

MATHEMATICS

A RANK-INVARIANT METHOD OF LINEAR AND POLYNOMIAL REGRESSION ANALYSIS

I

BY

H. THEIL

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

0. INTRODUCTION

0.0 Regression analysis is usually carried out under the hypothesis that one of the variables is normally distributed with constant variance, its mean being a function of the other variables. This assumption is not always satisfied, and in most cases difficult to ascertain.

In recent years attention has been paid to problems of estimating the parameters of regression equations under more general conditions (see the references at the end of this paper: A. WALD (1940), K. R. NAIR and M. P. SHRIVASTAVA (1942), K. R. NAIR and K. S. BANERJEE (1942), G. W. HOUSNER and J. F. BRENNAN (1948) and M. S. BARTLETT (1949)). Confidence regions, however, were obtained under the assumption of normality only; to obtain these without this assumption will be the main object of this paper.

0.1. In section 1. confidence regions will be given for the parameters of linear regression equations in two variables. In the sequel of this paper we hope to deal with equations in more variables, polynomial equations, systems of equations and problems of prediction.

1. CONFIDENCE REGIONS FOR THE PARAMETERS OF LINEAR REGRESSION EQUATIONS IN TWO VARIABLES

The probability set.

1.0. Throughout this section the probability set Γ ("Wahrscheinlichkeitsfeld" in the sense of A. KOLMOGOROFF) underlying the probability statements will be the $3n$ -dimensional Cartesian space R_{3n} with coordinates $u_1, \dots, u_n, v_1, \dots, v_n, w_1, \dots, w_n$. Every random variable mentioned is supposed to be defined on this probability set.

In the first place we suppose $3n$ random variables $\mathbf{u}_i, \mathbf{v}_i, \mathbf{w}_i$ ($i = 1, \dots, n$)¹⁾ to be defined on Γ , i.e. we suppose $\mathbf{u}_i, \mathbf{v}_i, \mathbf{w}_i$ to have a simultaneous probability distribution on Γ .

¹⁾ The distinction between a stochastic variable and the value it takes in a given observation (or system of observations) will be indicated by bold type for the former one.

If we now put:

$$\left. \begin{aligned} (1) \quad \theta_i &= a_0 + a_1 \xi_i \\ (2) \quad \eta_i &= \theta_i + \mathbf{w}_i \\ (3) \quad \mathbf{x}_i &= \xi_i + \mathbf{u}_i \\ (4) \quad \mathbf{y}_i &= \eta_i + \mathbf{v}_i \end{aligned} \right\} i = 1, \dots, n$$

then, for any set of values of the $(n+2)$ parameters ξ_i , a_0 and a_1 , the variables \mathbf{x}_i and \mathbf{y}_i have a simultaneous distribution on Γ , and are therefore random variables.

We shall call ξ_i the parameter values of the variable ξ . The equation (1) is the regression equation; this equation contains no stochastic variables. Furthermore we shall call \mathbf{w}_i the "true deviations from linearity"; hence the variable η is a linear function of ξ , but for the deviations \mathbf{w} . Finally \mathbf{u}_i and \mathbf{v}_i are called the "errors of observation" of the "true" values ξ_i and η_i respectively.

The problem then is, under certain conditions for the probability distribution of \mathbf{u}_i , \mathbf{v}_i , \mathbf{w}_i , to determine confidence intervals for the parameters a_0 and a_1 , given a sequence of observations $x_1, \dots, x_n, y_1, \dots, y_n$ of the random variables $\mathbf{x}_1, \dots, \mathbf{x}_n, \mathbf{y}_1, \dots, \mathbf{y}_n$.

Incomplete method: confidence interval for a_1 .²⁾

1.1. We suppose that the following conditions are satisfied:

Condition I: The n triples $(\mathbf{u}_i, \mathbf{v}_i, \mathbf{w}_i)$ are stochastically independent.

Condition II: 1. Each of the errors \mathbf{u}_i vanishes outside a finite interval $|\mathbf{u}_i| \leq g_i$.

2. For each $i \neq j$ we have: $|\xi_i - \xi_j| > g_i + g_j$.

From condition II it follows that either

$$P[\mathbf{x}_i < \mathbf{x}_j] = 1 \quad \text{and} \quad \xi_i < \xi_j$$

or

$$P[\mathbf{x}_i > \mathbf{x}_j] = 1 \quad \text{and} \quad \xi_i > \xi_j.$$

This condition means that the errors \mathbf{u}_i are sufficiently small in order that arrangement of the observed values x_i according to increasing magnitude be identical with the arrangement according to increasing values of ξ_i (cf. also A. WALD (1940), p. 294, seq., where a similar (weaker) condition is imposed). The arrangement of the \mathbf{x}_i is therefore uniquely determined. We therefore suppose the \mathbf{x}_i as well as the ξ_i to be arranged according to increasing order.

Put $n_1 = n - [\frac{1}{2}n]$; if n is odd, the observation with rank $\frac{1}{2}(n+1)$ is not used. We therefore omit this observation and write $n = 2n_1$.

²⁾ The author is indebted to Mr J. HEMELRIJK for his constructive criticism concerning some points of this section.

We determine the following n_1 statistics:

$$\Delta(i, n_1 + i) = \frac{y_{n_1+i} - y_i}{x_{n_1+i} - x_i} = \alpha_1 + \frac{z_{n_1+i} - z_i}{x_{n_1+i} - x_i},$$

in which $z_i = -\alpha_1 u_i + v_i + w_i$.

We now impose:

Condition III, which states:

$$P[z_i < z_{n_1+i}] = P[z_i > z_{n_1+i}] = \frac{1}{2}.$$

As all denominators $x_{n_1+i} - x_i$ are positive, it follows that

$$P[\Delta(i, n_1 + i) < \alpha_1] = P[\Delta(i, n_1 + i) > \alpha_1] = \frac{1}{2},$$

i.e. that $\Delta(i, n_1 + i)$ has a median α_1 and that its distribution function is continuous in the median.

The following conditions IIIa and IIIb are each sufficient in order that $P[z_i < z_{n_1+i}] = P[z_i > z_{n_1+i}] = \frac{1}{2}$:

Condition IIIa: the random variables z_i ($i = 1, \dots, n$) have the same continuous distribution function.

Condition IIIb: the random variables z_i have continuous distribution functions which are symmetrical with equal medians $\text{med}(\mathbf{z})$.

Proof: In case IIIa the simultaneous distribution of z_i and z_{n_1+i} is symmetrical about the line $z_i = z_{n_1+i}$, which proves the statement. In case IIIb it is symmetrical about the lines $z_i = \text{med}(\mathbf{z})$ and $z_{n_1+i} = \text{med}(\mathbf{z})$; hence the simultaneous distribution of $z_i - \text{med}(\mathbf{z})$ and $z_{n_1+i} - \text{med}(\mathbf{z})$ is symmetrical with respect to the origin, which proves the statement.

We now arrange the n_1 statistics $\Delta(i, n_1 + i)$ in increasing order:

$$\Delta_1 < \Delta_2 < \dots < \Delta_{n_1},$$

in which

$$\Delta_j = \Delta(i_j, n_1 + i_j).$$

The probability that exactly r among the n_1 values $\Delta(i, n_1 + i)$ are $< \alpha_1$, i.e. that $\Delta_r < \alpha_1 < \Delta_{r+1}$, is $2^{-n_1} \binom{n_1}{r}$ because of the conditions I and III. Hence:

$$\begin{aligned} P[\Delta_{r_1} \leq \alpha_1 \leq \Delta_{n_1-r_1+1} | \alpha_1] &= \\ &= 1 - 2^{1-n_1} \sum_{s=0}^{r_1-1} \binom{n_1}{s} \\ &= 1 - 2 I_1(r_1, n_1 - r_1 + 1) \end{aligned}$$

in which

$$I_1(r_1, n_1 - r_1 + 1) = \frac{\int_0^{\frac{1}{2}} x^{r_1-1} (1-x)^{n_1-r_1} dx}{\int_0^1 x^{r_1-1} (1-x)^{n_1-r_1} dx}$$

is the incomplete Beta-function for the argument $\frac{1}{2}$.

So we have proved:

Theorem I: under conditions I, II and III a confidence interval for α_1 is given by the largest but $(r_1 - 1)$ and the smallest but $(r_1 - 1)$ among the values $\Delta(i, n_1 + i)$, the level of significance being $2 I_1(r_1, n_1 - r_1 + 1)$.

We shall call this method an "incomplete method" because a limited use is made of the $\binom{n}{2}$ statistics

$$\Delta(ij) = \frac{y_i - y_j}{x_i - x_j} \quad (i < j).$$

Incomplete method: confidence region for α_0 and α_1 .

1.2. If the median of \mathbf{z}_i ($i = 1, \dots, n$) is numerically known, a confidence region for α_0 and α_1 can be found. We suppose that the following condition is satisfied:

Condition IV: the median of each \mathbf{z}_i ($i = 1, \dots, n$) is zero:

$$P[y_i - \alpha_1 x_i > \alpha_0] = P[y_i - \alpha_1 x_i < \alpha_0] = \frac{1}{2}.$$

For any value of α_1 we can arrange the n quantities $Z_i = y_i - \alpha_1 x_i$ according to increasing magnitude:

$$Z_1(\alpha_1) < Z_2(\alpha_1) < \dots < Z_n(\alpha_1).$$

Under the condition that α_1 has the value used in this arrangement, we can state that

$$\begin{aligned} P[\alpha_0 \in (Z_{r_0}(\alpha_1), Z_{n-r_0+1}(\alpha_1)) \mid \alpha_0, \alpha_1] &= \\ &= 1 - 2 I_1(r_0, n - r_0 + 1) = 1 - \varepsilon_0. \end{aligned}$$

On the other hand, if we write I_1 for the interval $(\Delta_{r_1}, \Delta_{n_1-r_1+1})$, we can state:

$$\begin{aligned} P[\alpha_1 \in I_1 \mid \alpha_1] &= \\ &= 1 - 2 I_1(r_1, n_1 - r_1 + 1) = 1 - \varepsilon_1. \end{aligned}$$

If we denote by I_0 the interval bounded by the lowest of the values $Z_{r_0}(\alpha_1)$ if α_1 varies through I_1 and by the largest of the values $Z_{n-r_0+1}(\alpha_1)$ if α_1 varies through I_1 we have

$$P[\alpha_0 \in I_0 \wedge \alpha_1 \in I_1 \mid \alpha_0, \alpha_1] \geq (1 - \varepsilon_0)(1 - \varepsilon_1).$$

So we have proved:

Theorem 2: under conditions I, II, III and IV a rectangular confidence region in the α_0, α_1 - plane is given by the intervals $\alpha_0 \in I_0$ and $\alpha_1 \in I_1$, the level of significance being $\leq \varepsilon_0 + \varepsilon_1 - \varepsilon_0 \varepsilon_1$.

If all observed points (x_i, y_i) obey the inequality $x_i \geq 0$, all quantities $y_i - \alpha_1 x_i$ are decreasing functions of α_1 . It follows that I_0 is bounded by $Z_{r_0}(\Delta_{n_1-r_1+1})$ and by $Z_{n-r_0+1}(\Delta_{r_1})$. The converse holds if every point satisfies the inequality $x_i \leq 0$.

Complete method.

1.3. We suppose that the conditions I, II and IIIa are satisfied and consider two arrangements of the points (x_i, y_i) : the arrangement according to increasing values of x and that according to $z = y - a_0 - a_1x$.

The arrangement according to z is possible for any assumed value of a_1 . The hypothesis that this value is the true one is rejected if and only if there is a significant rank correlation between the arrangements.

Consider the statistics

$$\Delta(ij) = \frac{y_i - y_j}{x_i - x_j} = a_1 + \frac{z_i - z_j}{x_i - x_j},$$

in which $i < j$, so that (if the ordering is according to x) $x_i < x_j$ and $\xi_i < \xi_j$. It follows that $\Delta(ij) > a_1$, if and only if $z_i < z_j$.

Now, under the null hypothesis that the arrangements of the points according to x and according to z are independent, the distribution of Kendall's "rank correlation coefficient"

$$\frac{S}{\binom{n}{2}}$$

is known, in which S is the number of cases in which the ordering according to z is the same as the ordering according to x ($z_k < z_{k'}$, and $x_k < x_{k'}$) minus the number of cases in which the ordering according to z is the inverse as compared with the one according to x ($z_k > z_{k'}$, while $x_k < x_{k'}$).

For any value of a_1 the number of cases $z_i > z_j$ can be found. Suppose this to be q ; it will be clear that

$$S = \binom{n}{2} - 2q.$$

The distribution function of S for any value of n has been given by M. G. Kendall (see M. G. KENDALL (1947), p. 403–407 and (1948), p. 55–62) by means of a recurrence formula. So the probability $P[q|n]$ that $q' \leq q$ cases $z_i > z_j$ are found can be determined. If this probability is below the level of significance chosen, we reject the hypothesis that a_1 has the value used in the arrangement according to z .

Hence, if we arrange the statistics $\Delta(ij)$ in increasing order:

$$\Delta_1 < \Delta_2 < \dots < \Delta_{\binom{n}{2}}$$

we find by symmetry

$$P[\Delta_q \leq a_1 \leq \Delta_{\binom{n}{2}-q+1} | a_1] = 1 - 2P[q-1 | n]$$

so that we have proved:

Theorem 3: under conditions I, II and IIIa a confidence interval for a_1 is given by the largest but $(q-1)$ and the smallest but $(q-1)$ among the values $\Delta(ij)$, the level of significance being $2P[q-1|n]$.

The method of 1.2. can be applied here to find a simultaneous confidence region for α_0 and α_1 , I_1 now being the interval $(\Delta_q, \Delta_{\binom{n}{2}-q+1})$.

A comparison.

1.4. The second method may be called a "complete method", because all statistics $\Delta(i, j)$ are used. It requires only 5 points in order to reach the level of significance 0,05 whereas the limited method needs 12 points. However, if the number of points is large, the computational labor of the complete method is considerably greater than that of the incomplete method. Moreover, the conditions under which the complete method is valid are more stringent; the fact that the set of conditions I, II and III is sufficient for the incomplete method is important in view of the general occurrence of "heteroscedastic" distribution, i.e. distributions in which the variance (if finite) of η_i is larger for higher values of ξ than for lower ones if $\alpha_1 > 0$ and conversely if $\alpha_1 < 0$.

Testing linearity.

1.5. Suppose that the set of conditions I, II and IIIa is valid. Then the hypothesis that the regression curve for two variables is linear can be tested against the alternative composite hypothesis that it is either positive- or negative-convex,³⁾ i.e. in the set of equations (1), (2), (3), (4) the equation $\theta_i = \alpha_0 + \alpha_1 \xi_i$ is tested against any equation $\theta_i = \theta(\xi_i)$ with

either
$$\frac{d^2\theta}{d\xi^2} > 0 \quad \text{for all } \xi$$

or
$$\frac{d^2\theta}{d\xi^2} < 0 \quad \text{for all } \xi,$$

the equations (2), (3), (4) remaining unchanged.

Consider the n_1 statistics

$$\Delta(1, n_1 + 1), \dots, \Delta(n_1, 2n_1)$$

in this arrangement. If this ordering has a significant rank correlation with the ordering of these statistics according to increasing magnitude, we reject the hypothesis that the regression curve is linear.

REFERENCES

- BARTLETT, M. S., Fitting a straight line when both variables are subject to error. *Biometrics* 5, 207–212 (1949).
 DANTZIG, D. VAN, Capita selecta der waarschijnlijkheidsrekening, caput II, (stenciled) (1947).

³⁾ A function $f(x)$ is positive-convex (cf. e.g. D. VAN DANTZIG, 93–94 (1947)) in an interval if for every x_1 and x_2 of this interval and for every real positive number $a < 1$ the following inequality is satisfied

$$a f(x_1) + (1-a) f(x_2) > f(ax_1 + \overline{1-a} x_2).$$

- HOUSNER, G. W. and J. F. BRENNAN, The estimation of linear trends. *Annals of Mathematical Statistics*, 19, 380—388 (1948).
- KENDALL, M. G., *The advanced theory of statistics*, London, 1, 3rd edition (1947).
- , *Rank correlation methods*, London (1948).
- NAIR, K. R. and K. S. BANERJEE, A note on fitting of straight lines if both variables are subject to error. *Sankhya*, 6, 331 (1942).
- and M. P. SHRIVASTAVA, On a simple method of curve fitting. *Sankhya*, 6, 121—132 (1942).
- WALD, A., The fitting of straight lines if both variables are subject to error. *Annals of mathematical statistics*, 11, 284—300 (1940).

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