenen voor een gegeven steekproef, moeten eerst de getallen a_α , b_α en c_α voor elk element van de steekproef berekend worden, dan A , B , en C uit (F.9) en dit substitueren in (F.10).

(III) De variantie van \underline{D} in het geval van onafhankelijkheid.

Omdat $F(x, y)$ volgens aanname continu is, is ook $F(x, \infty)$ en $F(\infty, y)$ continu. Dan zijn de ongelijkheden $x_1 < x_2$ en $F(x_1, \infty) < F(x_2, \infty)$ equivalent behalve als $F(x_1, \infty) = F(x_2, \infty)$. Hetzelfde geldt voor $y_1 < y_2$ en $F(\infty, y_1) < F(\infty, y_2)$. Dit toont dus aan dat de functie Φ , (F.3), niet van waarde verandert indien x_i, y_i door $F(x_i, \infty), F(\infty, y_i)$ vervangen wordt (zpr 0, behoudens een verzameling met waarschijnlijkheid nul). Δ en D zijn dus invariant voor de transformatie

$$u = F(x, \infty), \quad v = F(\infty, y); \quad \underline{u} = F(x, \infty), \quad \underline{v} = F(\infty, y).$$

In het geval van onafhankelijkheid is dus $F(x, y) = uv$ en

$$\chi_k^2 = \int \dots \int \left\{ \Phi'_k(u_i, v_i; \dots; u_k, v_k) \right\}^2 du_i dv_i \dots du_k dv_k$$

waarbij Φ'_k gedefinieerd wordt door Φ_k met x_i, y_i en $F(x_i, y_i)$ door u_i, v_i en $u_i v_i$ resp. (Voor onafhankelijkheid is $\Delta = 0$.)

Hier is dus

$$\begin{aligned}\Phi'_k &= \int \dots \int \Phi'(u_1, v_1; \dots; u_5, v_5) du_{k+1} dv_{k+1} \dots du_5 dv_5 \\ &= \int \dots \int \frac{1}{5!} \sum'' \Phi(u_{\alpha_1}, v_{\alpha_1}; \dots; u_{\alpha_5}, v_{\alpha_5}) du_{k+1} dv_{k+1} \dots du_5 dv_5\end{aligned}$$

en uit (F 3) is

$$\begin{aligned}&= \int \dots \int \frac{1}{4 \cdot 5!} \sum'' \left\{ c(u_{\alpha_1} - u_{\alpha_2}) - c(u_{\alpha_1} - u_{\alpha_3}) \right\} \left\{ c(u_{\alpha_1} - u_{\alpha_4}) - c(u_{\alpha_1} - u_{\alpha_5}) \right\} \times \\ &\quad \times \left\{ c(v_{\alpha_1} - v_{\alpha_2}) - c(v_{\alpha_1} - v_{\alpha_3}) \right\} \left\{ c(v_{\alpha_1} - v_{\alpha_4}) - c(v_{\alpha_1} - v_{\alpha_5}) \right\} \times \\ &\quad \times du_{k+1} dv_{k+1} \dots du_5 dv_5\end{aligned}$$

$$\begin{aligned}\therefore (4 \cdot 5!)^2 \gamma_k &= \int \dots \int \left\{ \int \dots \int \sum'' \left[c(u_{\alpha_1} - u_{\alpha_2}) - c(u_{\alpha_1} - u_{\alpha_3}) \right] \left[c(u_{\alpha_1} - u_{\alpha_4}) - c(u_{\alpha_1} - u_{\alpha_5}) \right] \times \right. \\ &\quad \times \left. \left[c(v_{\alpha_1} - v_{\alpha_2}) - c(v_{\alpha_1} - v_{\alpha_3}) \right] \left[c(v_{\alpha_1} - v_{\alpha_4}) - c(v_{\alpha_1} - v_{\alpha_5}) \right] du_{k+1} dv_{k+1} \dots du_5 dv_5 \right\}^2 \\ &\quad \times du_1 dv_1 \dots du_k dv_k\end{aligned}$$

De som \sum'' strekt zich uit over alle permutaties $(\alpha_1, \dots, \alpha_5)$ van de getallen $1, \dots, 5$ en $c(z) = \begin{cases} z & \text{voor } z \geq 0 \\ 0 & \text{voor } z < 0 \end{cases}$

Hoefding vindt de volgende waarden voor γ_k :

$$\gamma_1 = 0, \quad 200 \cdot 30^2 \gamma_2 = \frac{2}{9}, \quad 600 \cdot 30^2 \gamma_3 = \frac{14}{3},$$

$$600 \cdot 30^2 \gamma_4 = \frac{164}{9}, \quad 120 \cdot 30^2 \gamma_5 = 12.$$

(Ik kan er niet in slagen deze bepaalde integralen te verifiëren).

Uit (F.7) volgt nu

$$\text{Var } \underline{D}_n = \frac{5!}{n(n-1)(n-2)(n-3)(n-4)} \left\{ 5 \binom{n-5}{4} \gamma_1 + \frac{20}{2} \binom{n-5}{3} \gamma_2 + \frac{60}{6} \binom{n-5}{2} \gamma_3 + \frac{120}{24} (n-5) \gamma_4 + \gamma_5 \right\}$$

$$\therefore \text{Var } (30 \underline{D}_n) = \frac{120 \left\{ \frac{10}{6} (n-5)(n-6)(n-7) \frac{1}{100} + \frac{10}{2} (n-5)(n-6) \frac{14 \cdot 3}{600} + \frac{5(n-5)164}{600} + \frac{9}{10} \right\}}{9 n(n-1)(n-2)(n-3)(n-4)}$$

$$(F.10') = \frac{2n^3 + 6n^2 - 84n + 120}{9n(n-1)(n-2)(n-3)(n-4)} = \frac{2(n^2 + 5n - 32)}{9n(n-1)(n-3)(n-4)}$$

Een andere methode om de coëfficiënten γ_k in het geval van onafhankelijkheid te bepalen is om $\text{Var } \underline{D}_n$ voor $n=5, 6, 7$ uit de exacte verdelingen, als gegeven in de volgende paragraaf, te berekenen en $\lim_{n \rightarrow \infty} n^2 \text{Var } \underline{D}_n$ uit de asymptotische verdeling van $n \underline{D}_n$ (zie Hoefding [6], blz. 552, par. 8)

(IV) De exacte verdeling van \underline{D} in het geval van onafhankelijkheid voor $n = 5, 6, 7$.

Stel $S = \{(x_1, y_1), \dots, (x_n, y_n)\}$ is een steekproef uit een verzameling met een continue verdelingsfunctie. Wij beperken

ons dus tot steekproeven met $x_i \neq x_j$ en $y_i \neq y_j$ als $i \neq j$. Stel $(x'_1, y'_{\beta_1}), \dots, (x'_n, y'_{\beta_n})$ is een herrangschikking van $(x_1, y_1), \dots, (x_n, y_n)$ zodanig dat $x'_1 < x'_2 < \dots < x'_n$ en $y'_1 < y'_2 < \dots < y'_n$. De permutatie $\pi = (\beta_1, \dots, \beta_n)$ van $(1, \dots, n)$ wordt de rangschikking van de steekproef S genoemd.

D_n hangt slechts van de rangschikking van de steekproef af. Wij schrijven dus $D_n = D_n(\pi) = D_n(\beta_1, \dots, \beta_n)$. Als $(\beta'_{\alpha_1}, \dots, \beta'_{\alpha_m})$ een permutatie van $m (< n)$ van de gehele getallen $1, \dots, n$ is zodanig dat $\beta'_1 < \beta'_2 < \dots < \beta'_m$, dan definiëren wij $D_m(\beta'_1, \dots, \beta'_{\alpha_m}) = D_m(\alpha_1, \dots, \alpha_m)$. Als wij in (F.5) (x_{α}, y_{α}) vervangen door (α, β_{α}) krijgen wij

$$(F.11) \quad D_n = D_n(\beta_1, \dots, \beta_n) = \binom{n}{5}^{-1} \sum' D_5(\beta_{\alpha_1}, \dots, \beta_{\alpha_5})$$

waarbij \sum' de sommatie over alle α aanduidt zodanig dat $1 \leq \alpha_1 < \alpha_2 < \dots < \alpha_5 \leq n$.

Als we met $\pi^{(i)}$ de permutatie aanduiden, die verkregen wordt uit $\pi = (\beta_1, \dots, \beta_n)$ door weglating van β_i , dan krijgen wij de herleidingsformule

$$D_n(\pi) = \binom{n}{5}^{-1} \sum' D_5(\beta_{\alpha_1}, \dots, \beta_{\alpha_5}) \quad 1 \leq \alpha_1 < \dots < \alpha_5 \leq n$$

$$D_{n-1}(\pi) = \binom{n-1}{5}^{-1} \sum' D_5(\beta_{\alpha_1}, \dots, \beta_{\alpha_5}) \quad 1 \leq \alpha_1 < \dots < \alpha_5 \leq n-1$$

$$(F.12) \quad \text{d i.} \quad n D_n(\pi) = (n-5) \sum_{i=1}^n D_{n-1}(\pi^{(i)})$$

Uit (F.5) en (F.3) volgt

$$D_5(\beta_1, \dots, \beta_5) = \frac{1}{5!} \sum'' \phi(\alpha_1, \beta_{\alpha_1}; \dots; \alpha_5, \beta_{\alpha_5})$$

$$\text{Dus} \quad 60 D_5(\beta_1, \dots, \beta_5) = \frac{1}{2} \cdot \frac{1}{4} \sum'' \psi(\alpha_1, \alpha_2, \alpha_3) \cdot \psi(\alpha_1, \alpha_4, \alpha_5) \cdot \psi(\beta_{\alpha_1}, \beta_{\alpha_2}, \beta_{\alpha_3}) \cdot \psi(\beta_{\alpha_1}, \beta_{\alpha_4}, \beta_{\alpha_5})$$

waarbij \sum'' uitgestrekt is over alle permutaties van $\alpha_1, \dots, \alpha_5$.

Dus

$$60 D_5(\beta_1, \dots, \beta_5) = \sum'' \frac{1}{8} [c(\alpha_1, -\alpha_2) - c(\alpha_1, -\alpha_3)] [c(\alpha_1, -\alpha_4) - c(\alpha_1, -\alpha_5)] \psi(\beta_{\alpha_1}, \beta_{\alpha_2}, \beta_{\alpha_3}) \cdot \psi(\beta_{\alpha_1}, \beta_{\alpha_4}, \beta_{\alpha_5})$$

$\alpha_1, \dots, \alpha_5$ zijn dus alle mogelijke permutaties van de getallen $1, \dots, 5$. Alleen wanneer $\alpha_1 = 3$ is

$$[c(\alpha_1, -\alpha_2) - c(\alpha_1, -\alpha_3)] [c(\alpha_1, -\alpha_4) - c(\alpha_1, -\alpha_5)] \neq 0$$

en wel met de volgende permutaties:

3, 1, 4, 2, 5 (+1)	3, 2, 4, 1, 5 (+1)	3, 1, 5, 2, 4 (+1)
3, 1, 4, 5, 2 (-1)	3, 2, 4, 5, 1 (-1)	3, 1, 5, 4, 2 (-1)
3, 4, 1, 2, 5 (-1)	3, 4, 2, 1, 5 (-1)	3, 5, 1, 2, 4 (-1)
3, 4, 1, 5, 2 (+1)	3, 4, 2, 5, 1 (+1)	3, 5, 1, 4, 2 (+1)

3, 2, 5, 1, 4 (+1)

3, 2, 5, 4, 1 (-1)

3, 5, 2, 1, 4 (-1)

3, 5, 2, 4, 1 (+1)

Nu is $\psi(\beta_3 \cdot \beta_1 \cdot \beta_4) \cdot \psi(\beta_3 \cdot \beta_5 \cdot \beta_2) = -\psi(\beta_3 \cdot \beta_1 \cdot \beta_4) \cdot \psi(\beta_3 \cdot \beta_2 \cdot \beta_5)$ enz.,
zodat

$$60 D_5(\beta_1, \dots, \beta_5) = \frac{1}{8} \left\{ \delta \psi(\beta_3 \cdot \beta_1 \cdot \beta_4) \cdot \psi(\beta_3 \cdot \beta_2 \cdot \beta_5) + \delta \psi(\beta_3 \cdot \beta_2 \cdot \beta_5) \cdot \psi(\beta_3 \cdot \beta_1 \cdot \beta_4) \right\}$$

$$= \psi(\beta_3 \cdot \beta_1 \cdot \beta_4) \cdot \psi(\beta_3 \cdot \beta_2 \cdot \beta_5) + \psi(\beta_3 \cdot \beta_2 \cdot \beta_5) \cdot \psi(\beta_3 \cdot \beta_1 \cdot \beta_4)$$

of

$$(F.13) \quad 60 D_5(\beta_1, \dots, \beta_5) = \begin{cases} 0 & \text{as } \beta_3 \neq 3 \\ 2 & \text{as } \beta_3 = 3 \text{ en } \beta_1 \cdot \beta_2 < 3 \text{ of } \beta_1 \cdot \beta_2 > 3 \\ -1 & \text{as } \beta_3 = 3 \text{ en } \beta_1 < 3, \beta_2 > 3 \text{ of } \beta_1 > 3, \beta_2 < 3 \end{cases}$$

Nu is

$$(F.14) \quad D_n(\beta_1, \dots, \beta_n) = D_n(\beta_2 \cdot \beta_1 \cdot \beta_3, \dots, \beta_n) = D_n(\beta_1, \dots, \beta_{n-2} \cdot \beta_n \cdot \beta_{n-1})$$

$$= D_n(\beta_n \cdot \beta_{n-1}, \dots, \beta_1)$$

Voor $n=5$ volgt dit uit (F.13) en voor algemene n uit (F.11).
Uit de symmetrie van D_n t.o.v. x en y volgt ook dat D_n niet van waarde verandert, indien in de permutatie $(\beta_1, \dots, \beta_n)$ de getallen 1, 2 of de getallen $n-1, n$ verwisseld worden of als de permutatie door zijn inverse vervangen wordt.

In het geval van onafhankelijkheid hebben alle $n!$ mogelijke rangschikkingen dezelfde waarschijnlijkheid $\frac{1}{n!}$. Om dus de verdeling van D_n te vinden, moet het aantal rangschikkingen bepaald worden dat een bijzondere waarden van D_n geeft.

Voor $n=5$ zijn er $5! = 120$ rangschikkingen. Volgens (F.14) behoeven wij slechts de rangschikkingen te beschouwen met $\beta_1 < \beta_2$, $\beta_4 < \beta_5$, $\beta_1 < \beta_4$. Hun aantal is $\frac{120}{2 \cdot 2 \cdot 2} = 15$. Die rangschikkingen onder deze 15, waarbij $\beta_3 \neq 3$, geven $D_5 = 0$. Er blijven dus de volgende drie permutaties over

(1, 2, 3, 4, 5), (1, 4, 3, 2, 5), (1, 4, 3, 5, 2)

Uit (F.13) vinden we als respectievelijke waarden van $60 D_5$ 2, -1, -1. Dus is

$$P\{60 D_5 = 2\} = \frac{1}{15} ; \quad P\{60 D_5 = -1\} = \frac{2}{15} ; \quad P\{60 D_5 = 0\} = \frac{12}{15}$$

x	$15 P\{60 D_5 = x\}$	$P\{60 D_5 \geq x\}$
-1	2	1.0000
0	12	0.8667
2	1	0.0667

Analoog kunnen de verdelingen van $\underline{D}_6, \underline{D}_7, \dots$ berekend worden door gebruik te maken van de betrekkingen (F.11) en (F.14). Zo is voor $n = 6$:

x	$90 P\{180 \underline{D}_6 = x\}$	$P\{180 \underline{D}_6 \geq x\}$
-2	4	1.0000
-1	23	0.9556
0	36	0.6444
1	16	0.2444
2	1	0.0667
3	4	0.0556
6	1	0.0111

Hoefding geeft ook een tabel voor het geval $n = 7$.

Uit (F.13) en (F.11) volgt dat $-\frac{1}{60} \leq \underline{D}_n \leq \frac{2}{60} = \frac{1}{30}$ voor $n = 5, 6, \dots$

De bovenste grens $\frac{1}{30}$ wordt bereikt voor $\pi = (1, \dots, n)$ en iedere n . Volgens de gevallen $n = 5, 6, 7$ lijkt het alsof het maximum van \underline{D}_n schijnbaar toeneemt met n , $(-\frac{1}{60}, -\frac{1}{90}, -\frac{1}{115})$.

Uit $\mathcal{E} \underline{D}_n = \Delta$, volgt dat $\Delta \leq \frac{1}{30}$.

(V) Uitvoering van de D -toets om te toetsen op onafhankelijkheid.

Als een aselechte steekproef uit een tweedimensionale verdeling met een continue v.f. gegeven is, dan wordt als volgt op onafhankelijkheid getoetst:

$$H_0: F(x, y) = F(x, \infty) \cdot F(\infty, y)$$

Stel \int_n is het kleinste getal dat voldoet aan

$$(F.15) \quad P\{\underline{D}_n > \int_n | H_0\} \leq \alpha,$$

waarbij α ($0 < \alpha < 1$) de gekozen onbetrouwbaarheidsdrempel is.

Bereken \underline{D}_n zoals in par. F(II) is aangetoond. Als $\underline{D}_n > \int_n$ zoals uit de exacte verdeling verkregen is, dan wordt de nulhypothese H_0 , dat de twee stochastische grootheden x en y onafhankelijk zijn, verworpen. Voor $n = 5, 6, 7$ kan \int_n uit de berekende tabellen verkregen worden.

Uit de ongelijkheid van Tchebycheff en (F.10') volgt dat

$$P\{|30 \underline{D}_n - 30 \Delta| \geq \sqrt{\text{Var}(30 \underline{D}_n)} \cdot \frac{1}{\sqrt{\alpha}}\} \leq (\sqrt{\alpha})^2$$

maar $\Delta = 0$ bij onafhankelijkheid, dus

$$P\{30 \underline{D}_n > \sqrt{\frac{2(n^2 + 5n - 32)}{n(n-1)(n-3)(n-4)\alpha}}\} \leq \alpha$$

waaruit volgt

$$30 \rho_n \leq \sqrt{\frac{2(n^2 + 5n - 32)}{n(n-1)(n-3)(n-4)\alpha}}$$

omdat ρ_n het grootste getal is dat aan (F.15) voldoet. Dus is $\rho_n = O(\frac{1}{n})$.

Als $\Delta > 0$, is $\Delta - \rho_n > 0$ voor voldoende grote n . Dan is

$$P\{D_n > \rho_n | H\} \geq P\{|D_n - \Delta| \leq \Delta - \rho_n\} \geq 1 - \frac{\text{Var } D_n}{(\Delta - \rho_n)^2} \\ \rightarrow 1 \text{ als } n \rightarrow \infty \text{ (uit F.8),}$$

waarbij H de alternatieve hypothese $\Delta > 0$ voorstelt.

Dus volgt hieruit en uit stelling (F.1) dat de D -toets bruikbaar is ten opzichte van de klasse Ω'' , en dus asymptotisch zuiver is.

(VI) Enkele opmerkingen.

Volgens Hoefding lijkt het alsof het onderscheidingsvermogen van de D -toets in vergelijking met de productmomentcorrelatietoets voor een normale verdeling met een correlatie ρ , enigszins laag is voor kleine waarden van $|\rho|$ en voor $n \rightarrow \infty$. Het resultaat is misschien niet geheel bevredigend voor waarden van n waarin wij geïnteresseerd zijn. Aan de andere kant is te verwachten dat een toets, die bruikbaar is ten opzichte van een grote klasse van alternatieven, een lager onderscheidingsvermogen zal hebben met betrekking tot een deelklasse van alternatieven dan een toets, die optimale kenmerken ten opzichte van de deelklasse heeft. Dit werpt dus het probleem op om uit een gegeven klasse van parameter vrije toetsingsmethoden (zoals die, welke bruikbaar zijn ten opzichte van Ω'') een toetsingsmethode te kiezen, die het grootste onderscheidingsvermogen heeft ten opzichte van zekere parametrische alternatieven (b.v. normale verdelingen). (Zie Hoefding: Optimum Nonparametric Tests).

(VII) Een paar verdere stellingen.

Stel dat π_{nv} ($v = 1, 2, \dots, n!$) de $n!$ mogelijke rangschikkingen van steekproeven (grootte n) uit een tweedimensionale verdeling met continue v.f. $F(x, y)$ zijn (vergelijk par. F(IV)).

Als $F(x, y) = F(x, \infty) F(\infty, y)$, dan is

$$(F.16) \quad P\{\pi_{nv}\} = \frac{1}{n!} \quad (v = 1, \dots, n!)$$

voor iedere n .

Als (F.16) gegeven is, volgt dan daaruit dat, voor een zekere n , de veranderlijken x, y onafhankelijk zijn? Dit geldt niet voor $n=2$, want dan is (F.16) equivalent met $P\{(1,2)\} = \frac{1}{2}$. Indien de verdeling een waarschijnlijkheidsdichtheid $f(x, y)$ heeft, is

$$P\{(1,2)\} = \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} \left[\int_{-\infty}^x \int_{-\infty}^y f(u,v) du dv + \int_x^{\infty} \int_y^{\infty} f(u,v) du dv \right] f(x,y) dx dy$$

hetgeen $= \frac{1}{2}$ is wanneer $f(x,y) = f(-x,y)$.

Nu hebben wij echter

Stelling A. Als $F(x,y)$ tot de klasse Ω'' behoort en (F.16) voor een zekere $n \geq 5$ geldt, dan is

$$(F.17) \quad F(x,y) = F(x,\infty) \cdot F(\infty,y)$$

Bewijs: Omdat $\sum_{\nu=1}^n D_{n\nu} = \Delta$, kunnen wij schrijven

$$(F.18) \quad \sum_{\nu=1}^{n!} D_n(\pi_{n\nu}) P\{\pi_{n\nu}\} = \Delta.$$

Als (F.16) geldt, heeft het linkerlid van (F.18) dezelfde waarde als wanneer (F.17) geldt. Maar wanneer (F.17) geldt, is $\Delta = 0$.

Uit (F.16) volgt dus dat $\Delta = 0$ en uit stelling (F.1) volgt dat dit voldoende is voor (F.17). Hiermee is de stelling bewezen. Als dus de oorspronkelijke verdeling tot Ω'' behoort en als voor een zekere $n \geq 5$ de waarschijnlijkheden van de $n!$ rangpermutaties gelijk zijn, zijn de stochastische veranderlijken onafhankelijk.

Stelling B. Er bestaan geen toetsingsmethoden voor onafhankelijkheid, die alleen van de rangvolgorde van de waarnemingen afhangen en zuiver zijn bij elke onbetrouwbaarheidsdrempel ten opzichte van de klassen Ω' of Ω'' .

Bewijs: Elke critieke zone van een rangordetoets (d.i. een toets die alleen van de rangvolgorde van de waarnemingen afhangt) voor onafhankelijkheid is een verzameling $S_m = \{\pi_{n\nu_1}, \dots, \pi_{n\nu_m}\}$ van m rangschikkingen. In het geval van onafhankelijkheid is

$$P[S_m] = P[\pi_{n\nu} \in S_m] = \frac{m}{n!}$$

Wij kunnen ons dus beperken tot onbetrouwbaarheidsdrempels $\frac{m}{n!}$, $m = 1, 2, \dots, n! - 1$. Teneinde de stelling te bewijzen, zal het voldoende zijn om aan te tonen dat voor iedere $n = 2, 3, \dots$, voor een bepaalde m ($1 \leq m \leq n! - 1$) en voor iedere S_m , een v.f. F in Ω'' bestaat, zodanig dat

$$P[S_m | F] < \frac{m}{n!}$$

Wij bewijzen echter een iets meer algemene stelling dat dit geldt voor $m = 1, 2, 3$.

Wij definiëren de tweedimensionale verdeling A_n zodanig dat de waarschijnlijkheidsmassa gelijkmatig (uniform) verdeeld is over de $n-1$ segmenten

$$(F.19) \quad \frac{k-1}{n-1} < x \leq \frac{k}{n-1} \quad ; \quad y - x = \frac{n-2k}{n-1}$$

$$(k=1, 2, \dots, n-1)$$

en nul is in elk gebied dat geen deel van deze segmenten bevat.

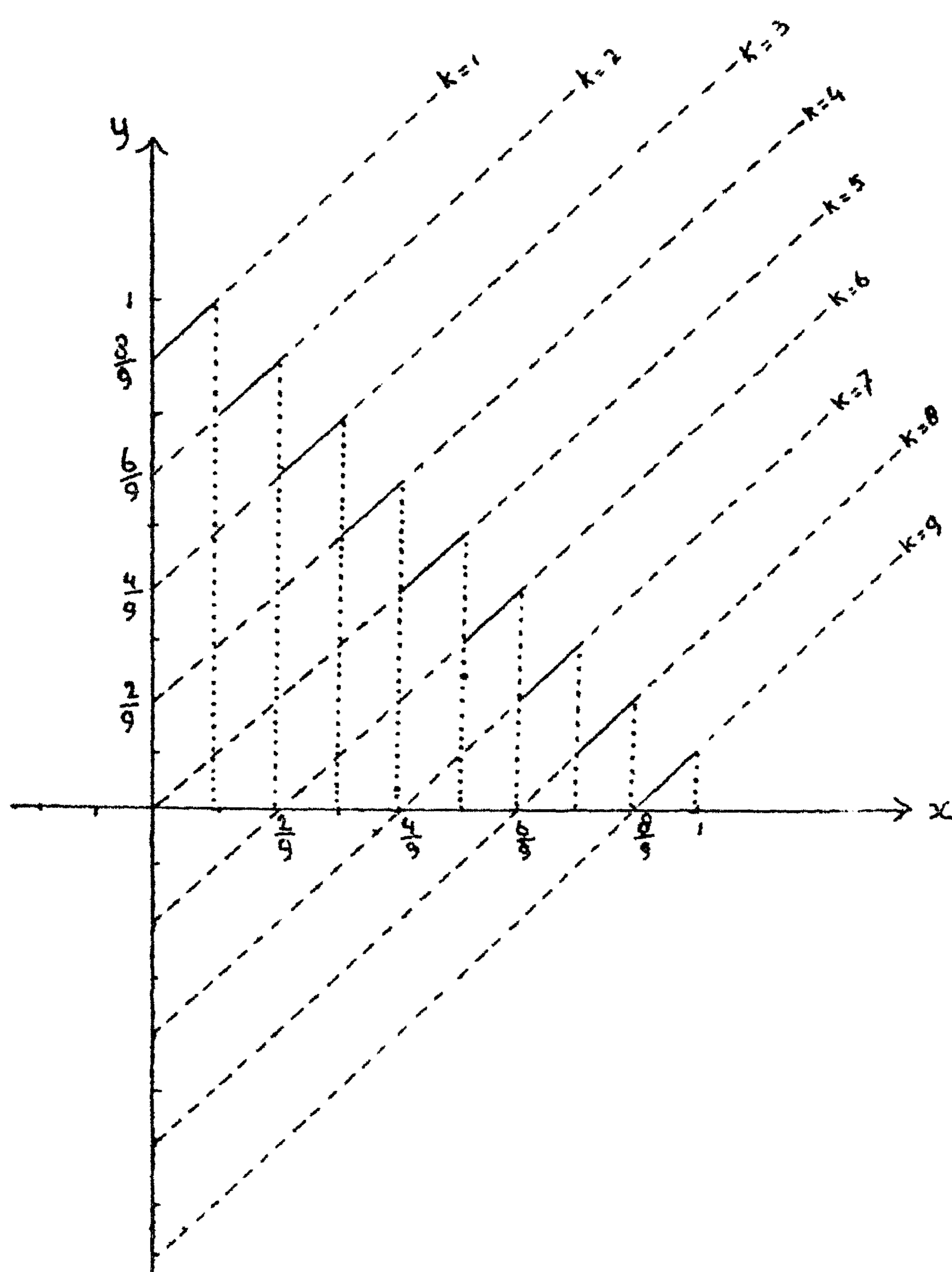
Stel B_n is een verdeling die gelijkmatig verdeeld is over de $n-1$ segmenten

$$(F.20) \quad \frac{k-1}{n-1} < x \leq \frac{k}{n-1} \quad ; \quad x + y = \frac{2k-1}{n-1} \quad (k=1, 2, \dots, n-1)$$

en een verdelingsdichtheid gelijk nul buiten deze segmenten.

De verdeling van de segmenten (F.19) ziet er dus als volgt uit:

Stel $n=10$:



k	x	y
1	$0 < x \leq \frac{1}{9}$	$x + \frac{8}{9}$
2	$\frac{1}{9} < x \leq \frac{2}{9}$	$x + \frac{6}{9}$
3	$\frac{2}{9} < x \leq \frac{3}{9}$	$x + \frac{4}{9}$
4	$\frac{3}{9} < x \leq \frac{4}{9}$	$x + \frac{2}{9}$
5	$\frac{4}{9} < x \leq \frac{5}{9}$	x
6	$\frac{5}{9} < x \leq \frac{6}{9}$	$x - \frac{2}{9}$
7	$\frac{6}{9} < x \leq \frac{7}{9}$	$x - \frac{4}{9}$
8	$\frac{7}{9} < x \leq \frac{8}{9}$	$x - \frac{6}{9}$
9	$\frac{8}{9} < x \leq 1$	$x - \frac{8}{9}$

Analoog kan de verdeling van de segmenten (F.20) overzichtelijk gemaakt worden.

De v.f.'s van A_n en B_n is continu, met

$$F(x, \infty) = F(\infty, x) = x \quad (0 \leq x \leq 1)$$

De waarschijnlijkheid dat (x, y) in elk van de segmenten (F.19) of (F.20) ligt, is $\frac{1}{n-1}$. De waarschijnlijkheden $P(\pi|A_n)$ en $P(\pi|B_n)$ kunnen gemakkelijk verkregen worden in termen van de multinomiale verdeling met $(n-1)$ gelijke waarschijnlijkheden.

In het bijzonder is

$$(F.21) \quad P(1, 2, \dots, n|A_2) = 1 \quad \text{en} \quad P(n, n-1, \dots, 1|B_2) = 1$$

(Want voor $n=2$ in (F.19) is $0 < x \leq 1$ en $y=x$. x_1, \dots, x_n liggen dus tussen 0 en 1. Als de x -waarden in volgorde $1, \dots, n$ vol-

gens rangnummers geranschikt zijn, zullen de y -waarden ook in volgorde staan ($y = x$). Let wel dat π de rangschikking is van de rangnummers van de y -waarden als de x -waarden volgens hun rangnummers geranschikt zijn. Dus $P(1, 2, \dots, n | A_2)$ is de waarschijnlijkheid dat $y_1 < y_2 < \dots < y_n$ wanneer $x_1 < x_2 < \dots < x_n$, voor de verdeling A_2 . Voor B_2 is $0 < x \leq 1$ en $y = 1 - x$.)

Ook is

$$(F.22) \quad P(1, 2, \dots, n | A_n) = P(n, n-1, \dots, 1 | B_n) = (n-1) \left(\frac{1}{n-1}\right)^n = \left(\frac{1}{n-1}\right)^{n-1}$$

(Er zijn $(n-1)$ segmenten waarover de n punten $(\underline{x}, \underline{y})$ verdeeld kunnen worden. De waarschijnlijkheid dat het punt $(\underline{x}, \underline{y})$ in een bepaald segment voor x valt is $\frac{1}{n-1}$. De waarschijnlijkheid dat alle n punten $(\underline{x}, \underline{y})$ in hetzelfde segment voor x valt, is $\left(\frac{1}{n-1}\right)^n$. Omdat er $(n-1)$ segmenten zijn, is de waarschijnlijkheid dat alle n punten $(\underline{x}, \underline{y})$ in eenzelfde segment voor x vallen, $(n-1) \left(\frac{1}{n-1}\right)^n$. Een rangschikking $(1, 2, \dots, n)$, d.w.z.

$y_1 < y_2 < \dots < y_n$ als $x_1 < x_2 < \dots < x_n$, is alleen dan mogelijk als alle n punten $(\underline{x}, \underline{y})$ in hetzelfde segment voor x vallen — dit volgt onmiddellijk uit de tabel op de vorige bladzijde voor $n = 10$. Dus $P(1, 2, \dots, n | A_n) = \left(\frac{1}{n-1}\right)^{n-1}$. Evenzo $P(n, n-1, \dots, 1 | B_n) = \left(\frac{1}{n-1}\right)^{n-1}$.)

Verder is $P(n, n-1, \dots, 1 | A_n) = P(1, 2, \dots, n | B_n) = 0$.

(Want n punten $(\underline{x}, \underline{y})$ (neem b.v. $n = 10$) kunnen als A_n gegeven is, niet zodanig over $n-1$ (d.i. 9) segmenten worden verdeeld dat de rangnummers van de y -waarden in een dalende volgorde staan bij stijgende volgorde van de rangnummers van de x -waarden.)

In het algemeen krijgen wij dus: als π_n alle permutaties van $1, 2, \dots, n$ voorstelt, is òf $P(\pi_n | A_n) = 0$, òf $P(\pi_n | B_n) = 0$. Want elke π_n met $P(\pi_n | A_n) \neq 0$ bevat ten minste een stijging in rangvolgorde van twee of meer opeenvolgende getallen (algemeen $i, i+1, \dots, i+k$) die niet voorafgegaan worden door kleinere getallen of gevolgd worden door grotere getallen. Aan de andere kant, als een π_n' met $P(\pi_n' | B_n) \neq 0$ een stijging in rangvolgorde bevat, wordt die of voorafgegaan door kleinere getallen of gevolgd door grotere getallen. Als dus $P(\pi_n | A_n) \neq 0$, is $P(\pi_n | B_n) = 0$ en soortgelijk als $P(\pi_n | B_n) \neq 0$, moet $P(\pi_n | A_n) = 0$ zijn.

Uit (F.21) volgt dat voor iedere verzameling S_m van m rangschikkingen die niet $(1, 2, \dots, n)$ of $(n, n-1, \dots, 1)$ bevat, is $P(S_m | A_2) = 0$ of $P(S_m | B_2) = 0$. Wij behoeven dus alleen critieke zones te beschouwen die zowel $(1, 2, \dots, n)$ als $(n, n-1, \dots, 1)$ bevatten. Voor $m = 1$ bestaan er geen zodanige zones. Voor $m = 2$ is er één. Uit (F.22) volgt dat voor $n > 2$

$$P(1, 2, \dots, n | A_n) + P(n, n-1, \dots, 1 | A_n) = \left(\frac{1}{n-1}\right)^{n-1} \\ < \frac{2}{n} \left(\frac{1}{n-1}\right)^{n-2} < \frac{2}{n!}$$

Ten slotte, als π_n iedere permutatie met uitzondering van $(1, 2, \dots, n)$ of $(n, n-1, \dots, 1)$ voorstelt, krijgen wij uit het voorgaande betoog voor A_n of ook voor B_n

$$P(1, 2, \dots, n) + P(n, n-1, \dots, 1) + P(\pi_n) = \left(\frac{1}{n-1}\right)^{n-1} < \frac{3}{n!}$$

Dus is ook $P(S_3 | F) < \frac{3}{n!}$. Hiermee is dus het bewijs voor de v.f.'s in Ω' voltooid.

Teneinde de stelling te bewijzen voor v.f.'s in Ω'' , kunnen wij de verdelingen A_n en B_n vervangen door verdelingen A_n' en B_n' die continue marginale en gezamenlijke dichtheden hebben en zodanig zijn dat de waarschijnlijkheden $P(\pi | A_n')$ en $P(\pi | B_n')$ willekeurig weinig verschillen van $P(\pi | A_n)$ respectievelijk $P(\pi | B_n)$. Zo b.v. kan A_2' gedefinieerd worden door de continue dichtheid

$$f(x, y) = K(\varepsilon - y + x) \quad \text{voor } 0 \leq y - x \leq \varepsilon, \quad x \leq 1 - \varepsilon, \quad y \geq \varepsilon; \\ = K(\varepsilon - x + y) \quad \text{voor } -\varepsilon \leq y - x \leq 0, \quad x \geq \varepsilon, \quad y \leq 1 - \varepsilon; \\ = K(x + y - \varepsilon) \quad \text{voor } x + y \geq \varepsilon, \quad x \leq \varepsilon, \quad y \leq \varepsilon; \\ = K(2 - \varepsilon - x - y) \quad \text{voor } x + y \leq 2 - \varepsilon, \quad x \geq 1 - \varepsilon, \quad y \geq 1 - \varepsilon; \\ = 0 \quad \text{overal elders,}$$

waarbij $K = \frac{3}{3\varepsilon^2 - 4\varepsilon^3}$ en $0 < \varepsilon \leq \frac{1}{2}$. Voor voldoende kleine ε , voldoet de verdeling aan de eisen.

Het bewijs van de stelling voor Ω'' toont ook aan dat er geen zuivere toetsingsmethode voor onafhankelijkheid, die alleen van de rangvolgorde van de waarnemingen afhangt, voor $n = 2$ en iedere betrouwbaarheidsdrempel (wij behoeven slechts $\varepsilon = \frac{1}{2}$ te nemen) bestaat. Volgens Hoefding kan ook bewezen worden dat voor $n = 3$ iedere $m = 1, 2, \dots, 5$ en iedere S_m , de ongelijkheid $P(S_m) < \frac{m}{3!}$ ten minste voor een van de verdelingen A_2, A_3, B_2, B_3 geldt. Het is nog onbekend of er rangvolgorde toetsingsmethoden voor onafhankelijkheid bestaan die zuiver zijn voor bepaalde steekproefgrootten n en bepaalde onbetrouwbaarheidsdrempels $\frac{m}{n!}$, ten opzichte van de klasse Ω'' van verdelingsfuncties F .

(VIII) Enkele opmerkingen.

Er bestaan dus geen toetsingsmethoden voor onafhankelijkheid die alleen van de rangvolgorde van de waarnemingen afhangen en die zuiver zijn op iedere onbetrouwbaarheidsdrempel ten

opzichte van de klasse Ω' of Ω'' van verdelingsfuncties (zie par. F(VII)).

Er bestaan echter wel toetsingsmethoden voor onafhankelijkheid die alleen van de rangvolgorde van de waarnemingen afhangen (dit zijn parameter vrije toetsingsmethoden) die bruikbaar zijn en dus asymptotisch zuiver, ten minste ten opzichte van de klasse Ω'' van verdelingsfuncties (zie einde van par. F(V)).

G. Rangcorrelatie en graadcorrelatie.

Hoefding [5] is er verder in geslaagd om de rangcorrelatiecoëfficiënt van Spearman (ρ) (zie Kendall [7]) met behulp van de graadcorrelatie K (zie Yule+Kendall [12]) en de rangcorrelatiecoëfficiënt τ van Kendall als een \mathcal{U} -statistische grootte uit te drukken.

Als een steekproef $(x_{\alpha}^{(1)}, x_{\alpha}^{(2)})$, ($\alpha = 1, \dots, n$) gegeven is met alle $x_{\alpha}^{(1)}$'s en alle $x_{\alpha}^{(2)}$'s verschillend, dan definieert Hoefding

$$k' = \frac{12}{(n^3 - n)} \sum_{\alpha=1}^n \left(R_{\alpha}^{(1)} - \frac{n+1}{2} \right) \left(R_{\alpha}^{(2)} - \frac{n+1}{2} \right) \quad (\text{zie ook par. 5.6})$$

hetgeen equivalent is met Spearman's ρ .

Substitutie van (C.3) hierin geeft

$$k' = \frac{3}{n^3 - n} \sum_{\alpha=1}^n \sum_{\beta=1}^n \sum_{\gamma=1}^n s(x_{\alpha}^{(1)} - x_{\beta}^{(1)}) s(x_{\alpha}^{(2)} - x_{\gamma}^{(2)})$$

hetgeen geschreven kan worden als

$$(G.1) \quad k' = \frac{(n-2)k + 3t}{n+1}$$

waarbij t de rangcorrelatiecoëfficiënt is (verg. D.3) en

$$k = \frac{3}{n(n-1)(n-2)} \sum'' s(x_{\alpha}^{(1)} - x_{\beta}^{(1)}) s(x_{\alpha}^{(2)} - x_{\gamma}^{(2)})$$

en de sommatie zich over alle onderling verschillende α, β, γ uitstrekt.

k is een \mathcal{U} -statistische grootte, en als een functie van een aselechte steekproef uit een verdeling met v.f. F , is k een zuivere schatting van de 3e graad regelmatige functioneel

$$(G.2) \quad K = 3 \int \dots \int s(x_1^{(1)} - x_2^{(1)}) s(x_1^{(2)} - x_3^{(2)}) dF(x_1) dF(x_2) dF(x_3)$$

Hoefding vindt voor continue en onafhankelijke veranderlijken $x^{(1)}, x^{(2)}$

$$(G.3) \quad \sigma^2(k) = \frac{n^2 - 3}{n(n-1)(n-2)}$$

$$(G.4) \quad \sigma^2(\underline{t}, \underline{k}) = \frac{2(n+2)}{3n(n-1)}$$

waaruit volgt

$$\begin{aligned} \sigma^2(\underline{k}') &= \frac{(n-2)^2 \sigma^2(\underline{k}) + 6(n-2) \sigma(\underline{t}, \underline{k}) + 9 \sigma^2(\underline{t})}{(n+1)^2} \\ &= \frac{1}{n-1} \end{aligned}$$

hetgeen overeenstemt met het resultaat door Kendall [7] afgeleid voor Spearman's ρ .

Volgens stelling 1.1 is $\sqrt{n}(\underline{k} - \mathcal{K})$ asymptotisch normaal verdeeld met gemiddelde nul en $9 \zeta_1(\mathcal{K})$ als variantie. Hetzelfde geldt voor de verdeling van de rangcorrelatiecoëfficiënt \underline{k}' , zoals volgt uit stelling 1.3 tezamen met (G.1).

Uit stelling 1.3 volgt ook dat de **gezamenlijke** verdeling van $\sqrt{n}(\underline{t} - \mathcal{T})$ en $\sqrt{n}(\underline{k} - \mathcal{K})$ (of $\sqrt{n}(\underline{k}' - \mathcal{K})$) tot de normale verdeling nader met varianties $4 \zeta_1(\mathcal{T})$ en $9 \zeta_1(\mathcal{K})$ en covariantie $6 \zeta_2(\mathcal{K}, \mathcal{T})$.

Literatuur.

- [1] H.Cramér, Mathematical methods of statistics, Princeton University Press, 1946.
- [2] H.E.Daniels and M.G.Kendall, The significance of rank correlations where parental correlation exists, Biometrika, vol 34, (1947), pp. 197-208.
- [3] P.R.Halmos, The theory of unbiased estimation, Annals of Math. Stat., vol. 17 (1946), pp. 34-43.
- [4] W.Hoeffding, On the distribution of the rank correlation coefficient ρ , when the variates are not independent, Biometrika, vol. 34 (1947), pp. 183-196.
- [5] W.Hoeffding, A class of statistics with asymptotically normal distribution, Annals of Math. Stat. vol. 19 (1948), pp. 293-325.
- [6] W.Hoeffding, A non-parametric test of independence, Annals of Math. Stat., vol. 19 (1948), pp. 546-557
- [7] M.G.Kendall, Rank correlation methods, Charles Griffin and Co. Ltd., London (1948).
- [8] M.G.Kendall, Advanced statistics, Charles Griffin and Co Ltd., London (1948), volume I.
- [9] H.B.Mann, Non-parametric tests against trend, Econometrica, vol. 13 (1945), pp. 245-259.
- [10] U.Nair, The standard error of Gini's mean difference, Biometrika, vol. 28 (1936), pp. 428-436.
- [11] M.G.Kendall, A new measure of rank correlation, Biometrika, vol. 30 (1938), pp. 81-93.
- [12] G.U.Yule and M.G.Kendall, Introduction to the theory of statistics, Charles Griffin and Co. Ltd., (1950), pp. 268-270.