

0269

SP5A

KONINKLIJKE NEDERLANDSE AKADEMIE  
VAN WETENSCHAPPEN

---

---

A rank-invariant method of linear and polynomial  
regression analysis

III

BY

H. THEIL

Reprinted from Proceedings Vol. LIII, No. 9, 1950

Reprinted from Indagationes Mathematicae, Vol. XII, Fasc. 5, 1950

1950

N.V. NOORD-HOLLANDSCHE UITGEVERS MAATSCHAPPIJ

(NORTH-HOLLAND PUBLISHING COMPANY)

AMSTERDAM



# A RANK-INVARIANT METHOD OF LINEAR AND POLYNOMIAL REGRESSION ANALYSIS. III <sup>1)</sup>

BY

H. THEIL

(Communicated by Prof. D. VAN DANTZIG at the meeting of Sept. 30, 1950)

## 3. CONFIDENCE REGIONS FOR THE PARAMETERS OF POLYNOMIAL REGRESSION EQUATIONS

*The probability set*

3. 0. The probability set  $\Gamma$  underlying the probability statements of this section is the  $n(\nu + 2)$ -dimensional Cartesian space  $R_{n(\nu+2)}$  with coordinates

$$u_{11}, \dots, u_{1n}, \dots, u_{\nu 1}, \dots, u_{\nu n}, v_1, \dots, v_n, w_1, \dots, w_n.$$

Every random variable mentioned is supposed to be defined on this probability set.

We suppose  $n(\nu + 2)$  random variables  $u_{\lambda i}$ ,  $v_i$ ,  $w_i$  ( $\lambda = 1, \dots, \nu$ ;  $i = 1, \dots, n$ ) to have a simultaneous probability distribution on  $\Gamma$ . Furthermore we consider  $n\nu$  parameters  $\xi_{\lambda i}$  and  $N$  parameters  $\alpha_{p_1 \dots p_\nu}$  for all sets of non-negative integers  $p_1, \dots, p_\nu$  satisfying

$$0 \leq \sum_{\lambda=1}^{\nu} p_{\lambda} \leq h.$$

Now we put <sup>2)</sup>

$$\left. \begin{aligned} (10) \quad \theta_i &= \sum \alpha_{p_1 \dots p_\nu} \xi_{1i}^{p_1} \dots \xi_{\nu i}^{p_\nu} \\ (11) \quad \eta_i &= \theta_i + w_i \\ (12) \quad x_{\lambda i} &= \xi_{\lambda i} + u_{\lambda i} \\ (13) \quad y_i &= \eta_i + v_i \end{aligned} \right\} \begin{cases} i = 1, \dots, n \\ \lambda = 1, \dots, \nu \end{cases}$$

So, for any set of values of the  $(N + n\nu)$  parameters  $\alpha_{p_1 \dots p_\nu}$ ,  $\xi_{\lambda i}$ , the variables  $x_{\lambda i}$  and  $y_i$  have a simultaneous distribution on  $\Gamma$ , and are therefore random variables.

The parameters  $\xi_{\lambda i}$  ( $i = 1, \dots, n$ ) are interpreted as values assumed by the variable  $\xi_{\lambda}$ . The equation (10) is the polynomial regression equation. The random variables  $w_i$  are called "the true deviations" from the poly-

<sup>1)</sup> This paper is the third of a series of papers, the first of which appeared in these Proceedings, 53, 386–392 (1950); the second appeared in these Proceedings, 53, 521–525 (1950).

<sup>2)</sup>  $\Sigma$  in equation (10) denotes summation over all sets  $p_1, \dots, p_\nu$ .

nomial of degree  $h$ ; the random variables  $\mathbf{u}_{\lambda i}$  and  $\mathbf{v}_i$  are called "the errors of observation" of the "true" values  $\xi_{\lambda i}$  and  $\eta_i$  respectively.<sup>3)</sup>

*Conditions; approximation*

3.1. In order to give confidence regions for the parameters  $\alpha_{p_1, \dots, p}$  we consider the following conditions:

*Condition I:* All  $n$   $(\nu + 2)$ -uples  $(\mathbf{u}_{\lambda i}, \mathbf{v}_i, \mathbf{w}_i)$  are stochastically independent.

*Condition IIa:* 1. Each of the errors  $\mathbf{u}_{\lambda i}$  vanishes outside a finite interval  $|\mathbf{u}_{\lambda i}| < g_{\lambda i}$ .

2. For each  $i \neq j$  we have  $|\xi_{\lambda i} - \xi_{\lambda j}| > g_{\lambda i} + g_{\lambda j}$ .

*Condition IIb:* 1. Each of the errors  $\mathbf{u}_{\lambda i}$  vanishes outside a finite interval  $|\mathbf{u}_{\lambda i}| < g_{\lambda i}$ .

2. For each  $i \neq j$ , for each set  $p_1, \dots, p_\nu$  and for any real  $h_{\lambda i}$  such that  $|h_{\lambda i}| \leq g_{\lambda i}$  we have

$$\operatorname{sgn} \left\{ \prod_{\lambda=1}^{\nu} (\xi_{\lambda i} + h_{\lambda i})^{\nu_{\lambda}} - \prod_{\lambda=1}^{\nu} (\xi_{\lambda j} + h_{\lambda j})^{\nu_{\lambda}} \right\} = \operatorname{sgn} \left\{ \prod_{\lambda=1}^{\nu} \xi_{\lambda i}^{\nu_{\lambda}} - \prod_{\lambda=1}^{\nu} \xi_{\lambda j}^{\nu_{\lambda}} \right\}.$$

*Condition III:* For all fixed values of the constants  $\varrho_{\lambda i}$  the  $n$  random variables  $\sum_{\lambda=1}^{\nu} \varrho_{\lambda i} \mathbf{u}_{\lambda i} + \mathbf{v}_i + \mathbf{w}_i = \mathbf{z}_i$  have continuous distribution functions, which are symmetrical with the median  $\operatorname{med}(\mathbf{z})$ .

Finally we mention that the solution will be given subject to the following

*Approximation:* For any positive  $s$  the quantities

$$\mathbf{u}_{\lambda i}^s \mathbf{u}_{\lambda' i}, \quad \mathbf{u}_{\lambda i}^s \mathbf{v}_i, \quad \mathbf{u}_{\lambda i}^s \mathbf{w}_i \quad (\lambda, \lambda' = 1, \dots, \nu; i = 1, \dots, n)$$

are neglected.<sup>4)</sup>

*Confidence regions*

3.2. We consider the case  $\nu = 1$ , so that equation (10) can be written as

$$\theta_i = \sum_{p=0}^h \alpha_p \xi_i^p.$$

Let us arrange the  $n$  observed points  $(x_i, y_i)$  according to increasing values of  $x$ :

$$x_1 < \dots < x_n.$$

<sup>3)</sup> It is clear that the random variables  $\mathbf{v}_i$  and  $\mathbf{w}_i$  cannot be separated in one sample of observations; if, however, the experiment is repeated for the same "true" values  $\xi_{\lambda i}, \eta_i$  (e.g. if — when the relation between income and consumption is investigated — for the same families and the same period the amounts of their incomes and outlays are repeatedly calculated), then the errors  $\mathbf{v}_i$  can be mitigated by averaging, whereas the deviations  $\mathbf{w}_i$  cannot.

<sup>4)</sup> The approximation implies that the errors  $\mathbf{u}_{\lambda i}$  are sufficiently small. This restriction is not very serious, because, unless the number of points  $n$  is very large, large values of  $\mathbf{u}_{\lambda i}$  will cause the confidence region for the parameters of the polynomial to be so large as to render the method useless.

We leave  $0, 1, \dots$  or  $h$  points out of consideration until the remaining number  $n'$  is such that  $n'/(h + 1)$  is an integer, and write  $n_h = n'/(h + 1)$ . (It seems advisable with respect to the power of the method to omit the points with rank  $n_h + 1, 2n_h + 1, \dots$  and / or  $hn_h + 1$ ). From now on we write  $n$  for the remaining number  $n'$ , so that  $(h + 1)n_h = n$ .

We define the following quantities:

$$\begin{aligned} \Delta(i, n_h + i) &= \frac{Y_i - Y_{n_h+i}}{x_i - x_{n_h+i}} \\ \Delta^{(2)}(i, n_h + i, 2n_h + i) &= \frac{\Delta(i, n_h + i) - \Delta(n_h + i, 2n_h + i)}{x_i - x_{2n_h+i}} \\ &\vdots \\ \Delta^{(h)}(i, n_h + i, \dots, hn_h + i) &= \\ &= \frac{\Delta^{(h-1)}(i, n_h + i, \dots, (h-1)n_h + i) - \Delta^{(h-1)}(n_h + i, 2n_h + i, \dots, hn_h + i)}{x_i - x_{hn_h+i}} \end{aligned}$$

We arrange the observed quantities  $\Delta^{(h)}(i, \dots, hn_h + i)$  according to increasing magnitude:

$$\Delta_1^{(h)} < \dots < \Delta_{n_h}^{(h)},$$

in which

$$\Delta_j^{(h)} = \Delta^{(h)}(i_j, \dots, hn_h + i_j).$$

3.3. Then we have the following theorem:

*Theorem 6:* Under conditions I, IIa, and III the interval  $(\Delta_{r_h}^{(h)}, \Delta_{n_h - r_h + 1}^{(h)})$  is a confidence interval for  $a_h$  to the approximate level of significance  $2I_{\frac{1}{2}}(n_h - r_h + 1, r_h)$ .<sup>5)</sup>

In order to prove this theorem we shall use the following *Lemma*. Define for all non-negative integers  $s$  and for all positive integers  $c$  and  $i$

$$P_{i, n_h+i, \dots, (c-1)n_h+i}^s = \sum_{s_i} \dots \sum_{\substack{s_{(c-1)n_h+i} \\ s_j \geq 0, \sum s_j = s}} x_i^{s_i} \dots x_{(c-1)n_h+i}^{s_{(c-1)n_h+i}}.$$

Then we have

$$\frac{P_{i, \dots, (c-1)n_h+i}^s - P_{n_h+i, \dots, cn_h+i}^s}{x_i - x_{cn_h+i}} = P_{i, \dots, cn_h+i}^{s-1}.$$

*Proof of the lemma:* We have

$$\begin{aligned} P_{i, \dots, (c-1)n_h+i}^s - P_{n_h+i, \dots, cn_h+i}^s &= \\ &= \sum_{s_i} \dots \sum_{s_{(c-1)n_h+i}} x_{n_h+i}^{s_{n_h+i}} \dots x_{(c-1)n_h+i}^{s_{(c-1)n_h+i}} (x_i^{s_i} - x_{cn_h+i}^{s_i}), \end{aligned}$$

<sup>5)</sup> In the first and second part of this paper the arguments of the incomplete Beta-function must be reversed.

in which  $\sum s_j = s$ . It follows that

$$\begin{aligned} & \frac{P_{i, \dots, (c-1)n_h+i}^s - P_{n_h+i, \dots, cn_h+i}^s}{x_i - x_{cn_h+i}} = \\ & = \sum_{s_i} \dots \sum_{s_{(c-1)n_h+i}} x_i^{s_{n_h+i}} \dots x_{(c-1)n_h+i}^{s_{(c-1)n_h+i}} (x_i^{s_i-1} + x_i^{s_i-2} x_{cn_h+i} + \dots + x_{cn_h+i}^{s_i-1}) = \\ & = P_{i, \dots, cn_h+i}^{s-1}. \end{aligned}$$

*Proof of Theorem 6:* The relation between  $\mathbf{x}_i$  and  $\mathbf{y}_i$  is given by

$$\begin{aligned} \mathbf{y}_i &= \sum_{p=0}^h a_p (\mathbf{x}_i - \mathbf{u}_i)^p + \mathbf{v}_i + \mathbf{w}_i \\ &\approx \sum_{p=0}^h a_p \mathbf{x}_i^p - \mathbf{u}_i (a_1 + 2a_2 \xi_i + \dots + ha_h \xi_i^{h-1}) + \mathbf{v}_i + \mathbf{w}_i, \end{aligned}$$

in which we neglected (in accordance with the Approximation)  $\mathbf{u}_i^s$  for  $s > 1$ . Putting  $\mathbf{z}_i = \rho_i \mathbf{u}_i + \mathbf{v}_i + \mathbf{w}_i$ , in which

$$-\rho_i = a_1 + 2a_2 \xi_i + \dots + ha_h \xi_i^{h-1},$$

we get

$$\mathbf{y}_i \approx \sum_{p=0}^h a_p \mathbf{x}_i^p + \mathbf{z}_i.$$

Now we have according to the lemma:

$$\begin{aligned} \Delta(i, n_h+i) &\approx a_1 + a_2 \frac{\mathbf{x}_i^2 - \mathbf{x}_{n_h+i}^2}{\mathbf{x}_i - \mathbf{x}_{n_h+i}} + \dots + a_h \frac{\mathbf{x}_i^h - \mathbf{x}_{n_h+i}^h}{\mathbf{x}_i - \mathbf{x}_{n_h+i}} + \frac{\mathbf{z}_i - \mathbf{z}_{n_h+i}}{\mathbf{x}_i - \mathbf{x}_{n_h+i}} \\ &= a_1 + a_2 \mathbf{P}_{i, n_h+i}^1 + \dots + a_h \mathbf{P}_{i, n_h+i}^{h-1} + \frac{\mathbf{z}_i - \mathbf{z}_{n_h+i}}{\mathbf{x}_i - \mathbf{x}_{n_h+i}}. \\ \Delta^{(2)}(i, n_h+i, 2n_h+i) &\approx a_2 + a_3 \mathbf{P}_{i, n_h+i, 2n_h+i}^1 + \dots + a_h \mathbf{P}_{i, n_h+i, 2n_h+i}^{h-2} + \\ &\quad \dots + \frac{\mathbf{z}_i - \mathbf{z}_{n_h+i}}{\mathbf{x}_i - \mathbf{x}_{n_h+i}} - \frac{\mathbf{z}_{n_h+i} - \mathbf{z}_{2n_h+i}}{\mathbf{x}_{n_h+i} - \mathbf{x}_{2n_h+i}} \\ &\quad + \frac{\mathbf{z}_i - \mathbf{z}_{2n_h+i}}{\mathbf{x}_i - \mathbf{x}_{2n_h+i}}. \end{aligned}$$

$$\Delta^{(h)}(i, n_h+i, \dots, hn_h+i) \approx a_h + \mathbf{Z}_i,$$

in which  $\mathbf{Z}_i$  is a random variable depending on

$$\mathbf{x}_i, \dots, \mathbf{x}_{hn_h+i}, \quad \mathbf{z}_i, \dots, \mathbf{z}_{hn_h+i}.$$

$\mathbf{Z}_i$  can be written as a fraction, the denominator being a product of terms  $(\mathbf{x}_{cn_h+i} - \mathbf{x}_{c'n_h+i})$  ( $c, c' = 0, \dots, h; c \neq c'$ ); according to condition IIa this denominator has a definite sign. The numerator consists of a sum of terms

$$(\mathbf{z}_{cn_h+i} - \mathbf{z}_{c'n_h+i}) \prod_{c'' \neq c'} (\mathbf{x}_{c''n_h+i} - \mathbf{x}_{c'n_h+i}).$$

But to our order of approximation this is equal to

$$(\mathbf{z}_{c'n_h+i} - \mathbf{z}_{c'n_h+i}) \prod_{c', c''} (\xi_{c'n_h+i} - \xi_{c''n_h+i}),$$

so that the numerator can be written as

$$\sum_{c=0}^h \tau_c \mathbf{z}_{c'n_h+i} \text{ with } \sum_c \tau_c = 0.$$

It follows from condition III that this quantity has zero median. From this and from the above-mentioned property of the denominator it follows that

$$P[\Delta^{(h)}(i, \dots, hn_h + i) < a_h | a_h] = P[\Delta^{(h)}(i, \dots, hn_h + i) > a_h | a_h] = \frac{1}{2}.$$

From this and from condition I the theorem immediately follows.

*Additional remarks*

3. 5. If  $\alpha_h$  is known, a confidence interval for  $\alpha_{h-1}$  can be found. Consider the equation

$$\mathbf{y}_i - \alpha_h \mathbf{x}_i^h \approx \sum_0^{h-1} \alpha_p \mathbf{x}_i^p + \mathbf{z}_i,$$

which shows that the problem is reduced to the case of a polynomial of degree  $(h-1)$ . So, if a confidence interval for  $\alpha_h$  is given, a confidence region for  $\alpha_h$  and  $\alpha_{h-1}$  can be found. This can be generalized to an  $(h+1)$ -dimensional confidence region for the parameters  $\alpha_0, \dots, \alpha_h$  in a way analogous to the one described in 2. 2. and 2. 3.

3. 6. If  $\nu > 1$ , an  $N$ -dimensional confidence region for the parameters  $\alpha_{p_1 \dots p_\nu}$  can be found in the following way:

1. Given the other parameters, a confidence region for the parameters  $\alpha_{p_1 0 \dots 0}$  ( $p_1 = 0, \dots, h$ ) in the  $N$ -dimensional parameter space can be constructed, the level of significance being  $\varepsilon_1$  (cf. 3. 5.).

2. In the same way one can proceed with the parameters

$$\alpha_{0p_2 \dots 0}, \dots, \alpha_{00 \dots p_\nu} \quad (p_\lambda = 1, \dots, h),$$

the levels of significance being  $\varepsilon_2, \dots, \varepsilon_\nu$ .

3. Finally the parameters  $\alpha_{p_1^1 \dots p_\nu^1}$  which have at least two indices  $p_i^1 \neq 0$ . We suppose that (apart from the conditions I and III) condition IIIb is valid, which is a more stringent condition than condition IIa. Consider the equation

$$\mathbf{R}_i \equiv \mathbf{y}_i - \sum \alpha_{p_1 \dots p_\nu} \mathbf{x}_{1i}^{p_1} \dots \mathbf{x}_{\nu i}^{p_\nu} \approx \alpha_{p_1^1 \dots p_\nu^1} \mathbf{x}_{1i}^{p_1^1} \dots \mathbf{x}_{\nu i}^{p_\nu^1} + \mathbf{z}_i,$$

in which  $\sum$  denotes summation over all sets  $p_1, \dots, p_\nu$  except the set  $p_1^1, \dots, p_\nu^1$ . We can then state regarding the quantities

$$S_{ij} = \frac{\mathbf{R}_i - \mathbf{R}_j}{\mathbf{x}_{1i}^{p_1^1} \dots \mathbf{x}_{\nu i}^{p_\nu^1} - \mathbf{x}_{1j}^{p_1^1} \dots \mathbf{x}_{\nu j}^{p_\nu^1}}$$

that

$$P[\mathbf{S}_{ij} < \alpha_{p_1^1 \dots p_\nu^1} | \alpha_{p_1^1 \dots p_\nu^1}] = P[\mathbf{S}_{ij} > \alpha_{p_1^1 \dots p_\nu^1} | \alpha_{p_1^1 \dots p_\nu^1}] = \frac{1}{2},$$

so that in a well-known way confidence regions for each of the parameters  $\alpha_{p_1^1 \dots p_\nu^1}$  can be found with levels of significance

$$\varepsilon_{\nu+1}, \dots, \varepsilon_{N-(h-1)\nu-1}.$$

The common part of the  $N-(h-1)\nu-1$  regions is a confidence region for the "true parameter point" in the  $N$ -dimensional parameter space, the level of significance being

$$\leq \sum_{q=1}^{N-(h-1)\nu-1} \varepsilon_q.$$

#### 4. CONFIDENCE REGIONS FOR THE PARAMETERS OF SYSTEMS OF REGRESSION EQUATIONS.

4. 0. In recent years considerable work has been done on the subject of systems of regression equations (see e.g. T. HAAVELMO (1943, 1944), T. KOOPMANS (1945, 1950), R. BENTZEL and H. WOLD (1946), M. A. GIRSHICK and T. HAAVELMO (1947)). In this section we shall give a brief investigation into the application of the methods considered on this subject.

##### *The probability set*

4. 1. Our probability set  $\Gamma$  will be the  $n(\nu + 2\tau)$ -dimensional Cartesian space  $R_{n(\nu+2\tau)}$  with coordinates

$$\begin{aligned} u_{11}, \dots, u_{1n}, \dots, u_{\nu 1}, \dots, u_{\nu n} \\ v_{11}, \dots, v_{1n}, \dots, v_{\tau 1}, \dots, v_{\tau n} \\ w_{11}, \dots, w_{1n}, \dots, w_{\tau 1}, \dots, w_{\tau n}. \end{aligned}$$

We suppose  $n(\nu + 2\tau)$  random variables  $\mathbf{u}_{\lambda i}$ ,  $\mathbf{v}_{\kappa i}$ ,  $\mathbf{w}_{ti}$  ( $i = 1, \dots, n$ ;  $\lambda = 1, \dots, \nu$ ;  $\kappa, t = 1, \dots, \tau$ ) to have a simultaneous probability distribution on  $\Gamma$ . Furthermore we consider  $n\nu + m\tau$  parameters  $\xi_{\lambda i}$ ,  $\alpha_{ij}$  ( $i = 1, \dots, n$ ;  $j = 1, \dots, m$ ;  $\lambda = 1, \dots, \nu$ ;  $t = 1, \dots, \tau$ ). Finally we consider the following equations

$$\begin{aligned} (14) \quad & f_t(\eta_{1i}, \dots, \eta_{\tau i}, \xi_{1i}, \dots, \xi_{\nu i}, \alpha_{i1}, \dots, \alpha_{im}) = \mathbf{w}_{ti} \\ (15) \quad & \mathbf{x}_{\lambda i} = \xi_{\lambda i} + \mathbf{u}_{\lambda i} \\ (16) \quad & \mathbf{y}_{\kappa i} = \eta_{\kappa i} + \mathbf{v}_{\kappa i}. \end{aligned} \left\{ \begin{array}{l} i = 1, \dots, n \\ \lambda = 1, \dots, \nu \\ \kappa, t = 1, \dots, \tau \end{array} \right.$$

The equations (14) are supposed to have a unique solution for  $\eta_{\kappa i}$  ( $\kappa = 1, \dots, \tau$ ;  $i = 1, \dots, n$ ) on every element of  $\Gamma$ , except possibly on a set of elements with zero probability.

The equations (14) are called the "stochastic regression equations". The parameters  $\xi_{\lambda i}$  ( $i = 1, \dots, n$ ) are interpreted as the values which the variable  $\xi_{\lambda}$  assumes ( $\lambda = 1, \dots, \nu$ ). The random variables  $\mathbf{w}_{ti}$  are called



the "true deviations" in the stochastic regression equations. Finally, the random variables  $u_{\lambda i}$  and  $v_{\kappa i}$  are called the "errors of observation" of the "true" values  $\xi_{\lambda i}$  and  $\eta_{\kappa i}$  respectively.

The problem is again, to determine confidence regions for the parameters  $a_{ij}$ .

### Confidence regions

4. 2. We reduce the equations (14), (15) and (16) to the forms

$$y_{\kappa i} = g_{\kappa}(x_{1i}, \dots, x_{vi}, u_{1i}, \dots, u_{vi}, v_{1i}, \dots, v_{vi}, w_{1i}, \dots, w_{vi}, \\ a_{11}, \dots, a_{1m}, \dots, a_{\tau 1}, \dots, a_{\tau m}). \quad (\kappa = 1, \dots, \tau)$$

Consider e.g. the case

$$f_t \equiv H_t(\xi_{1i}, \dots, \xi_{vi}) + \sum_{\kappa=1}^{\tau} \beta_{t\kappa} \eta_{\kappa i}, \quad (t = 1, \dots, \tau)$$

in which  $\beta_{t\kappa}$  are real numbers and  $H_t$  are polynomials of degree  $h$  in the  $\xi$ 's. Suppose that the errors  $u_{\lambda i}$  are sufficiently small in order that terms containing  $u_{\lambda i} u_{\lambda' i}$  ( $\lambda, \lambda' = 1, \dots, \nu$ ) can be neglected (cf. the Approximation of section 3. 1.); then we have

$$H_t(x_{1i}, \dots, x_{vi}) + \sum_{\kappa=1}^{\tau} \beta_{t\kappa} y_{\kappa i} \approx z_{ti}, \quad (t = 1, \dots, \tau)$$

in which  $z_{ti}$  are linear functions of  $u_{\lambda i}, v_{\kappa i}, w_{ii}$  ( $\lambda = 1, \dots, \nu; \kappa, t = 1, \dots, \tau$ ). The random variables  $(z_{1i}, \dots, z_{\tau i})$  ( $i = 1, \dots, n$ ) have a simultaneous probability distribution, while the  $n$   $\tau$ -uples  $(z_{1i}, \dots, z_{\tau i})$  are supposed to be stochastically independent. So we have

$$(17) \quad y_{\kappa i} \approx \sum_{t=1}^{\tau} \frac{B_{t\kappa}}{B} \{-H_t(x_{1i}, \dots, x_{vi}) + z_{ti}\},$$

in which

$$B = \begin{vmatrix} \beta_{11} & \dots & \beta_{1\tau} \\ \vdots & & \vdots \\ \beta_{\tau 1} & \dots & \beta_{\tau \tau} \end{vmatrix}$$

and  $B_{t\kappa}$  is the cofactor of the element  $\beta_{t\kappa}$ .

Then the problem is reduced to the case considered in section 3. Call  $N$  the number of parameters of a polynomial of degree  $h$ . Then, in a way and under conditions which are analogous to those stated in section 3, a confidence region for  $\tau N$  parameters of the equations (17) can be given. But the original equations contain  $\tau(N + \tau - 1)$  parameters. This means that, if  $\tau(\tau - 1)$  parameters of the original equations are given, a confidence region for the remaining  $\tau N$  parameters can be constructed. If the level of significance of the confidence regions for the parameters of the equations (17) are  $\varepsilon_{\kappa}$  ( $\kappa = 1, \dots, \tau$ ), this level of the confidence region for the  $\tau N$  parameters of the original equations is  $\leq \sum_{\kappa=1}^{\tau} \varepsilon_{\kappa}$ .

4. 3. We shall now elaborate a simple example, which is due to T. HAAVELMO (1944), p. 99 seq. Suppose we have the following equations:

$$\left. \begin{aligned} \eta_{1i} - \beta \eta_{2i} &= \mathbf{w}_{1i} \\ \alpha \xi_{1i} + \eta_{1i} - \alpha \eta_{2i} &= \mathbf{w}_{2i} \\ x_{1i} &= \xi_{1i} \\ \mathbf{y}_{1i} &= \eta_{1i} + \mathbf{v}_{1i} \\ \mathbf{y}_{2i} &= \eta_{2i}, \end{aligned} \right\} i = 1, \dots, n$$

in which  $\alpha$  and  $-\beta$  are positive.

We obtain:

$$\mathbf{y}_{1i} = \frac{\alpha\beta}{\alpha-\beta} x_{1i} + \frac{\alpha\mathbf{w}_{1i} - \beta\mathbf{w}_{2i}}{\alpha-\beta} + \mathbf{v}_{1i}$$

$$\mathbf{y}_{2i} = \frac{\alpha}{\alpha-\beta} x_{1i} + \frac{\mathbf{w}_{1i} - \mathbf{w}_{2i}}{\alpha-\beta}.$$

Suppose that the complete or the incomplete method gives two confidence intervals

$$a_1 \leq \frac{\alpha\beta}{\alpha-\beta} \leq a_2$$

$$b_1 \leq \frac{\alpha}{\alpha-\beta} \leq b_2$$

with levels of significance  $\varepsilon_1$  and  $\varepsilon_2$  respectively. Then we obtain two confidence regions in the  $\alpha, \beta$ -plane, bounded by hyperbolas and by straight lines respectively. The probability that the common part contains the "true" point  $(\alpha, \beta)$  is  $\geq 1 - \varepsilon_1 - \varepsilon_2$ . (See fig. 2).

#### *On multicollinearity*

4. 4. As a final application we consider the following case. The following equations are given (cf. section 2. 0.):

$$\left. \begin{aligned} \theta_i &= \alpha_0 + \alpha_1 \xi_{1i} + \alpha_2 \xi_{2i} \\ \eta_i &= \theta_i + \mathbf{w}_i \\ \mathbf{x}_{\lambda i} &= \xi_{\lambda i} + \mathbf{u}_{\lambda i} \\ \mathbf{y}_i &= \eta_i + \mathbf{v}_i. \end{aligned} \right\} \begin{cases} i = 1, \dots, n \\ \lambda = 1, 2. \end{cases}$$

Hence

$$\mathbf{y}_i = \alpha_0 + \alpha_1 \mathbf{x}_{1i} + \alpha_2 \mathbf{x}_{2i} + \mathbf{z}_i$$

with

$$\mathbf{z}_i = \mathbf{v}_i + \mathbf{w}_i - \alpha_1 \mathbf{u}_{1i} - \alpha_2 \mathbf{u}_{2i}.$$

Suppose that the observed values  $x_{1i}, x_{2i}$  ( $i = 1, \dots, n$ ) are such that the following *condition* is satisfied:

For each pair  $j$  ( $i, j = 1, \dots, n$ ) the quotient

$$\frac{x_{1i} - x_{1j}}{x_{2i} - x_{2j}} \quad (i \neq j)$$

has the same sign.

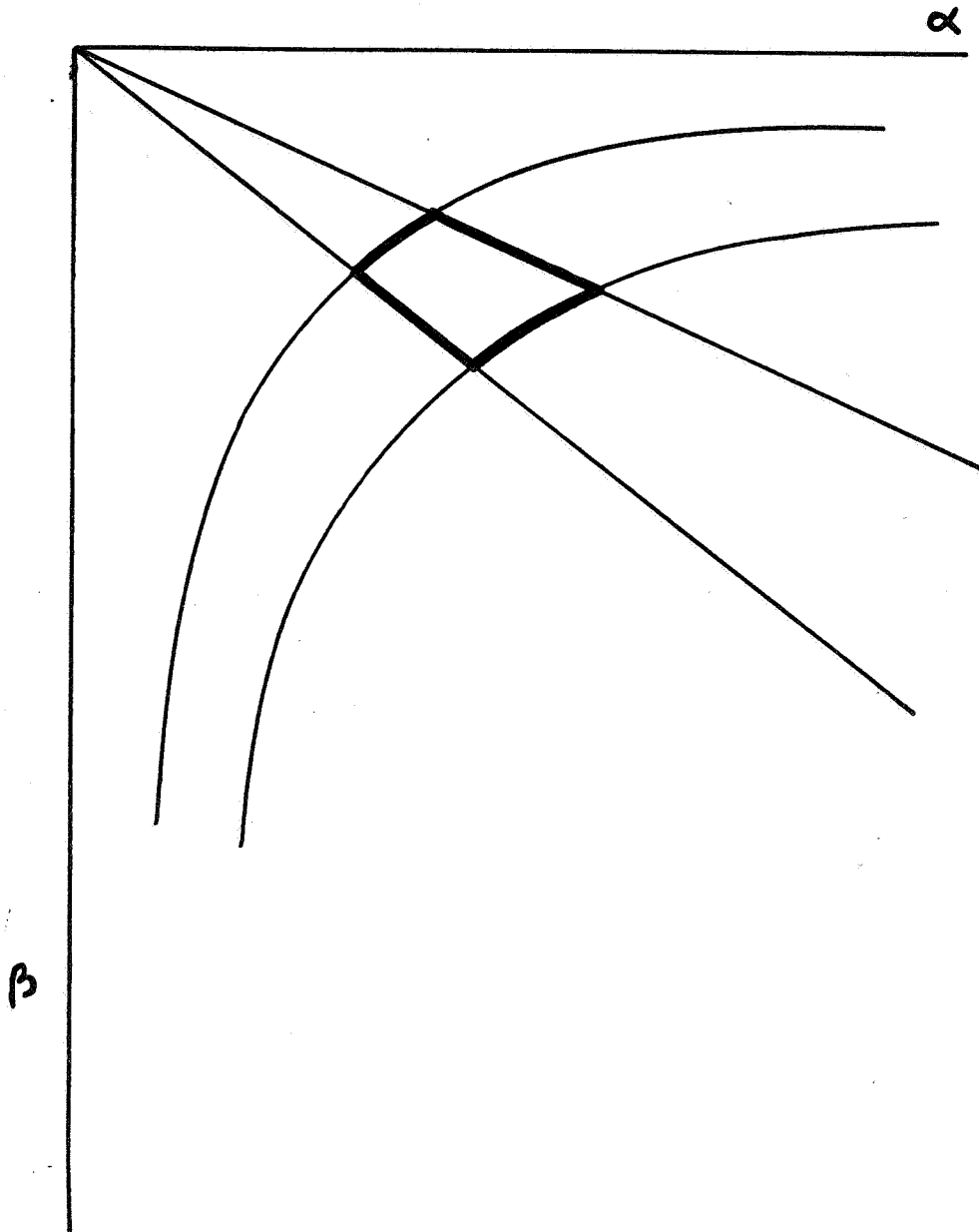


Fig. 2

This condition implies that, apart from the above-mentioned linear relation between  $x_{1i}$ ,  $x_{2i}$  and  $y_i$ , we have an additional monotonic relation between the observed values  $x_{1i}$  and  $x_{2i}$  (if this relation also is — approximately — linear, we have a case of “multicollinearity”).

We now have the following

*Theorem 7.* Under the above-mentioned condition the regions  $A_1$  and  $A_2$  (cf. section 2. 2.) are identical, and their common part  $A$  is unbounded.

Proof. If the condition is satisfied the arrangement of the observed points  $(x_{1i}, x_{2i}, y_i)$  according to increasing values of  $x_1$  is the same as (or just the reverse of) the arrangement according to increasing values of  $x_2$ . Moreover (cf. section 2. 2.) the quantities  $K^{(1)}(ij)$  and  $K^{(2)}(ij)$  which are estimates of  $a_1$ , given  $a_2$ , and of  $a_2$ , given  $a_1$ , respectively are represented by the same set of straight lines in the  $a_1, a_2$ -plane:

$$\begin{aligned}(x_{1i} - x_{1j}) K^{(1)}(ij) + (x_{2i} - x_{2j}) a_2 &= y_i - y_j \\ (x_{1i} - x_{1j}) a_1 + (x_{2i} - x_{2j}) K^{(2)}(ij) &= y_i - y_j.\end{aligned}$$

As the slopes of these straight lines  $-(x_{1i} - x_{1j})/(x_{2i} - x_{2j})$  have the same sign, the regions  $A_1$  and  $A_2$  are identical, from which the theorem follows.

If the incomplete method instead of the complete method is used, the same theorem holds with respect to the regions  $A'_1$  and  $A'_2$ , whereas the condition that all quantities  $(x_{1i} - x_{1j})/(x_{2i} - x_{2j})$  have the same sign is weakened to the condition that all quantities  $(x_{1i} - x_{1, n_1+i})/(x_{2i} - x_{2, n_1+i})$  have the same sign ( $i = 1, \dots, n_1$ ).

## 5. PROBLEMS OF PREDICTION

### *The probability set*

5. 0. For the probability set and the random variables defined on it we refer to section 4. 1. We assume, however, that all errors  $u_{\lambda i}, v_{\kappa i}$  are identically equal to zero ( $\lambda = 1, \dots, \nu; \kappa = 1, \dots, \tau; i = 1, \dots, n$ ).

### *Conditions*

5. 1. We impose the following conditions:

*Condition I:* All  $n$   $\tau$ -uples  $(w_{1i}, \dots, w_{\tau i})$  are distributed independently of each other.

*Condition IIIa:* All  $n$   $\tau$ -uples  $(w_{1i}, \dots, w_{\tau i})$  have the same continuous simultaneous distribution function.

Apart from these conditions we shall use the additional conditions, which are necessary for the determination of a confidence region for the parameters of the regression equations.

### *The problem*

5. 2. Suppose that the following  $n$  points are observed:

$$(\xi_{1i}, \dots, \xi_{\nu i}, \eta_{1i}, \dots, \eta_{\tau i}).$$

Suppose further that the following  $\nu$  parameters are given:

$$\xi_{1, n+1}, \dots, \xi_{\nu, n+1}.$$

These parameters are interpreted as the  $\xi$ -coordinates of an  $(n+1)$ -th point, which is not observed. The problem is to determine a confidence region for the  $\eta$ -coordinates of this point, i.e. for

$$\eta_{1, n+1}, \dots, \eta_{\nu, n+1}.$$

*Confidence regions*

5.3. Consider again the case (cf. section 4. 2.):

$$f_t \equiv H_t(\xi_{1t}, \dots, \xi_{\tau t}) + \sum_{i=1}^{\tau} \beta_{ti} \eta_{xi}, \quad (t = 1, \dots, \tau)$$

so that we have

$$\eta_{xi} = \sum_{t=1}^{\tau} \frac{B_{tx}}{B} \{-H_t(\xi_{1i}, \dots, \xi_{\tau i}) + \mathbf{w}_{ti}\}. \quad (\kappa = 1, \dots, \tau; i = 1, \dots, n)$$

Putting  $i = n + 1$  we can write

$$\eta_{\kappa, n+1} = g_{\kappa}(\xi_{1, n+1}, \dots, \xi_{\tau, n+1}, \beta) + h_{\kappa}(\mathbf{w}_{1, n+1}, \dots, \mathbf{w}_{\tau, n+1}, \beta),$$

in which

$$g_{\kappa}(\xi_{1, n+1}, \dots, \xi_{\tau, n+1}, \beta) = - \sum_{t=1}^{\tau} \frac{B_{t\kappa}}{B} H_t(\xi_{1, n+1}, \dots, \xi_{\tau, n+1})$$

$$h_{\kappa}(\mathbf{w}_{1, n+1}, \dots, \mathbf{w}_{\tau, n+1}, \beta) = \sum_{t=1}^{\tau} \frac{B_{t\kappa}}{B} \mathbf{w}_{ti}$$

and in which  $\beta$  is the "true parameter point";  $\beta$  may be considered as a vector, the components of which are  $\beta_{t\kappa}(t, \kappa = 1, \dots, \tau)$  and all parameters determining the polynomials  $H_t(t = 1, \dots, \tau)$ .

Suppose  $\beta$  is known. Then we can arrange the  $n$  quantities  $h_{\kappa}(w_{1i}, \dots, w_{\tau i}, \beta)$  according to increasing magnitude:

$$h_{\kappa 1} < \dots < h_{\kappa n}, \quad (\kappa = 1, \dots, \tau)$$

in which

$$h_{\kappa j} = h_{\kappa}(w_{1i_j}, \dots, w_{\tau i_j}, \beta).$$

We have the following

*Theorem 8:* Under conditions I and IIIa a confidence interval for  $\eta_{\kappa, n+1}$  is given by

$$(g_{\kappa}(\xi_{1, n+1}, \dots, \xi_{\tau, n+1}, \beta) + h_{\kappa s}, g_{\kappa}(\xi_{1, n+1}, \dots, \xi_{\tau, n+1}, \beta) + h_{\kappa, n-s+1})$$

if  $\beta$  is the known "true parameter point"; the level of significance is  $2s(n+1)^{-1}$ .

In order to prove this theorem, we shall use the following lemma (see W. R. THOMPSON (1936)):

*Lemma:* If a random sample of size  $n$  is drawn from a universe with continuous distribution function; if the sample values are arranged in ascending order; if an  $(n+1)$ -th draw from the same universe is to be effected; then the probability that the stochastic interval bounded by the  $s$ -th and the  $(n-s+1)$ -th of these values will contain the  $(n+1)$ -th is equal to  $1 - 2s/(n+1)$ .

*Proof of Theorem 8:* As  $g_{\kappa}(\xi_{1, n+1}, \dots, \xi_{\tau, n+1}, \beta)$  is ex hypothesi a known quantity, the problem is to determine a confidence interval for

$h_x(w_{1,n+1}, \dots, w_{\tau,n+1}, \beta)$ . But  $n$  sample values  $h_x(w_{1i}, \dots, w_{\tau i}, \beta)$  from the same universe (cf. condition IIIa) are obtained; hence the lemma is sufficient in order to show the validity of the theorem.

5. 4. Generally, however,  $\beta$  is unknown, and we can only calculate a confidence region  $R$  for  $\beta$ . Let now  $\beta$  vary through  $R$ , and denote by  $J_x$  the interval bounded by the lowest of all lower limits of the interval considered in Theorem 8 and by the highest of all upper limits ( $x = 1, \dots, \tau$ ). If the level of significance of  $R$  is  $\varepsilon$ , the following theorem immediately follows:

*Theorem 9:*

$$P[\eta_{1,n+1} \in J_1, \dots, \eta_{\tau,n+1} \in J_\tau, \beta \in R] \geq (1-\varepsilon) \left(1 - \frac{2\tau s}{n+1}\right).$$

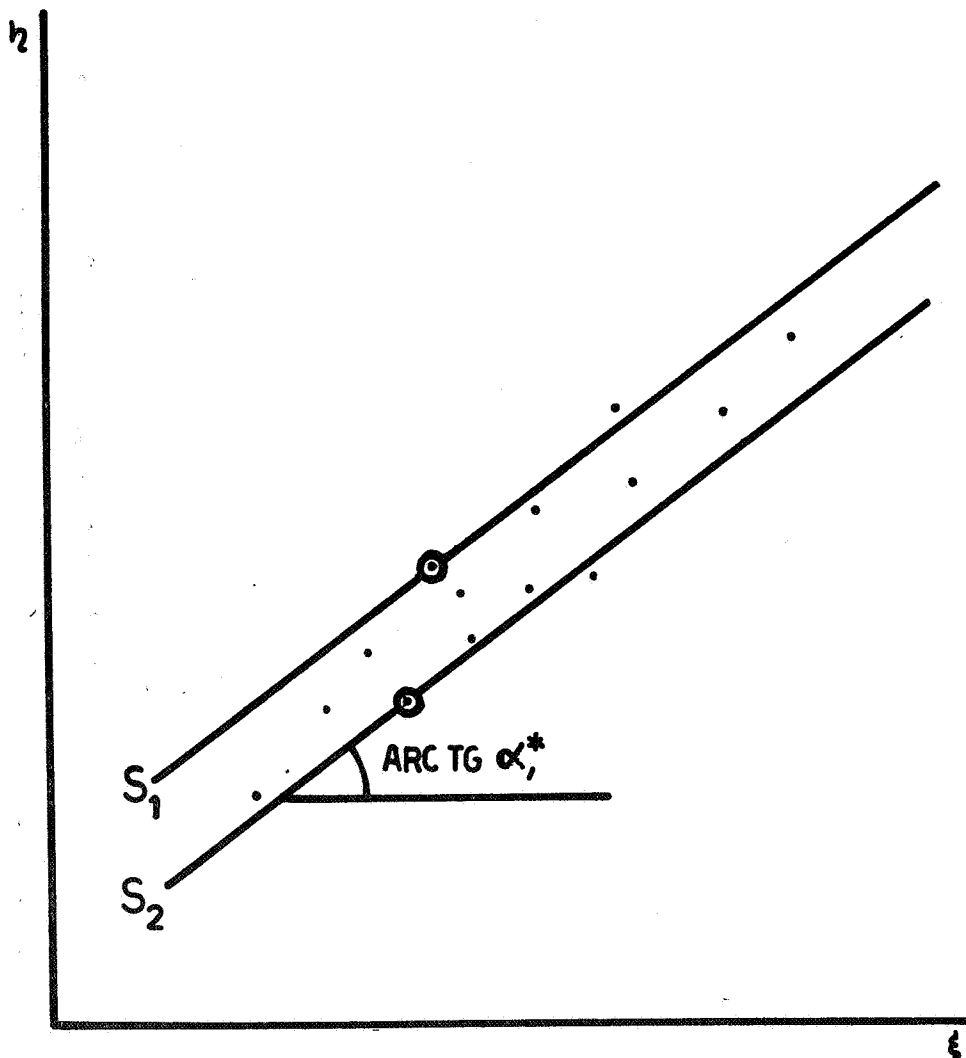


Fig. 3.  $n = 14$   $s = 2$

The linear case in two variables

5. 5. For the linear case of two variables

$$\eta_i = \alpha_0 + \alpha_1 \xi_i + w_i$$

a simple graphical representation can be given. Suppose that  $\alpha_1^*$  is the "true"  $\alpha_1$ ; after arranging the sample values

$$w_i(\alpha_1^*) = \eta_i - \alpha_0 - \alpha_1^* \xi_i$$

in order we find two straight lines:

$S_1$ : 
$$\eta = \alpha_0 + \alpha_1^* \xi + w_s(\alpha_1^*)$$

and  $S_2$ : 
$$\eta = \alpha_0 + \alpha_1^* \xi + w_{n-s+1}(\alpha_1^*).$$

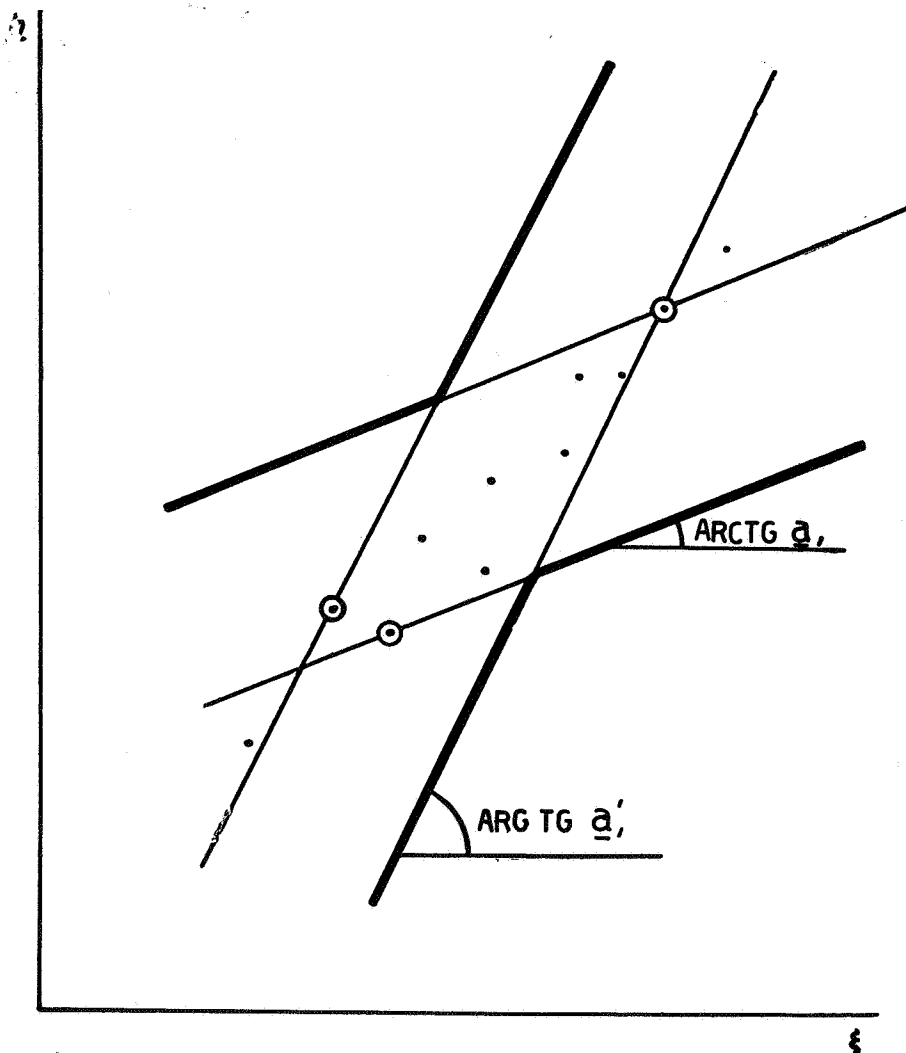


Fig. 4.  $n = 11$   $s = 2$

The probability that the region bounded above by  $S_1$  and below by  $S_2$  will contain an  $(n + 1)$ -th sample point  $D$  is, under the condition that  $\alpha_1^* = \alpha_1$ , equal to  $1 - 2s/(n + 1)$ . (See fig. 3).

Suppose that the confidence interval for  $\alpha_1$  is  $(\alpha_1, \alpha_1')$ . When  $\alpha_1^*$  varies through this interval the lines  $S_1$  and  $S_2$  revolve around the observed points  $(\xi_s, \eta_s)$  with ranks  $s$  and  $(n - s + 1)$  respectively with respect to increasing values of  $w$ . As long as the observed points having these properties remain the same,  $S_1$  and  $S_2$  revolve around one point; but as soon as variation of  $\alpha_1^*$  causes another point to have this property, the revolution takes place around this point. The figures 4 and 5 elucidate the fact that sometimes the region is bounded by the straight lines  $S_1$  and  $S_2$  for  $\alpha^* = \alpha_1$  and  $\alpha_1'$  only, whereas it is sometimes necessary to consider values between  $\alpha_1$  and  $\alpha_1'$  as well.

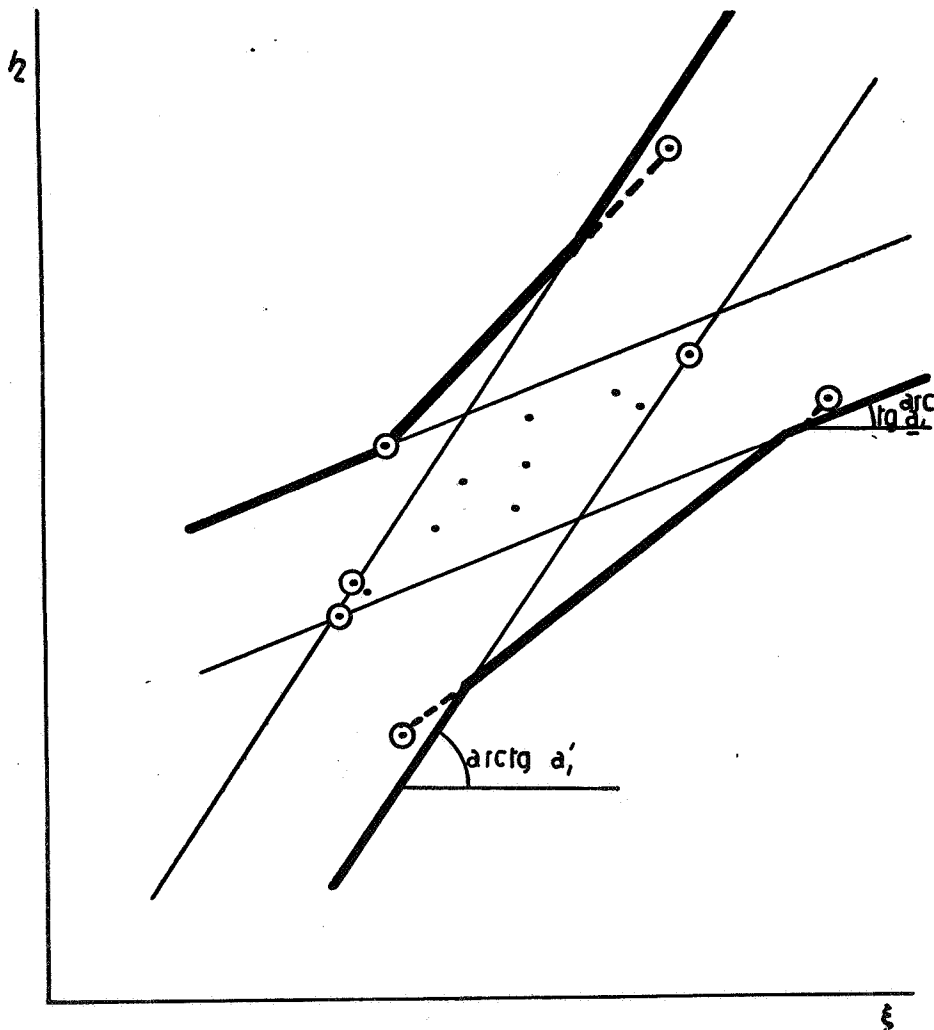


Fig. 5.  $n = 15$   $s = 2$



## 6. CONCLUDING REMARKS

6.0. The methods of determining confidence regions which may be derived from this kind of analysis have not been exhaustively treated. In order to elucidate this statement we shall give a confidence interval for  $\alpha$  in the stochastic regression equation

$$\eta_i = A \xi_i^\alpha + w_i, \quad (i = 1, \dots, n)$$

in which all  $\xi_i$  are positive, and in which

$$P[A \eta_i \leq 0] = 0$$

holds for  $i = 1, \dots, n$ .  $\xi_1, \dots, \xi_n$  are known,  $A$  and  $\alpha$  are unknown parameters, and  $w_1, \dots, w_n$  are random variables, which are supposed (1) to be distributed stochastically independent, (2) to have continuous symmetrical distribution functions with zero median.

We arrange the observed points  $(\xi_i, \eta_i)$  according to increasing magnitude of  $\xi$  and define

$$d_i = \frac{\log |\eta_{n_1+i}| - \log |\eta_i|}{\log \xi_{n_1+i} - \log \xi_i}, \quad (i = 1, \dots, n_1)$$

in which  $n_1 = \frac{1}{2}n$  (if  $n$  is odd the point  $(\xi_{i(n+1)}, \eta_{i(n+1)})$  is neglected). After arranging the observed quantities  $d_i$  according to increasing magnitude:

$$d_{(1)} < \dots < d_{(n_1)}$$

we have the following

*Theorem 10.* Under conditions (1) and (2) the interval  $(d_{(r_1)}, d_{(n_1-r_1+1)})$  is a confidence interval for  $\alpha$  to the level of significance  $2I_1(n_1 - r_1 + 1, r_1)$ .

*Proof.* We have

$$\frac{\eta_{n_1+i} - w_{n_1+i}}{\eta_i - w_i} = \left( \frac{\xi_{n_1+i}}{\xi_i} \right)^\alpha$$

or:

$$\eta_{n_1+i} - \left( \frac{\xi_{n_1+i}}{\xi_i} \right)^\alpha \eta_i = w_{n_1+i} - \left( \frac{\xi_{n_1+i}}{\xi_i} \right)^\alpha w_i.$$

It follows from condition (2), that

$$w_{n_1+i} - \left( \frac{\xi_{n_1+i}}{\xi_i} \right)^\alpha w_i$$

has a continuous distribution function with zero median. Hence:

$$\begin{aligned} P \left[ \eta_{n_1+i} - \left( \frac{\xi_{n_1+i}}{\xi_i} \right)^\alpha \eta_i < 0 \right] &= \\ &= P [\log \eta_{n_1+i} - \log \eta_i < \alpha (\log \xi_{n_1+i} - \log \xi_i)] = \\ &= P \left[ \frac{\log \eta_{n_1+i} - \log \eta_i}{\log \xi_{n_1+i} - \log \xi_i} < \alpha \right] = P \left[ \frac{\log \eta_{n_1+i} - \log \eta_i}{\log \xi_{n_1+i} - \log \xi_i} > \alpha \right] = \frac{1}{2}, \end{aligned}$$

if  $\eta_1, \dots, \eta_n$  are positive; if they are negative we have to replace  $\eta_i$  and  $\eta_{m_1+i}$  by  $-\eta_i$  and  $-\eta_{m_1+i}$  respectively. From this and from condition (1) the theorem follows.

The theorem shows that this method of determining a confidence interval for  $\alpha$  is identical with the incomplete method for  $\alpha$  in the linear equation

$$\log \eta_i = \log A + a \log \xi_i + w'_i$$

(which can be written as

$$\eta_i = A \xi_i^a e^{w'_i},$$

if  $w'_1, \dots, w'_n$  satisfy the same conditions (1) and (2).

6. 1. Finally we mention that it is possible to find estimates instead of confidence intervals. Consider e.g. the statistics  $\Delta(ij)$ ; each of these  $\binom{n}{2}$  statistics has the property that its sampling median is equal to  $\alpha_1$  (cf. section 1. 3.). Hence one can use the sample median of the observed quantities  $\Delta(ij)$  as an estimate of  $\alpha_1$ .

It is a pleasure to acknowledge my indebtedness to Professor Dr D. VAN DANTZIG for his stimulating interest and to Mr J. HEMELRIJK for his valuable and constructive criticism.

#### REFERENCES

- BENTZEL, R. and H. WOLD, On statistical demand analysis from the viewpoint of simultaneous equations. *Skand. Aktuarietidskr.* 29, 95—114 (1946).
- GIRSHICK, M. A. and T. HAAVELMO, Statistical analysis of the demand for food: examples of simultaneous estimation of structural equations. *Econometrica*, 15, 79—110 (1947).
- HAAVELMO, T., The statistical implications of a system of simultaneous equations. *Econometrica*, 11, 1—12 (1943).
- , The probability approach in econometrics. *Econometrica*, 12, suppl. (1944).
- KOOPMANS, T., Statistical estimation of simultaneous economic relations. *Journal of the Amer. Statist. Assoc.*, 40, 448—466 (1945).
- , (ed.), *Statistical inference in dynamic economic models* (New York, 1950).
- THOMPSON, W. R., On confidence ranges for the median and other expectation distributions for populations of unknown distribution form. *Annals of Math. Statist.*, 7, 122—128 (1936).

*Publication of the Statistical Department of the  
"Mathematisch Centrum", Amsterdam*