

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 189 (M 73)

Exacte toetsen tegen kettingcorrelatie  
van M. OGAWARA en E.J. HANNAN

Verslag van een colloquium-voordracht over een  
artikel van E.J. HANNAN, (Biometrika, 1955).

door

Gerda Klerk-Grobbe

1956

The Mathematical Centre at Amsterdam, founded the 11th of February 1946, is a non-profit institution aiming at the promotion of pure mathematics and its applications, and is sponsored by the Netherlands Government through the Netherlands Organization for Pure Research (Z.W.O.) and the Central National Council for Applied Scientific Research in the Netherlands (T.N.O.), by the Municipality of Amsterdam and by several industries.

## Inhoud

1. Inleiding
  2. Methode van M. OGAWARA (1951)
  3. Methode van E.J. HANNAN (1955)
  4. Voor- en nadelen van de methode van OGAWARA en HANNAN
- Literatuur

## 1. Inleiding

M. OGAWARA (1951) leidde een exacte toets tegen kettingcorrelatie in een reeks waarnemingen af, terwijl E.J. HANNAN (1955) dezelfde methode toepaste wanneer aan deze waarnemingen een lineair regressiemodel ten grondslag ligt. HANNAN bespreekt beide gevallen en vergelijkt de afgeleide toetsen met de toets gebaseerd op een rechtstreekse schatting van de correlatiecoëfficiënt en met de toets van J. DURBIN en G.S. WATSON (1950, 1951). De procedure van HANNAN geeft bovendien zuivere schattingen en exacte betrouwbaarheidsgrenzen voor de regressiecoëfficiënten ook indien kettingcorrelatie aanwezig is. Deze schattingen zijn echter in vele gevallen niet de meest doeltreffende.

Voor een uitvoerigere bespreking van het lineaire regressiemodel en van de wenselijkheid van een toets tegen kettingcorrelatie verwijzen wij naar het colloquium-verslag S 174 (M 62) over de toets van DURBIN en WATSON.

## 2. Methode van M. OGAWARA (1951)

2.1. OGAWARA onderstelt dat voor de waarnemingen  $y_1, \dots, y_{2n+1}$  het model:

$$(2.1;1) \quad y_i = \alpha + u_i \quad (i=1, \dots, 2n+1) \quad 1)$$

voldoet, waarbij de grootheden  $u_i$  normaal verdeeld zijn met verwachting 0 en variantie  $\sigma^2$ ; ze behoeven echter niet onderling onafhankelijk te zijn, maar kunnen een kettingcorrelatie van de eerste orde vertonen, d.w.z.:

$$(2.1;2) \quad u_i = \rho u_{i-1} + v_i \quad (i=2, \dots, 2n+1),$$

waarin  $|\rho| < 1$  is en de  $v_i$  onderling onafhankelijk  $N(0, (1-\rho^2)\sigma^2)$  verdeeld zijn.

Een gebruikelijke toets voor de hypothese  $\rho=0$  is die gebaseerd op de schatting (M.G. KENDALL (1944)).

$$(2.1;3) \quad z_1 = \frac{\sum_{j=2}^{2n+1} (y_j - \bar{y}_0)(y_{j-1} - \bar{y}_0)}{\sqrt{\sum_{j=2}^{2n+1} (y_j - \bar{y}_0)^2 \sum_{j=2}^{2n+1} (y_{j-1} - \bar{y}_0)^2}}$$

voor de eerste orde kettingcorrelatie coëfficiënt  $\rho$ . Hierin is  $\bar{y}_0 = \frac{\sum_{i=1}^{2n+1} y_i}{2n+1}$ . De verdeling van  $z_1$  is echter nogal ingewikkeld, zodat hiervoor een benadering gebruikt moet worden.

-----  
1) Stochastische grootheden worden door onderstreepte letters weergegeven.

2)  $N(\mu, \sigma^2)$  gebruiken wij als symbool voor een normale verdeling met gemiddelde  $\mu$  en variantie  $\sigma^2$ .

$$x = \rho u + v$$

$$y = \rho z + w$$

$x$	$y$	$z$
$\rho u$	$\rho z$	$\rho z$
$v$	$w$	$w$
$\rho z$	$\rho z$	$\rho z$

$$-A_{11} = \frac{1+\rho^2}{1-\rho^2} \sigma^2 \quad -A_{13} = -\frac{\rho}{1-\rho^2} \sigma^2 \quad -A_{22} = \frac{1}{1-\rho^2} \sigma^2$$

$$-A_{12} = -\frac{\rho}{1-\rho^2} \sigma^2 \quad -A_{23} = 0 \quad -A_{33} = \frac{1}{1-\rho^2} \sigma^2$$

$$f(x, y, z) = \frac{1}{(2\pi)^{3/2} (1-\rho^2)^{3/2}} e^{-\frac{1}{2} \left( \frac{1+\rho^2}{1-\rho^2} \sigma^2 \right) x^2 + \frac{\rho \sigma^2}{1-\rho^2} xy + \frac{\rho^2 \sigma^2}{1-\rho^2} y^2 - \frac{\rho}{1-\rho^2} \sigma^2 xz - \frac{\rho^2 \sigma^2}{1-\rho^2} z^2}$$

$$f(x, y, z) = \frac{1}{\sqrt{2\pi} \sigma \sqrt{1-\rho^2}} e^{-\frac{1}{2} \left( \frac{1+\rho^2}{1-\rho^2} \right) \left[ x - \frac{\rho}{1+\rho^2} y - \frac{\rho}{1+\rho^2} z \right]^2}$$

$$f(x) = \frac{\rho}{1+\rho^2} y + \frac{\rho}{1+\rho^2} z$$

$$\left| \begin{array}{ccc} \frac{1+\rho^2}{1-\rho^2} & -\frac{\rho}{1-\rho^2} & -\frac{\rho}{1-\rho^2} \\ -\frac{\rho}{1-\rho^2} & \frac{1}{1-\rho^2} & 0 \\ -\frac{\rho}{1-\rho^2} & 0 & \frac{1}{1-\rho^2} \end{array} \right|$$

Let  $x = u_{22}$   
 $y = u_{21}$   
 $z = u_{23}$

$$\rightarrow E[u_{22} / u_{21}, u_{23}] = \frac{\rho}{1+\rho^2} (u_{21} + u_{23})$$

$$\sigma^2 [u_{22} / u_{21}, u_{23}] = \frac{1-\rho^2}{1+\rho^2} \sigma^2$$

2.2. De methode van OGAWARA om een exacte toets voor de hypothese  $\rho = 0$  af te leiden berust nu hierop, dat wij de voorwaardelijke verdeling van  $y_2, y_4, \dots, y_{2n}$  bekijken onder de voorwaarde dat  $y_1, y_3, \dots, y_{2n+1}$  (nl. de waargenomen waarden) bezitten. Bij gegeven  $y_{2t+1}$  ( $t=0, \dots, n$ ) mogen wij ook de  $u_{2t+1}$  (zie (2.1;1)) als constant, dus niet-stochastisch beschouwen. Volgens (2.1;2) zijn dan de grootheden  $u_{2t}$  ( $t=1, \dots, n$ ) op een constante na gelijk aan  $\underline{v}_{2t}$  en dus onderling onafhankelijk normaal verdeeld. Dit zelfde geldt dan voor  $y_{2t}$  ( $t=1, \dots, n$ ). Volgens (2.1;2) gelden nu de betrekkingen:

$$\underline{u}_{2t} = \rho u_{2t-1} + \underline{v}_{2t}$$

en

$$u_{2t+1} = \rho \underline{u}_{2t} + \underline{v}_{2t+1} \quad (t=1, \dots, n),$$

of

$$\underline{u}_{2t} = \rho u_{2t-1} + \underline{v}_{2t}$$

en

$$\rho^2 \underline{u}_{2t} = \rho u_{2t+1} - \rho \underline{v}_{2t+1}.$$

Door optelling vinden wij hieruit:

$$\underline{u}_{2t} = \rho(1+\rho^2)^{-1}(u_{2t-1} + u_{2t+1}) + (1+\rho^2)^{-1}(\underline{v}_{2t} - \rho \underline{v}_{2t+1}).$$

De verwachting en de variantie van  $\underline{u}_{2t}$  onder voorwaarde van gegeven  $u_1, u_3, \dots, u_{2n+1}$  zijn dus:

$$\mathcal{E}\{\underline{u}_{2t} | u_1, u_3, \dots, u_{2n+1}\} = \rho(1+\rho^2)^{-1}(u_{2t-1} + u_{2t+1})$$

en

$$\begin{aligned} \sigma^2\{\underline{u}_{2t} | u_1, u_3, \dots, u_{2n+1}\} &= (1+\rho^2)^{-2}(\sigma^2\{\underline{v}_{2t}\} + \rho^2\sigma^2\{\underline{v}_{2t+1}\}) = \\ &= (1+\rho^2)^{-1}(1-\rho^2)\sigma^2. \end{aligned}$$

Voor de verwachting van  $y_{2t}$  volgt hieruit (met (2.1;1)):

$$(2.2;1) \quad \mathcal{E}\{y_{2t} | y_1, y_3, \dots, y_{2n+1}\} = \alpha(1-2\rho(1+\rho^2)^{-1}) + \rho(1+\rho^2)^{-1}(y_{2t-1} + y_{2t+1}),$$

$(t=1, \dots, n).$

Noemen wij

$$\gamma_0 \stackrel{\text{def}}{=} \alpha [1 - 2\rho(1+\rho^2)^{-1}]$$

(2.2;2)

$$\gamma_1 \stackrel{\text{def}}{=} 2\rho(1+\rho^2)^{-1}$$

en

$$z_t \stackrel{\text{def}}{=} 2^{-1}(y_{2t-1} + y_{2t+1}) \quad (t=1, \dots, n)$$

dan wordt (2.2;1)

$$(2.2;3) \quad \mathcal{E}\{y_{2t} | y_1, \dots, y_{2n+1}\} = \gamma_0 + \gamma_1 z_t \quad (t=1, \dots, n).$$

Nu is volgens (2.1;1)

$$\mathcal{E}\{y_{2t} | y_1, y_3, \dots, y_{2n+1}\} = \alpha + \mathcal{E}\{u_{2t} | u_1, u_3, \dots, u_{2t+1}\}$$

en dus

$$y_{2t} = \mathcal{E}\{y_{2t} | y_1, y_3, \dots, y_{2n+1}\} + u_{2t} - \mathcal{E}\{u_{2t} | u_1, u_3, \dots, u_{2t+1}\}.$$

Noemen wij hierin nog

$$v_t' \stackrel{\text{def.}}{=} u_{2t} - \mathcal{E}\{u_{2t} | u_1, u_3, \dots, u_{2t+1}\}$$

dan is dus

$$(2.2;4) \quad y_{2t} = \gamma_0 + \gamma_1 z_t + v_t', \quad (t=1, \dots, n),$$

waarin de  $v_t'$  (onder voorwaarde van gegeven  $y_{2t+1}$  ( $t=0, \dots, n$ )) onderling onafhankelijk  $N(0; (1+\rho^2)^{-1}(1-\rho^2)\sigma^2)$  verdeeld zijn. Dit is dus een lineaire regressiebetrekking, welke aan alle eisen voldoet voor de geldigheid van de toetsen en schattingen die bij de regressieanalyse gebruikelijk zijn ( $F$ -toetsen en kleinste kwadraten-schattingen, zie S 174 (M 62)).

2.3. Zo vinden wij met de methode der kleinste kwadraten als meest aannemelijke (maximum likelihood) zuivere schatting voor  $\gamma_1$  ( $= 2\rho(1+\rho^2)^{-1}$ ):

$$(2.3;1) \quad \hat{\gamma}_1 = \frac{\sum_{t=1}^n (y_{2t} - \bar{y})(z_t - \bar{z})}{\sum_{t=1}^n (z_t - \bar{z})^2},$$

waarin  $\bar{y} \stackrel{\text{def.}}{=} \sum_{t=1}^n y_{2t} / n$  is en  $\bar{z} \stackrel{\text{def.}}{=} \sum_{t=1}^n z_t / n$ .

De hypothese  $\rho=0$  in model (2.1;2) komt overeen met de hypothese  $\gamma_1=0$  in model (2.2;3). Voor het toetsen van deze hypothese gebruikt men de toetsingsgrootte

$$(2.3;2) \quad \hat{\gamma}_1 \sqrt{(n-2) \sum_{t=1}^n (z_t - \bar{z})^2 / \sum_{t=1}^n [y_{2t} - \bar{y} - \hat{\gamma}_1 (z_t - \bar{z})]^2},$$

die een Student-verdeling met  $(n-2)$  vrijheidsgraden bezit, of het kwadraat hiervan:

$$(2.3;3) \quad \hat{\gamma}_1^2 (n-2) \sum_{t=1}^n (z_t - \bar{z})^2 / \sum_{t=1}^n [y_{2t} - \bar{y} - \hat{\gamma}_1 (z_t - \bar{z})]^2,$$

die een  $F$ -verdeling met 1 en  $(n-2)$  vrijheidsgraden bezit.

Ook andere waarden van  $\gamma_1$  kunnen op deze wijze getoetst worden: voor de hypothese  $\gamma_1 = \gamma_{10}$  is dan de toetsingsgrootte gelijk aan (2.3;2) of (2.3;3) als wij hierin  $\hat{\gamma}_1$  door  $(\hat{\gamma}_1 - \gamma_{10})$  vervangen.

2.4. Voor de hypothese  $\rho=0$  hebben wij nu de beschikking over twee toetsen: de toets gebaseerd op  $\underline{z}_1$ , en de toets met  $\hat{y}_1$ . Dit zijn dus twee mogelijke toetsen voor eenzelfde hypothese. Welke toets verdient de voorkeur? ook in de volgende paragrafen zullen wij dit probleem van een vergelijking tussen twee toetsen tegenkomen. De methode welke HANNAN voor deze vergelijking gebruikt kan met de volgende overwegingen plausibel gemaakt worden.

Stel voor het toetsen van de hypothese  $\theta = \theta_0$ , voor een onbekende parameter  $\theta$ , staan ons twee toetsingsgrootheden,  $\underline{t}_1$  en  $\underline{t}_2$ , ter beschikking. Stel verder dat deze toetsingsgrootheden beide een asymptotisch normale verdeling bezitten (als  $n$ , de steekproefuitgebreidheid, naar oneindig gaat) en dat ze als schatting voor de onbekende  $\theta$  asymptotisch zuiver zijn, d.w.z. dat

$$\lim_{n \rightarrow \infty} E\{\underline{t}_1\} = \theta \quad \text{en} \quad \lim_{n \rightarrow \infty} E\{\underline{t}_2\} = \theta$$

en dat hun varianties voor grote  $n$  van de orde  $n^{-1}$  zijn. Onder deze onderstellingen ligt het voor de hand de verhouding

$$\sigma^2\{\underline{t}_1 | \theta_0\} / \sigma^2\{\underline{t}_2 | \theta_0\}$$

te beschouwen en  $\underline{t}_1$  boven  $\underline{t}_2$  te verkiezen als deze verhouding  $< 1$  is (natuurlijk afgezien van andere meer technische overwegingen, zoals verschil in de bewerkelijkheid van de berekening van  $\underline{t}_1$  en  $\underline{t}_2$  en hun verdeling onder de nulhypothese).

De grootheid

$$E \stackrel{\text{def.}}{=} \lim_{n \rightarrow \infty} \sigma^2\{\underline{t}_2 | \theta_0\} / \sigma^2\{\underline{t}_1 | \theta_0\},$$

welke de relatieve asymptotische doeltreffendheid van  $\underline{t}_1$  en  $\underline{t}_2$  genoemd wordt, gebruikt HANNAN dan ook als criterium bij de vergelijking van de toetsen. Is dus  $E < 1$  dan is  $\underline{t}_2$  (asymptotisch) doeltreffender dan  $\underline{t}_1$ , is  $E = 1$  dan zijn zij gelijkwaardig en is  $E > 1$  dan is  $\underline{t}_1$  beter dan  $\underline{t}_2$ .

De keuze van  $E$  als criterium wordt nog gerechtvaardigd door een stelling van PITMAN (zie STUART (1954)), die zegt, dat (onder zekere regulariteitseisen voor de optredende waarschijnlijkheidsverdelingen en bovengenoemde onderstellingen) het quotient  $n_2/n_1$  van de aantallen waarnemingen welke bij gebruik van  $\underline{t}_2$  resp.  $\underline{t}_1$  nodig zijn om eenzelfde onderscheidingsvermogen te bereiken voor grote steekproeven tot  $E$  nadert (het onderscheidingsvermogen wordt hierbij bepaald voor alternatieven waarvan het verschil met  $\theta_0$  van de orde van  $n^{-1/2}$  is).



2.5. Wij zullen nu de beide toetsen voor  $\rho=0$  met elkaar vergelijken. De toetsingsgrootte  $\underline{x}_1$  is asymptotisch normaal verdeeld met gemiddelde  $\rho$  (zie WALKER (1954)), terwijl de variantie in grote steekproeven ( $2n+1$  waarnemingen,  $n$  groot) nadert tot:

$$(2.5;1) \quad \sigma^2\{\underline{x}_1\} \simeq (2n)^{-1}(1-\rho^2),$$

zie BARTLETT (1946).

Willen wij nu de stelling van PITMAN toepassen, dan moeten wij de toets met  $\hat{y}_1$  omzetten in een aequivalente toets met een schatting van  $\rho$  als toetsingsgrootte, daar  $\underline{x}_1$  ook een schatting van  $\rho$  is.

De toetsingsgrootte  $\hat{y}_1$  heeft als verwachting  $y_1 = 2\rho(1+\rho^2)^{-1}$  (zie (2.2;2)) Voor  $\rho$  volgt hieruit:  $\rho = y_1^{-1}(1 \pm \sqrt{1-y_1^2})$  waarbij alleen de oplossing met het -teken kleiner dan 1 is en dus zin heeft. Als schatting voor  $\rho$  kunnen wij gebruiken

$$(2.5;2) \quad \hat{\rho} = \hat{y}_1^{-1}(1 - \sqrt{1-\hat{y}_1^2}),$$

hetgeen een monotone functie van  $\hat{y}_1$  is, zodat het gebruik van  $\hat{\rho}$  als toetsingsgrootte aequivalent is aan het gebruik van  $\hat{y}_1$ . Deze schatting zal niet zuiver zijn en niet normaal verdeeld zijn zoals met  $\hat{y}_1$  het geval is. Wij zullen echter zien dat ook  $\hat{\rho}$  voor grote  $n$  bij benadering een normale verdeling met verwachting  $\rho$  bezit, zodat wij op  $\hat{\rho}$  de stelling van PITMAN kunnen toepassen ter vergelijking van de op  $\hat{y}_1$  gebaseerde toets met de op  $\underline{x}_1$  gebaseerde.

Om de asymptotische verdeling van  $\hat{\rho}$  te bepalen ontwikkelen wij  $\hat{\rho}$  naar opklimmende machten van  $(\hat{y}_1 - y_1)$ . Wij vinden hiervoor (gebruik makende van (2.5;2) en (2.2;2)):

$$(2.5;3) \quad \begin{aligned} \hat{\rho} &= \rho + \left\{ \frac{d\hat{\rho}}{d\hat{y}_1} \right\}_{\hat{y}_1=y_1} (\hat{y}_1 - y_1) + 2^{-1} \left\{ \frac{d^2\hat{\rho}}{d\hat{y}_1^2} \right\}_{\hat{y}_1=y_1} (\hat{y}_1 - y_1)^2 + \dots = \\ &= \rho + 2^{-1}(1+\rho^2)^2(1-\rho^2)^{-1}(\hat{y}_1 - y_1) + \dots \end{aligned}$$

Deze reeks is alleen convergent indien  $|\hat{y}_1 - y_1| < 1$  is. Echter  $\hat{y}_1$  is normaal verdeeld met verwachting  $y_1$ , terwijl in grote steekproeven de variantie  $\sigma^2\{\hat{y}_1\}$  van de orde  $n^{-1}$  is (hierop komen wij verderop nog terug). Dus de kans dat  $|\hat{y}_1 - y_1| > \varepsilon$ , bij willekeurig kleine gegeven  $\varepsilon$ , wordt willekeurig klein indien wij  $n$  maar groot genoeg kiezen. Voor een onderzoek naar het asymptotisch gedrag van  $\hat{\rho}$  zullen wij dus van deze reeks gebruik mogen maken en de hogere machten van  $(\hat{y}_1 - y_1)$  weg laten. Wij zien

dus hieruit dat  $\hat{\rho}$  bij benadering (voor grote  $n$ ) een normale verdeling bezit met verwachting

$$E\{\hat{\rho}\} \approx \rho + 2^{-1}(1+\rho^2)^2(1-\rho^2)^{-1} E\{\hat{y}_1 - y_1\} = \rho,$$

en variantie

$$(2.5;4) \quad \sigma^2\{\hat{\rho}\} \approx 4^{-1}(1+\rho^2)^4(1-\rho^2)^{-2} \sigma^2\{\hat{y}_1\}.$$

Uit (2.3;1) volgt dat

$$(2.5;5) \quad \begin{aligned} \sigma^2\{\hat{y}_1\} &= \sigma^2\{y_{2t}\} \left(\sum_{t=1}^n (x_t - \bar{x})^2\right)^{-1} \\ &= \sigma^2\{u_{2t}\} \left(\sum_{t=1}^n (x_t - \bar{x})^2\right)^{-1} \end{aligned}$$

is. Nu zal voor grote  $n$  de uitdrukking  $\sum_{t=1}^n (x_t - \bar{x})^2$  naderen tot

$$2^{-1}n\sigma^2(1+\rho^2), \quad 3)$$

zodat voor grote  $n$

$$(2.5;6) \quad \begin{aligned} \sigma^2\{\hat{y}_1\} &\approx (1+\rho^2)^{-1}(1-\rho^2)\sigma^2 / \{2^{-1}n\sigma^2(1+\rho^2)\} = \\ &= 2n^{-1}(1-\rho^2)(1+\rho^2)^{-2} \end{aligned}$$

is. Uit (2.5;4) en (2.5;6) volgt dus dat

$$\sigma^2\{\hat{\rho}\} \approx (2n)^{-1}(1-\rho^2)^{-1}(1+\rho^2)^2$$

is, voor grote  $n$ .

Voor de relatieve doeltreffendheid van  $\hat{\rho}$  en  $\underline{z}_1$ , dus van de toets met  $\hat{y}_1$  en die met  $\underline{z}_1$ , geldt dus

$$E = \lim_{n \rightarrow \infty} \sigma^2\{\underline{z}_1\} / \sigma^2\{\hat{\rho}\} = \lim_{n \rightarrow \infty} \frac{(2n)^{-1}(1-\rho^2)}{(2n)^{-1}(1-\rho^2)^{-1}(1+\rho^2)^2} = \left(\frac{1-\rho^2}{1+\rho^2}\right)^2.$$

Voor  $\rho = 0$  is  $E = 1$ ; voor het toetsen van de hypothese  $\rho = 0$  zijn dus beide toetsen gelijkwaardig. Voor het toetsen van hypothesen  $\rho = \rho_0 \neq 0$  zou echter de toets met  $\underline{z}_1$  doeltreffender zijn dan die met  $\hat{y}_1$  (want  $E = (1-\rho_0^2)^2 / (1+\rho_0^2)^2$  wordt kleiner naarmate  $\rho_0^2$  groter wordt) indien de verdeling van  $\underline{z}_1$ , voor  $\rho = \rho_0$ , bekend was. Dit is echter niet het geval. De hypothese  $\rho = \rho_0$  kan dus met  $\underline{z}_1$  nog niet getoetst worden, doch wel met  $\hat{y}_1$ .

3) Bij gegeven  $\underline{z}_t$  geldt voor de variantie van  $\hat{y}_1$  natuurlijk de uitdrukking (2.5;5). Voor grote  $n$  zal echter in de onvoorwaardelijke verdeling  $n^{-1} \sum_{t=1}^n (x_t - \bar{x})^2$  behoudens een willekeurig kleine kans naderen tot  $\sigma^2\{\underline{z}_t\} = 2^{-1}\sigma^2(1+\rho^2)$ . In de voorwaardelijke verdeling mogen wij daarom  $\sum_{t=1}^n (x_t - \bar{x})^2$  door  $n\sigma^2\{\underline{z}_t\}$  vervangen, mits  $n$  voldoende groot is.

### 3. Methode van E.J. HANNAN (1955)

3.1. Op dezelfde wijze als M. OGAWARA leidt E.J. HANNAN een exacte toets tegen kettingcorrelatie af bij lineaire regressie. Hij onderstelt dan dat de waarnemingen  $y_1, \dots, y_{2n+1}$  passen in het model

$$(3.1;1) \quad y_i = \alpha + \beta_1 x_{1i} + \dots + \beta_k x_{ki} + u_i \quad (i=1, \dots, 2n+1)$$

waarin  $\alpha, \beta_1, \dots, \beta_k$  onbekende coëfficiënten zijn. Terwille van de overzichtelijkheid zullen wij ons beperken tot één regressiecoëfficiënt  $\beta$ , dus tot het model:

$$(3.1;2) \quad y_i = \alpha + \beta x_i + u_i \quad (i=1, \dots, 2n+1);$$

voor  $k$  regressiecoëfficiënten lopen de afleidingen volkomen analoog. In dit model zijn de  $x_i$  bekende (waargenomen) waarden van de onafhankelijke variabele; tussen de grootheden  $u_i$  kan een eerste orde kettingcorrelatie bestaan, dus:

$$(3.1;3) \quad u_i = \rho u_{i-1} + v_i \quad (i=2, \dots, 2n+1),$$

zodat  $|\rho| < 1$  is en de  $v_i$ , onderling onafhankelijk en onafhankelijk van de  $x_i$ ,  $N(0, (1-\rho^2)\sigma^2)$  verdeeld zijn.

Om een exacte toets voor de hypothese  $\rho=0$  af te leiden bepalen wij weer, evenals in paragraaf 2.2, de voorwaardelijke verdeling van  $y_2, y_4, \dots, y_{2n}$  onder voorwaarde van de waargenomen waarden  $y_1, y_3, \dots, y_{2n+1}$ . Voor de voorwaardelijke verwachting en variantie van  $u_{2t}$  ( $t=1, \dots, n$ ) geldt dan weer:

$$E \left\{ u_{2t} \mid y_1, y_3, \dots, y_{2n+1} \right\} = \rho(1+\rho^2)^{-1} (u_{2t-1} + u_{2t+1})$$

en

$$\sigma^2 \left\{ u_{2t} \mid y_1, y_3, \dots, y_{2n+1} \right\} = (1+\rho^2)^{-1} (1-\rho^2) \sigma^2.$$

Met (3.1;2) volgt hieruit dat, onder voorwaarde van de gevonden waarden van  $y_1, y_3, \dots, y_{2n+1}$ ,

$$(3.1;4) \quad y_{2t} = \alpha(1-2\rho(1+\rho^2)^{-1}) + \beta x_{2t} + \rho(1+\rho^2)^{-1}(y_{2t-1} + y_{2t+1}) + \\ - \rho(1+\rho^2)^{-1} \beta (x_{2t-1} + x_{2t+1}) + v'_t \quad (t=1, \dots, n)$$

is, waarbij de grootheden  $v'_t = u_{2t} - E \left\{ u_{2t} \mid y_1, \dots, y_{2n+1} \right\}$  onderling onafhankelijk  $N(0, (1+\rho^2)^{-1}(1-\rho^2)\sigma^2)$  verdeeld zijn.

Als afkorting definiëren wij weer:

$$(3.1;5) \quad \begin{aligned} \gamma_0 &= \alpha [1 - 2\rho(1+\rho^2)^{-1}], & x_{1,t} &= x_{2t}, \\ \gamma_1 &= \beta, & x_{2,t} &= 2^{-1}(y_{2t-1} + y_{2t+1}), \quad (t=1, \dots, n), \\ \gamma_2 &= 2\rho(1+\rho^2)^{-1}, & x_{3,t} &= 2^{-1}(x_{2t-1} + x_{2t+1}). \\ \gamma_3 &= -\gamma_2 \beta, \end{aligned}$$

Substitueren wij dit in (3.1;4) dan ontstaat:

$$(3.1;6) \quad y_{2t} = \gamma_0 + \gamma_1 x_{1,t} + \gamma_2 x_{2,t} + \gamma_3 x_{3,t} + v'_t \quad , \quad (t=1, \dots, n).$$

Dit lineaire regressiemodel bezit alle gewenste eigenschappen om er de gebruikelijke methoden uit de regressie-analyse op te kunnen toepassen, mits wij de relatie  $\gamma_3 = -\gamma_2 \gamma_1$  verwaarlozen, hetgeen betekent, dat wij van de schattingen niet eisen, dat zij aan deze relatie voldoen. De kleinste kwadraten schattingen die wij op deze wijze voor de onbekende coëfficiënten uit (3.1;6), nl.  $\gamma_0$ ,  $\gamma_1$ ,  $\gamma_2$  en  $\gamma_3$ , afleiden zullen daarom niet overeenstemmen met de meest aannemelijke schattingen. Deze zouden verkregen worden indien de genoemde relatie wel aan de schattingen wordt opgelegd; de schattingen worden dan echter moeilijk te bepalen en onhandelbaarder wat hun verdeling betreft. De kleinste kwadraten schattingen zijn wel zuiver en zij geven ook nu exacte toetsen voor hypothesen omtrent de  $\gamma$ 's.

3.2. De hypothese  $\rho = 0$  komt nu in model (3.1;6) overeen met de hypothese  $\gamma_2 = 0$ <sup>4)</sup>. Wij zullen de kleinste kwadraten schatting voor  $\gamma_2$  en de toets voor  $\gamma_2 = 0$  niet afleiden maar alleen de resultaten vermelden; zij kunnen met de in de regressie-analyse gebruikelijke methoden worden gevonden.

Wij definiëren de gemiddelden:

$$\bar{y} = \sum_{t=1}^n y_{2t} / n \quad ; \quad \bar{x}_1 = \sum_{t=1}^n x_{1,t} / n \quad ; \quad \bar{x}_2 = \sum_{t=1}^n x_{2,t} / n \quad ; \quad \bar{x}_3 = \sum_{t=1}^n x_{3,t} / n$$

en verder de matrices:

$$Z = \begin{pmatrix} x_{1,1} - \bar{x}_1 & x_{2,1} - \bar{x}_2 & x_{3,1} - \bar{x}_3 \\ \vdots & \vdots & \vdots \\ x_{1,n} - \bar{x}_1 & x_{2,n} - \bar{x}_2 & x_{3,n} - \bar{x}_3 \end{pmatrix} \quad ; \quad y = \begin{pmatrix} y_2 - \bar{y} \\ y_4 - \bar{y} \\ \vdots \\ y_{2n} - \bar{y} \end{pmatrix} \quad ;$$

$$L = (Z'Z) = \begin{pmatrix} \sum_{t=1}^n (x_{1,t} - \bar{x}_1)^2 & \sum_{t=1}^n (x_{1,t} - \bar{x}_1)(x_{2,t} - \bar{x}_2) & \sum_{t=1}^n (x_{1,t} - \bar{x}_1)(x_{3,t} - \bar{x}_3) \\ \sum_{t=1}^n (x_{1,t} - \bar{x}_1)(x_{2,t} - \bar{x}_2) & \sum_{t=1}^n (x_{2,t} - \bar{x}_2)^2 & \sum_{t=1}^n (x_{2,t} - \bar{x}_2)(x_{3,t} - \bar{x}_3) \\ \sum_{t=1}^n (x_{1,t} - \bar{x}_1)(x_{3,t} - \bar{x}_3) & \sum_{t=1}^n (x_{2,t} - \bar{x}_2)(x_{3,t} - \bar{x}_3) & \sum_{t=1}^n (x_{3,t} - \bar{x}_3)^2 \end{pmatrix} ;$$

4) Eigenlijk volgt uit  $\rho = 0$  dat  $\gamma_2 = \gamma_3 = 0$  is. Om  $\rho = 0$  te toetsen kunnen wij dus ook deze hypothese toetsen. Welke van de twee toetsen de voorkeur verdient wat betreft hun onderscheidingsvermogen is nog niet onderzocht; de toets voor  $\gamma_2 = 0$  is wat de berekeningen betreft iets eenvoudiger.

$|L_{ij}|$  is de minor van het element in de  $i^e$  rij en de  $j^e$  kolom van  $L$ . Voor de kleinste kwadraten schattingen  $\hat{\gamma}_1$ ,  $\hat{\gamma}_2$  en  $\hat{\gamma}_3$  van resp.  $\gamma_1$ ,  $\gamma_2$  en  $\gamma_3$  geldt dan

$$\begin{pmatrix} \hat{\gamma}_1 \\ \hat{\gamma}_2 \\ \hat{\gamma}_3 \end{pmatrix} = L^{-1} Z' y$$

en dus is

$$(3.2;1) \quad \hat{\gamma}_2 = |L|^{-1} \begin{pmatrix} |L_{21}| & |L_{22}| & |L_{23}| \end{pmatrix} \begin{pmatrix} \sum_{t=1}^n (x_{1,t} - \bar{x}_1)(y_{2t} - \bar{y}) \\ \sum_{t=1}^n (x_{2,t} - \bar{x}_2)(y_{2t} - \bar{y}) \\ \sum_{t=1}^n (x_{3,t} - \bar{x}_3)(y_{2t} - \bar{y}) \end{pmatrix}$$

de kleinste kwadraten schatting voor  $\gamma_2 = 2\rho(1+\rho^2)^{-1}$  5).

De toets voor de hypothese  $\rho=0$ , dus  $\gamma_2=0$  heeft als toetsingsgrootheid

$$(3.2;2) \quad F = (n-4) \hat{\gamma}_2^2 |L| / (|L_{22}| \cdot \mathcal{Q}),$$

met 
$$\mathcal{Q} = \sum_{t=1}^n (y_{2t} - \bar{y} - \hat{\gamma}_1(x_{1,t} - \bar{x}_1) - \hat{\gamma}_2(x_{2,t} - \bar{x}_2) - \hat{\gamma}_3(x_{3,t} - \bar{x}_3))^2.$$

$F$  heeft een  $F$ -verdeling met 1 en  $(n-4)$  vrijheidsgraden.

De schatting  $\hat{\gamma}_1$  voor  $\gamma_1$ , waarvoor een overeenkomstige uitdrukking als (3.2;1) geldt, is tevens een zuivere schatting voor  $\beta$ . Analoog aan (3.2;2) (nl. door hierin  $\hat{\gamma}_2^2$  en  $|L_{22}|$  te vervangen door  $(\hat{\gamma}_1 - \gamma_{10})^2$  resp.  $|L_{11}|$ ) vinden wij een toets voor de hypothese  $\gamma_1 = \beta = \gamma_{10}$ . Deze toets is exact ook indien  $\rho \neq 0$  is.

3.3. HANNAN heeft weer, zoals in paragraaf 2.4 werd aangegeven, door middel van de relatieve asymptotische doeltreffendheid de toets (3.2;2) vergeleken met de toets van DURBIN en WATSON. Deze gebruiken als toetsingsgrootheid

$$d = \frac{\sum_{i=2}^{2n+1} (t_i - t_{i-1})^2}{\sum_{i=1}^{2n+1} t_i^2},$$

met 
$$t_i = y_i - \bar{y}_0 - b(x_i - \bar{x}_0) \quad , \quad (i=1, \dots, 2n+1),$$

5)  $\left( \frac{|L_{21}|}{|L|} \quad \frac{|L_{22}|}{|L|} \quad \frac{|L_{23}|}{|L|} \right)$  is de tweede rij uit  $L^{-1}$ .

waarin  $\bar{y}_o = \sum_{i=1}^{2n+1} y_i / (2n+1)$  is;  $\bar{x}_o = \sum_{i=1}^{2n+1} x_i / (2n+1)$  en  $\underline{b}$  de rechtstreekse kleinste kwadraten schatting voor  $\rho$  uit het oorspronkelijke materiaal.  $\underline{d}$  nadert asymptotisch tot  $2(1 - \underline{z}_1)$  als  $\underline{z}_1$  (vgl. (2.1;3)) berekend wordt uit de residuen  $t_i$  in plaats van  $y_i$ . Men kan afleiden dat de variantie van  $\underline{d}$  voor  $\rho = 0$  voor grote steekproeven van de orde  $(2n)^{-1}$  wordt.

Voor de variantie van de schatting voor  $\rho$ , die wij uit  $\hat{y}_2$  kunnen afleiden, wordt dan, als in paragraaf 2.5, gevonden  $(2n)^{-1}(1-\rho^2)^{-1}(1+\rho^2)^2$  dus voor  $\rho = 0$  eveneens  $(2n)^{-1}$ . Asymptotisch zijn beide toetsen voor  $\rho = 0$  dus even doeltreffend.

Op dezelfde wijze kunnen wij de schatting  $\hat{y}_1$  voor  $\beta$  met de rechtstreekse kleinste kwadraten schatting  $\underline{b}$  voor  $\beta$  vergelijken. Is de eerste orde kettingcorrelatiecoëfficiënt van de waarden  $x_i$  gelijk aan nul, dan is de verhouding tussen de asymptotische varianties van  $\underline{b}$  en van  $\hat{y}_1$ :

$$E = 2^{-1}(1+\rho^2)(1-\rho^2)^{-1}.$$

Deze verhouding is groter dan 1 voor  $|\rho\sqrt{3}| > 1$ ; voor deze waarden van  $\rho$  zal dus de schatting  $\hat{y}_1$  asymptotisch doeltreffender zijn (een kleinere variantie bezitten) dan de rechtstreekse schatting  $\underline{b}$ .

In het algemeen zijn de varianties zowel van  $\hat{y}_1$  als van  $\underline{b}$  ook nog afhankelijk van de correlatiecoëfficiënten tussen de  $x$ -en. Hebben in ons geval van één onafhankelijke variabele de  $x_i$  een eerste orde ketting-correlatie-coëfficiënt  $\underline{z}_x$  dan worden de varianties voor grote  $n$

$$\sigma^2\{\underline{b}\} = (2n)^{-1}\sigma^2 s_x^{-2} (1+\rho\underline{z}_x)(1-\rho\underline{z}_x)^{-1},$$

waarin  $s_x^2$  de variantie tussen de  $x_i$  ( $i = 1, \dots, 2n+1$ ) voorstelt,

en

$$\sigma^2\{\hat{y}_1\} \approx n^{-1}\sigma^2 s_x^{-2} (1-\rho^2)(1+\rho^2)^{-1}(1+\underline{z}_x^2)(1-\underline{z}_x^2)^{-1}$$

(zie WOLD (1953), pag. 211). De verhouding  $\sigma^2\{\underline{b}\}/\sigma^2\{\hat{y}_1\}$  wordt groter dan 1 voor grote  $\rho$  en  $\underline{z}_x$  met gelijk teken en veel kleiner dan 1 voor grote  $\rho$  en  $\underline{z}_x$  met verschillend teken. Bij kleine  $\rho$  is de verhouding steeds kleiner dan 1 (bij  $\rho = \underline{z}_x = 0$  wordt de verhouding  $1/2$ ); bij kleine  $\underline{z}_x$  en grote  $\rho$  groter dan 1. In tabel I zijn de door HANNAN berekende resultaten voor verschillende waarden van  $\rho$  en  $\underline{z}_x$  opgenomen.

Tabel I

$$E = \lim_{n \rightarrow \infty} \sigma^2\{\underline{b}\} / \sigma^2\{\underline{\hat{\beta}}_1\}.$$

$\rho \backslash z_x$	-0,8	-0,6	-0,4	-0,2	0	0,2	0,4	0,6	0,8
0	2,28	1,06	0,69	0,54	0,50	0,54	0,69	1,06	2,28
0,2	1,52	0,77	0,54	0,46	0,46	0,54	0,75	1,25	2,90
0,4	0,85	0,47	0,36	0,33	0,36	0,42	0,69	1,26	3,20
0,6	0,38	0,24	0,20	0,20	0,24	0,32	0,63	1,06	3,05
0,8	0,11	0,08	0,08	0,09	0,11	0,16	0,29	0,66	2,28

3.4. Een gebruikelijke methode om de regressiecoëfficiënten te schatten, indien positieve kettingcorrelatie verwacht wordt, is om niet van  $y_i$  en  $x_i$  uit te gaan, maar van de eerste verschillen  $y_i - y_{i+1}$  en  $x_i - x_{i+1}$  en uit deze waarden de kleinste kwadraten schattingen voor de coëfficiënten te berekenen (zie COCHRAN and ORCUTT (1949)).

Voor  $\underline{y}_i$  geldt het model (3.1;2) en (3.1;3) dus:

$$\underline{y}_i = \alpha + \beta x_i + \underline{u}_i, \quad (\text{zie paragraaf 3.1})$$

waarbij tussen de grootheden  $\underline{u}_i$  een eerste orde ketting-correlatie bestaat met correlatiecoëfficiënt  $\rho$ , dat wil dus zeggen dat

$$\mathcal{E}\{\underline{u}_i \underline{u}_{i-1}\} / \sigma^2\{\underline{u}_i\} = \rho$$

en

$$\mathcal{E}\{\underline{u}_i \underline{u}_{i-2}\} / \sigma^2\{\underline{u}_i\} = \rho^2 \text{ enz.}$$

Voor  $\underline{y}_i - \underline{y}_{i+1}$  geldt dan

$$\underline{y}_i - \underline{y}_{i+1} = \beta(x_i - x_{i+1}) + (\underline{u}_i - \underline{u}_{i+1})$$

De correlatiecoëfficiënt  $\rho'$  tussen  $\underline{v}'_i \stackrel{\text{def}}{=} \underline{u}_i - \underline{u}_{i+1}$  en  $\underline{v}'_{i-1} \stackrel{\text{def}}{=} \underline{u}_{i-1} - \underline{u}_i$  kunnen wij als volgt afleiden:

$$\begin{aligned} \rho' &= \mathcal{E}\{(\underline{u}_i - \underline{u}_{i+1})(\underline{u}_{i-1} - \underline{u}_i)\} / \sigma^2\{\underline{u}_i - \underline{u}_{i+1}\} = \\ &= (\mathcal{E}\{\underline{u}_i \underline{u}_{i-1}\} + \mathcal{E}\{\underline{u}_{i+1} \underline{u}_i\} - \mathcal{E}\{\underline{u}_{i+1} \underline{u}_{i-1}\} - \mathcal{E}\{\underline{u}_i^2\}) / \mathcal{E}\{(\underline{u}_i - \underline{u}_{i+1})^2\}. \end{aligned}$$

Nu is:

$$\mathcal{E}\{\underline{u}_i \underline{u}_{i-1}\} = \mathcal{E}\{\underline{u}_{i+1} \underline{u}_i\} = \rho \sigma^2,$$

$$\mathcal{E}\{\underline{u}_{i+1} \underline{u}_{i-1}\} = \mathcal{E}\{(\rho \underline{u}_i + \underline{v}_{i+1}) \underline{u}_{i-1}\} = \rho \mathcal{E}\{\underline{u}_i \underline{u}_{i-1}\} = \rho^2 \sigma^2,$$

$$\mathcal{E}\{\underline{u}_i^2\} = \sigma^2\{\underline{u}_i\} = \sigma^2,$$

en

$$\mathcal{E}\{(\underline{u}_i - \underline{u}_{i+1})^2\} = \mathcal{E}\{\underline{u}_i^2\} + \mathcal{E}\{\underline{u}_{i+1}^2\} - 2\mathcal{E}\{\underline{u}_i \underline{u}_{i+1}\} = 2\sigma^2 - 2\rho\sigma^2,$$

zodat dus

$$\rho' = \frac{\sigma^2(2\rho - \rho^2 - 1)}{2\sigma^2(1 - \rho)} = -2^{-1}(1 - \rho)$$

is. Hieruit volgt dat  $\rho'$  altijd negatief is, verder dat voor  $\rho > 0$  de waarde van  $\rho'$  tussen  $-\frac{1}{2}$  en 0 ligt en dat als  $\rho$  dichtbij 1 ligt  $\rho'$  ongeveer 0 wordt. Is dus  $\rho$  bijna 1 dan heeft men door de eerste verschillen te nemen een reeks verkregen waarvan de correlaties tussen opvolgende waarnemingen praktisch 0 is. De correlatie tussen de grootheden  $\underline{y}'_i$  is nu echter geen eerste orde ketting-correlatie meer: voor de correlatiecoëfficiënt tussen  $\underline{y}'_i$  en  $\underline{y}'_{i-2}$  vinden wij op de boven aangegeven manier  $\rho\rho'$ , voor de correlatiecoëfficiënt tussen  $\underline{y}'_i$  en  $\underline{y}'_{i-p}$  de waarde  $\rho^{p-1}\rho'$  in plaats van  $\rho^{p-1}$  zoals bij eerste orde ketting-correlatie het geval is.

HANNAN vergeleek weer door middel van hun relatieve asymptotische doeltreffendheid de schatting voor  $\beta$  uit deze eerste verschillen met  $\hat{\gamma}_1$  en vond dat voor  $\rho \geq 0$  de schatting uit de eerste verschillen van de waarnemingen asymptotisch doeltreffender is (een kleinere variantie bezit). Het nadeel bij deze methode is echter weer dat, indien  $\rho \neq 0$  de methode der kleinste kwadraten geen zuivere schatting voor de variantie van de regressiecoëfficiënt geeft en geen exacte toetsen voor hypothesen over  $\beta$ .

#### 4. Voor- en nadelen van de methode van OGAWARA en HANNAN

In de voorgaande paragrafen zijn al verschillende goede eigenschappen van de beschreven toetsings- en schattingsmethode naar voren gekomen. Daar was in de eerste plaats de exactheid van de verdelingen van de toetsingsgrootheden, ook indien  $\rho \neq 0$  is. Wij hebben zo niet alleen een exacte toets voor hypothesen betreffende  $\rho$  ( $\rho = 0$  of  $\rho = \rho_0$ ), maar ook exacte toetsen voor hypothesen over de regressiecoëfficiënten  $\beta_1, \dots, \beta_k$ . Dit laatste geeft tevens de mogelijkheid om exacte betrouwbaarheidsgrenzen voor de  $\beta$ 's te bepalen, hetgeen via de rechtstreekse kleinste kwadraten-schattingen alleen mogelijk is indien  $\rho = 0$  is.

Een nadeel van deze methode is het meerdere rekenwerk dat vereist wordt: komen in het oorspronkelijke model ( $k+1$ ) onbekende parameters voor ( $\alpha, \beta_1, \dots, \beta_k$ ), dan komen in het hieruit afgeleide model ( $2k+2$ ) parameters  $\gamma_0, \gamma_1, \dots, \gamma_{2k+1}$  voor. Om deze te schatten moet een  $(2k+1) \times (2k+1)$  -matrix geïnverteerd worden.



Voor de rechtstreekse kleinste kwadraten-schattingen, die voor de berekening van de grootheid  $d$  van DURBIN en WATSON gebruikt worden, kan met het inverteren van een  $k \times k$ -matrix worden volstaan. Bovendien is de schatting voor de regressiecoëfficiënten volgens de methode van HANNAN wel zuiver, maar lang niet altijd (ook niet als  $\rho \neq 0$ ) de meest doeltreffende.

Tenslotte zijn om de methode te kunnen toepassen vrij veel waarnemingen vereist. Bezit het oorspronkelijke model, buiten  $\rho$ ,  $k+1$  onbekende parameters, dus het hieruit afgeleide model  $2k+2$ , dan moeten er dus minstens  $2k+3$  waarden  $y_{2t}$  zijn, d.w.z.  $n$  moet minimaal  $2k+3$  zijn. Het oorspronkelijk waarnemingsmateriaal bevat  $2n+1$  waarnemingen, dus minstens  $2(2k+3)+1 = 4k+7$ .

-----

Literatuur

- M.S. BARTLETT (1946), On the theoretical specification and sampling properties of auto-correlated time series.  
J.R.S.S., Supl., 8, pp 27-41.
- D. COCHRAN and G.H. ORCUTT (1949), Application of least squares regression to relationships containing auto-correlated error terms.  
J.A.S.A., 44, pp. 32-61.
- J. DURBIN and G.S. WATSON (1950, 1951), Testing for serial correlation in least squares regression.  
Part I, Biometrika, 37, pp 409-428.  
Part II, Biometrika, 38, pp 159-178.
- E.J. HANNAN (1955), Exact tests for serial correlation.  
Biometrika, 42, pp 133-143.
- M.G. KENDALL (1944), On auto-regressive time series.  
Biometrika, 33, pp 105-122.
- M. OGAWARA (1951), A note on the test of serial correlation coefficients.  
Annals, 22, pp 115-118.
- A. STUART (1954), Asymptotic relative efficiencies of distribution-free tests of randomness against normal alternatives.  
J.A.S.A., 49, pp 147-157.
- A.M. WALKER, The asymptotic distribution of serial correlation coefficients for auto-regressive processes with dependent residuals.  
Proc. of the Cambridge Philosophical Society, 50, pp 60-64.
- H. WOLD (1953), Demand Analysis.  
New York, John Wiley & Sons.