

MATHEMATICS

A RANK-INVARIANT METHOD OF LINEAR AND POLYNOMIAL REGRESSION ANALYSIS

II ¹⁾

BY

H. THEIL

(Communicated by Prof. D. VAN DANTZIG at the meeting of March 25, 1950)

2. CONFIDENCE REGIONS FOR THE PARAMETERS OF LINEAR REGRESSION EQUATIONS IN THREE AND MORE VARIABLES.

The probability set.

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

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

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

In this first place we consider $n(\nu + 2)$ random variables $\mathbf{u}_{\lambda i}, \mathbf{v}_i, \mathbf{w}_i$ ($\lambda = 1, \dots, \nu; i = 1, \dots, n$). Furthermore we consider $(n + 1)\nu + 1$ parameters $a_0, a_\lambda, \xi_{\lambda i}$ ($i = 1, \dots, n; \lambda = 1, \dots, \nu$) and put:

$$\left. \begin{aligned} (5) \quad \theta_i &= a_0 + \sum_{\lambda=1}^{\nu} a_\lambda \xi_{\lambda i} \\ (6) \quad \eta_i &= \theta_i + \mathbf{w}_i \\ (7) \quad \mathbf{x}_{\lambda i} &= \xi_{\lambda i} + \mathbf{u}_{\lambda i} \\ (8) \quad \mathbf{y}_i &= \eta_i + \mathbf{v}_i \end{aligned} \right\} \begin{cases} i = 1, \dots, n \\ \lambda = 1, \dots, \nu. \end{cases}$$

So the variables $\mathbf{x}_{\lambda i}$ and \mathbf{y}_i have a simultaneous distribution on Γ , and are therefore random variables.

We call $\xi_{\lambda i}$ the parameter values of the variable ξ_λ . The equation (5) is the multiple regression equation. The random variables \mathbf{w}_i are called "the true deviations from linearity", while the random variables $\mathbf{u}_{\lambda i}$ and \mathbf{v}_i are called "the errors of observation" of the values $\xi_{\lambda i}$ and η_i respectively.

¹⁾ This paper is the second of a series of papers, the first of which appeared in these Proceedings, 53, 386–392 (1950).

Putting

$$\mathbf{z}_i = - \sum_{\lambda=1}^{\nu} \alpha_{\lambda} \mathbf{u}_{\lambda i} + \mathbf{v}_i + \mathbf{w}_i$$

we have

$$\mathbf{y}_i = \alpha_0 + \sum_{\lambda=1}^{\nu} \alpha_{\lambda} \mathbf{x}_{\lambda i} + \mathbf{z}_i,$$

the random variables \mathbf{z}_i being called “the apparent deviations from linearity”.

Confidence regions for $\alpha_0, \alpha_1, \dots, \alpha_{\nu}$.

2. 1. In order to give confidence regions for the $(\nu + 1)$ parameters $\alpha_0, \alpha_{\lambda}$ ($\lambda = 1, \dots, \nu$) we impose the following *conditions*:

Condition I: The $n(\nu + 2)$ -uples $(\mathbf{u}_{1i}, \dots, \mathbf{u}_{\nu i}, \mathbf{v}_i, \mathbf{w}_i)$ are stochastically independent.

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

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

Furthermore we impose for the *incomplete method* to be mentioned:

Condition III:

$$P[\mathbf{z}_i < \mathbf{z}_j] = P[\mathbf{z}_i > \mathbf{z}_j] = \frac{1}{2} \text{ for } i \neq j$$

and for the *complete method*:

Condition IIIa: Each \mathbf{z}_i has the same continuous distribution function.

2. 2. Secondly we *define* the following quantities:

$$\begin{aligned} \mathbf{G}^{(\lambda')} (i) &= \mathbf{y}_i - \sum_{\substack{\lambda=1 \\ \lambda \neq \lambda'}}^{\nu} \alpha_{\lambda} \mathbf{x}_{\lambda i} = \\ &= \alpha_0 + \alpha_{\lambda'} \mathbf{x}_{\lambda' i} + \mathbf{z}_i \quad (\lambda' = 1, \dots, \nu; i = 1, \dots, n). \end{aligned}$$

Furthermore, after arranging the n observed points $(y_i, x_{1i}, \dots, x_{\nu i})$ according to increasing values of $x_{\lambda'}$ (which, by condition II, is identical with the arrangement according to increasing values of $\xi_{\lambda'}$):

$$x_{\lambda'1} < x_{\lambda'2} < \dots < x_{\lambda'n}$$

we define the quantities

$$\begin{aligned} \mathbf{K}^{(\lambda')} (i, j) &= \frac{\mathbf{G}^{(\lambda')} (i) - \mathbf{G}^{(\lambda')} (j)}{\mathbf{x}_{\lambda' i} - \mathbf{x}_{\lambda' j}} = \\ &= \frac{\mathbf{y}_i - \mathbf{y}_j}{\mathbf{x}_{\lambda' i} - \mathbf{x}_{\lambda' j}} - \sum_{\substack{\lambda=1 \\ \lambda \neq \lambda'}}^{\nu} \alpha_{\lambda} \frac{\mathbf{x}_{\lambda i} - \mathbf{x}_{\lambda j}}{\mathbf{x}_{\lambda' i} - \mathbf{x}_{\lambda' j}} = \\ &= \alpha_{\lambda'} + \frac{\mathbf{z}_i - \mathbf{z}_j}{\mathbf{x}_{\lambda' i} - \mathbf{x}_{\lambda' j}} \quad (i = 1, \dots, n-1; j = i+1, \dots, n). \end{aligned}$$

For any set of values $\alpha_1, \dots, \alpha_{\lambda'-1}, \alpha_{\lambda'+1}, \dots, \alpha_v$ we arrange the quantities $K^{(\lambda')}(i, j)$ according to increasing magnitude; we define $K_i^{(\lambda')}$ as the quantity with rank i in this arrangement:

$$K_1^{(\lambda')} < K_2^{(\lambda')} < \dots < K_{\binom{n}{2}}^{(\lambda')}$$

Finally we define the intervals $I_{\lambda'}(\alpha_1, \dots, \alpha_{\lambda'-1}, \alpha_{\lambda'+1}, \dots, \alpha_v)$ as the intervals

$$\left(K_q^{(\lambda')}, K_{\binom{n}{2}-q+1}^{(\lambda')} \right)$$

with $2q \leq \binom{n}{2}$; $A_{\lambda'}$ as the union of

$$I_{\lambda'}(\alpha_1, \dots, \alpha_{\lambda'-1}, \alpha_{\lambda'+1}, \dots, \alpha_v) \text{ for all } \alpha_{\lambda'} (\lambda = 1, \dots, v; \lambda \neq \lambda');$$

and A as the union of all $A_{\lambda'} (\lambda' = 1, \dots, v)$.

2.3. We have the following theorem concerning the *complete method* for three and more variables:

Theorem 4: Under conditions I, II and IIIa the region A is a confidence region for the parameters $\alpha_1, \dots, \alpha_v$, the level of significance being $\leq 2v \cdot P[q-1 | n]^2$.

Proof: If the set of assumed parameters values $\alpha_1, \dots, \alpha_{\lambda'-1}, \alpha_{\lambda'+1}, \dots, \alpha_v$ is the "true" set, it follows from the analysis in section 1.3., that $I_{\lambda'}(\alpha_1, \dots, \alpha_{\lambda'-1}, \alpha_{\lambda'+1}, \dots, \alpha_v)$ is a confidence interval for $\alpha_{\lambda'}$ to the level of significance $2P[q-1 | n]$. Hence it follows that if $(\alpha_1, \dots, \alpha_v)$ represents the "true" point in the $\alpha_1, \dots, \alpha_v$ -space, we have

$$P[(\alpha_1, \dots, \alpha_v) \in A_{\lambda'}] = 1 - 2P[q-1 | n], \quad (\lambda' = 1, \dots, v),$$

which proves the theorem.

2.4. If condition III (but not necessarily IIIa) is fulfilled, the method mentioned above can be replaced by the following one. We replace the quantities

$$K^{(\lambda')}(i, j) \quad (\lambda' = 1, \dots, v; i = 1, \dots, n-1; j = i+1, \dots, n)$$

by

$$K^{(\lambda')}(i, n_1 + i) \quad (\lambda' = 1, \dots, v; i = 1, \dots, n_1).^3$$

The intervals $I'_{\lambda'}(\alpha_1, \dots, \alpha_{\lambda'-1}, \alpha_{\lambda'+1}, \dots, \alpha_v)$ are now defined as the intervals bounded by the values of $K^{(\lambda')}(i, n_1 + i)$ with rank r_1 and $(n_1 - r_1 + 1)$ respectively, if they are arranged in ascending order; whereas the definitions of $A'_{\lambda'}$ as the union of all $I'_{\lambda'}$ and of A' as the union of all $A'_{\lambda'}$,

²⁾ For the definition of $P[q-1 | n]$ the reader is referred to section 1.3. (part I of this paper).

³⁾ $n_1 = \frac{1}{2}n$. Cf. section 1.2.

remain unchanged. The following theorem of the *incomplete method* for three and more variables will now be obvious from the analysis of section 1. 1.:

Theorem 5. Under conditions I, II and III the region \mathbf{A}' is a confidence region for the parameters $\alpha_1, \dots, \alpha_v$, the level of significance being $\leq 2v \cdot I_{\frac{1}{2}}(r_1, n_1 - r_1 + 1)$.

2. 5. A confidence region for the parameters $\alpha_0, \alpha_1, \dots, \alpha_v$ can be constructed, if the median of \mathbf{z}_i is known, e.g. if the following condition is fulfilled:

Condition IV: The median of each \mathbf{z}_i is zero.

The method for the construction of this confidence region is analogous to the one given in section 1. 2.

An illustration for the special case $v = 2$.

2. 6. The form of the region \mathbf{A}_λ or \mathbf{A}'_λ will now be indicated for the case of three variables:

$$y_i = \alpha_0 + \alpha_1 x_{1i} + \alpha_2 x_{2i} + z_i.$$

Using the incomplete method we find n_1 functions of α_2 :

$$K^{(1)}(i, n_1 + i) = \frac{y_i - y_{n_1 + i}}{x_{1i} - x_{1, n_1 + i}} - \alpha_2 \frac{x_{2i} - x_{2, n_1 + i}}{x_{1i} - x_{1, n_1 + i}},$$

which are estimates of α_1 , given α_2 . They are represented by straight lines in the α_1, α_2 -plane. For any value of α_2 we can arrange these quantities in ascending order. As long as (under continuous variation of α_2) the numbers i_1 and i_2 for which the statistics $K^{(1)}(i_1, n_1 + i_1)$ and $K^{(1)}(i_2, n_1 + i_2)$

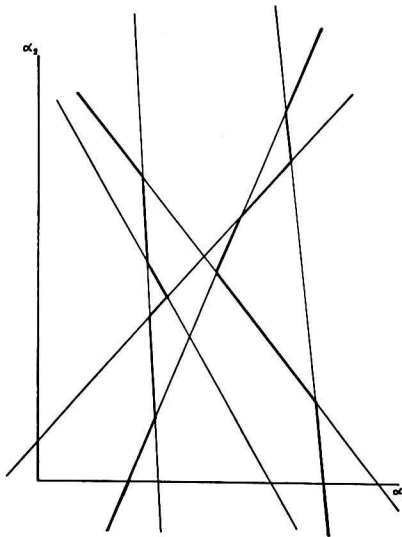


Fig. 1. $n_1 = 6, r_1 = 2$.

have the r_1 -th and $(n_1 - r_1 + 1)$ -th rank according to increasing order (with r_1 as defined in section 2. 4.) remain constant, the extreme points of the confidence intervals vary along straight lines. If, when passing some value a_2^* of a_2 either i_1 or i_2 changes, the corresponding straight line passes into another one, intersecting the first one in a point with $a_2 = a_2^*$.

So a diagram can be constructed, in which the n_1 straight lines are drawn in the a_1, a_2 -plane. This gives the stochastic region \mathbf{A}'_1 depending on the given observations and bounded to the left and to the right by broken lines.

According to Theorem 5 it contains the true point (a_1, a_2) with the probability

$$1 - 2 I_{\frac{1}{2}}(r_1, n_1 - r_1 + 1).$$

The region \mathbf{A}'_2 , bounded above and below, can be constructed in a similar way; then the observed points must be arranged in ascending order of x_2 .

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