

December 16, 1947

Remarks on Estimation

by George W. Rasch

1. In the problem of estimating α in the stochastic difference equation

$$x_t = \alpha x_{t-1} + u_t$$

where x_0 is fixed and u_t normally and independently distributed about 0 with a constant variance σ^2 we get 3 statistics which are together sufficient for the 2 parameters. Similar difference equations of second order contain 3 parameters for which we derive a set of 6 sufficient statistics. And so on.

Thus we are faced with the problem: How to utilize k statistics for the estimation of $l < k$ parameters? The following is an attempt to formulate this problem for $l=1$ and solve it approximately under the assumption that the variances of the statistics are fairly small as is more often the case for medium-sized samples. It should be noted, however, that the principle may be extended to more parameters.

The problem we are going to solve approximately may be stated in the following terms: Provided the vector $x = (x_1, \dots, x_k)$ follows some k -dimensional distribution with one parameter θ :

$$(1) \quad \phi\{x_1, \dots, x_k | \theta\}$$

we want an unbiased estimate of θ with a minimum variance.

Denote the mean vector and the variance matrix of x by

$$(2) \quad E\{x\} = E(\theta), \text{ resp. } V\{x\} = M(\theta).$$

Then if $f(x)$ is some reasonable function and if ϕ is some reasonable distribution we have with a tolerable approximation

$$(3) \quad E\{f(x)\} = f(E(\theta))$$

and

$$(4) \quad V\{f(x)\} = \left(\frac{\partial f}{\partial E_1}, \dots, \frac{\partial f}{\partial E_k} \right) M(\theta) \begin{pmatrix} \frac{\partial f}{\partial E_1} \\ \vdots \\ \frac{\partial f}{\partial E_k} \end{pmatrix} \\ = \sum \mu_{ij}(\theta) \frac{\partial f}{\partial E_i} \cdot \frac{\partial f}{\partial E_j}$$

For

$$(5) \quad t = f(\lambda)$$

to be an unbiased estimate of θ we therefore must have

$$(6) \quad f(\bar{E}(\theta)) = \theta$$

which on differentiation gives

$$(7) \quad \sum \frac{\partial f}{\partial \bar{E}_i} \cdot \bar{E}'_i(\theta) = 1.$$

To the degree of approximation we are working with our problem is a minimization of (4) under the condition (7). Putting for a moment

$$(8) \quad \frac{\partial f}{\partial \bar{E}_i} = \lambda_i$$

we differentiate

$$(9) \quad \sum \mu_{ij} \lambda_i \lambda_j - 2\kappa (\sum \bar{E}'_i \lambda_i - 1)$$

with respect to the λ'_i and equate to 0:

$$(10) \quad \sum \mu_{ij} \lambda_j = \kappa \cdot \bar{E}'_i$$

or¹

$$(11) \quad M \lambda^* = \kappa \bar{E}'^*$$

i.e.

$$(12) \quad \lambda^* = \kappa \bar{E}' M^{-1}$$

where κ from

$$(13) \quad \bar{E}' \lambda^* = 1$$

is determined as

$$(14) \quad \kappa = 1 / \bar{E}' M^{-1} \bar{E}'^*$$

The minimum value of (4) is found to be

$$(15) \quad V\{t\} = \lambda^* M \lambda^* = \kappa^2 \bar{E}' M^{-1} \bar{E}'^* = \kappa.$$

Now we are faced with the problem of determining such a function $f(\lambda)$

that

$$(16) \quad f(\bar{E}(\theta)) = \theta$$

and

$$(17) \quad f'_i(\theta) = \left(\frac{\partial f}{\partial \bar{E}_i} \right)_{\lambda = \bar{E}} = \lambda_i(\theta)$$

1 -- To avoid confusion with the differentiation sign we shall in this paper use * for transposition.

where

$$(18) \quad \lambda_i(\theta) = \frac{\sum \mu_{ij}^{(-1)}(\theta) \xi_j'(\theta)}{\sum \sum \mu_{ij}^{(-1)}(\theta) \xi_i'(\theta) \xi_j'(\theta)}$$

are known functions.

A solution is readily obtained if we consider the first terms of the Taylor expansion of $f(x)$ about $\xi(z)$:

$$(19) \quad \begin{cases} f(x) = f(\xi(z)) + (x - \xi(z)) \left(\frac{f}{f_k}(\xi(z)) \right) + \text{terms of second order} \\ = z + \sum (\pi_j - \xi_j(z)) \lambda_j(z) + \dots \end{cases}$$

Take in this expansion z as some function of x :

$$(20) \quad z = h(x)$$

and differentiate (19) with respect to x_i . Then we get

$$(21) \quad \frac{\partial f}{\partial x_i} = \frac{\partial h}{\partial x_i} + \lambda_i(z) - \sum \xi_j'(z) \lambda_j(z) \frac{\partial z}{\partial x_i} + \text{terms of first order}$$

which according to (13) reduces to

$$(22) \quad \frac{\partial f}{\partial x_i} = \lambda_i(z) + \text{terms of first order.}$$

If we now choose $h(x)$ such that

$$(23) \quad h(\xi(z)) = z$$

and put

$$(24) \quad f(x) = z + \sum (\pi_j - \xi_j(z)) \lambda_j(z), \quad z = h(x),$$

then we have a solution of our problem.

One way of choosing $h(x)$ is to invert the functions $\xi_1(\theta), \dots, \xi_k(\theta)$:

$$(25) \quad \eta_i(\xi_i(\theta)) = \theta$$

and take a linear combination of $\eta_i(x_i)$:

$$(26) \quad z = \sum \alpha_i \eta_i(x_i), \quad \sum \alpha_i = 1.$$

More generally we may take any function of x :

$$(27) \quad y = g(x),$$

find out which function

$$(28) \quad g(\xi(\theta)) = \psi(\theta)$$

is of θ , invert that function

$$(29) \quad \varphi(\varphi(\theta)) = \theta$$

and take

$$(30) \quad h(x) = \varphi(g(x)).$$

With such a multitude of solutions as given by (30) and (24) -- all of which are fairly easily obtainable -- it seems natural to ask: What can be gained by a sensible choice of h ? Could the bias still present, due to our approximation, be reduced? Or could we secure a very nearly normal distribution of \hat{t} ? These questions still have to be investigated.

2. It is generally realized that maximum likelihood estimates are often hard to handle in practice. In the following is suggested the application of a class of estimates which - after a once for all tabulation of some functions - seem to be really computable and nevertheless possess a high degree of efficiency.

For simplicity we shall in this note only consider the case of n independent observations x_1, \dots, x_n and one parameter θ :

$$(1) \quad p_n \{x_1, \dots, x_n | \theta\} = \prod_{v=1}^n p_v \{x_v | \theta\}$$

We consider the computation of the mean of some tabulated function - $\varphi(x)$ - of the observation as fairly easy:

$$(2) \quad a = \frac{1}{n} \sum_{v=1}^n \varphi(x_v)$$

If φ is a sensible function and the distribution of $y = \varphi(x)$ not too bad a is fairly nearly normally distributed about the mean

$$(3) \quad a = E\{\varphi(x)\} = \int \varphi(x) p_v \{x | \theta\} dx = \lambda(\theta), \text{ say}$$

and with the variance

$$(4) \quad V\{a\} = \frac{1}{n} V\{\varphi(x)\}$$

This will often hold for medium sized and large samples.

Now denote by $k(\alpha)$ the inverse function to $\lambda(\theta)$:

$$(5) \quad \lambda(k(\alpha)) = \alpha$$

and take

$$(6) \quad t = k(a)$$

as an estimate of θ which in any particular case has particular value θ_0 .

If $V(\hat{a})$ is not too large - how large it may be depends on the function we are going to choose and has got to be checked in each particular case - we can safely apply the approximations

$$(7) \quad E\{t\} \cong k(\alpha) = \theta_0$$

and

$$(8) \quad V\{t\} \cong V\{a\} k'^2(\alpha) = \frac{1}{n} \frac{V\{\varphi(x)\}}{\lambda'^2(\theta_0)} = \frac{1}{n} \mu(\theta_0)$$

The question we are going to consider is this: How to choose $\varphi(x)$ in order to get "the best computable estimate" of θ_0 ? Under certain safeguards t may be assumed to be approximately normally distributed about θ_0 with the variance (8) and our question therefore boils down to a minimization of $V\{t\}$.

Now, let $\varphi(x)$ be the minimizing function and consider the variation of (8) if $\varphi(x)$ is replaced by some neighbour function

$$(9) \quad \varphi(x) + \epsilon \psi(x)$$

It is found that

$$(10) \quad \frac{1}{2} \delta \mu(\theta_0) = \frac{\int_{-\infty}^{\infty} \varphi(x) \psi(x) r(x|\theta_0) dx - \int_{-\infty}^{\infty} \varphi(x) r(x|\theta_0) dx \int_{-\infty}^{\infty} \psi(x) r(x|\theta_0) dx}{\left(\int_{-\infty}^{\infty} \varphi(x) \frac{\partial r(x|\theta_0)}{\partial \theta_0} dx \right)^2} - \frac{\int_{-\infty}^{\infty} \varphi^2(x) r(x|\theta_0) dx - \left(\int_{-\infty}^{\infty} \varphi(x) r(x|\theta_0) dx \right)^2}{\left(\int_{-\infty}^{\infty} \varphi(x) \frac{\partial r(x|\theta_0)}{\partial \theta_0} dx \right)^3} \int_{-\infty}^{\infty} \psi(x) \frac{\partial r(x|\theta_0)}{\partial \theta_0} dx$$

On equating to 0 and choosing for $\psi(x)$ the particular function

$$(11) \quad \psi(x) = \begin{cases} 0, & x < z \\ 1, & x > z \end{cases}$$

where z is any real number we find that $\varphi(x)$ must satisfy the equation

$$(12) \quad \int_z^{\infty} \varphi(x) r(x|\theta_0) dx - c_1 \int_z^{\infty} r(x|\theta_0) dx = c_2 \int_z^{\infty} \frac{\partial r(x|\theta_0)}{\partial \theta_0} dx$$

where c_1 and c_2 are constants:

$$(13) \quad c_1 = \int_{-\infty}^{\infty} \varphi(x) r(x|\theta_0) dx$$

$$c_2 = \frac{\int_{-\infty}^{\infty} \varphi^2(x) r(x|\theta_0) dx - \left(\int_{-\infty}^{\infty} \varphi(x) r(x|\theta_0) dx \right)^2}{\int_{-\infty}^{\infty} \varphi(x) \frac{\partial r(x|\theta_0)}{\partial \theta_0} dx}$$

when differentiating and replacing z by x we get

$$(14) \quad \varphi(x) = c_1 + c_2 \frac{\partial \log r(x|\theta_0)}{\partial \theta_0}$$

for which function actually

$$(15) E\{\varphi(x)\} = c_1 + c_2 \int_{-\infty}^{\infty} \frac{\partial \log p\{x|\theta_0\}}{\partial \theta_0} p\{x|\theta_0\} dx = c_1$$

$$(16) V\{\varphi(x)\} = c_2^2 \int_{-\infty}^{\infty} \left(\frac{\partial \log p\{x|\theta_0\}}{\partial \theta_0}\right)^2 p\{x|\theta_0\} dx = c_2^2 \int_{-\infty}^{\infty} \left(\frac{\partial p\{x|\theta_0\}}{\partial \theta_0}\right)^2 \frac{dx}{p\{x|\theta_0\}}$$

$$(17) \int_{-\infty}^{\infty} \varphi(x) \frac{\partial p\{x|\theta_0\}}{\partial \theta_0} dx = c_2 \int_{-\infty}^{\infty} \frac{\partial \log p\{x|\theta_0\}}{\partial \theta_0} \frac{\partial p\{x|\theta_0\}}{\partial \theta_0} dx$$

$$= c_2 \int_{-\infty}^{\infty} \left(\frac{\partial p\{x|\theta_0\}}{\partial \theta_0}\right)^2 \frac{dx}{p\{x|\theta_0\}}$$

and consequently

$$(18) V\{t\} = \frac{1}{n} \frac{1}{\int_{-\infty}^{\infty} \left(\frac{\partial p\{x|\theta_0\}}{\partial \theta_0}\right)^2 \frac{dx}{p\{x|\theta_0\}}}$$

which is the maximum likelihood variance.

The conclusion of this development seems rather sad. On the assumption that a "best computable estimate" of θ_0 exists we find that such an estimate depends on θ_0 itself - and in such a way that part of the argument, in the form presented, breaks down since under this condition

$$(19) \int_{-\infty}^{\infty} \varphi(x) \frac{\partial p\{x|\theta_0\}}{\partial \theta_0} dx$$

cannot be taken as $\lambda'(\theta_0)$ -- but that if it did exist, then it would be just as good as the maximum likelihood estimate itself.

In practice these difficulties can however be pretty nearly avoided: Select a "suitable" --to be specified below--sequence of point on the θ -axis:

$\theta_0, \theta_{11}, \theta_{12}, \dots$ and take them as say midpoints/^{of} to some extent overlapping intervals: $(\theta'_i, \theta'_{i+1})$ and tabulate the functions

$$(20) \varphi_i(x) = \frac{\partial \log p\{x|\theta_i\}}{\partial \theta_i}$$

Then

$$(21) a_{ii} = \frac{1}{n} \sum \varphi_i(x_v)$$

has the expectation

$$(22) \alpha_i = \int \frac{\partial \log p_i \{x | \theta_i\}}{\partial \theta_i} p_i \{x | \theta\} dx = \lambda(\theta_i, \theta), \text{ say}$$

which vanishes for $\theta = \theta_i$:

$$(23) \lambda(\theta_i, \theta_i) = 0$$

Invert the function $\lambda(\theta_i, \theta)$:

$$(24) \theta = K(\alpha_i, \theta_i)$$

and take

$$(25) t_i = K(\alpha_i, \theta_i)$$

as an estimate of θ :

$$(26) E\{t_i\} \simeq K(\alpha_i, \theta_i) = \theta$$

The variance of t_i becomes

$$(27) V\{t_i\} = \frac{V\{\alpha_i\}}{\left(\frac{\partial \lambda(\theta_i, \theta)}{\partial \theta}\right)^2} = \frac{1}{n} \mu(\theta, \theta_i)$$

where

$$(28) \mu(\theta, \theta_i) = \frac{\int_{-\infty}^{\infty} \left(\frac{\partial \log p_i \{x | \theta_i\}}{\partial \theta_i}\right)^2 p_i \{x | \theta\} dx}{\left(\int_{-\infty}^{\infty} \frac{\partial \log p_i \{x | \theta_i\}}{\partial \theta_i} \frac{\partial \log p_i \{x | \theta\}}{\partial \theta} dx\right)^2}$$

For $\theta = \theta_i$ this reduces to the maximum likelihood estimate and in practice the intervals (θ'_i, θ''_i) should be such chosen that $\mu(\theta, \theta_i)$ does not deviate too much from $\mu(\theta_i, \theta_i)$ such that the function φ_i covers an interval in which the efficiency of α_i is not smaller than wanted. For the method to be really practicable it should be possible to choose such intervals fairly large.

The method leads inevitably to more than one estimate of θ in several situations, but as each estimate gives rise to a confidence interval this should present no real difficulty in practice.

It may be mentioned that the method has actually been used in a case which otherwise would seem rather difficult to handle (Hald, Jersild & Rasch, "On the

Determination of the Phagocytic Power of Lencocytes." Acta path. et microbiol.
20, 1942.)

One point more may be worth mentioning: The characteristic function of
is readily available. When von Neuman's machine starts working, it may therefore
be possible to solve the sampling problems even for small samples. Which of course
does not imply that such estimates are "best" in any sense for small