Printed at the Mathematical Centre, 49, 2e Boerhaavestraat 49, Amsterdam.

.

The Mathematical Centre, founded the 11-th of February 1946, is a nonprofit institution aiming at the promotion of pure mathematics and its applications. It is sponsored by the Netherlands Government through the Netherlands Organization for the Advancement of Pure Research (Z.W.O.), by the Municipality of Amsterdam, by the University of Amsterdam, by the Free University at Amsterdam, and by industries.

## MATHEMATICAL CENTRE TRACTS 26

# SELECTED **STATISTICAL PAPERS 1**

**EUROPEAN MEETING 1968** 

MATHEMATISCH CENTRUM AMSTERDAM 1968

AMS(MOS) subject classification scheme (1970): 60J10, 60K05, 62E10, 62F10, 62C05, 62G10, 62J05.

· · · · ·

.

## PREFACE

.

This volume of the series "Mathematical Centre Tracts" is published on the occasion of the European Meeting 1968 on Statistics, Econometrics and Management Science in Amsterdam. With permission of the Organizing Committee of this Meeting, the Statistical Department of the Mathematical Centre has invited some authors of papers on Statistics and Probability Theory to publish their work in the form of this Tract. This first volume contains eight papers. The authors come from seven countries and their subjects vary from renewal processes to slippage tests. The papers appear in an almost random order, determined mainly by the date of arrival of the manuscript. It is hoped that a second volume of the same kind will be published shortly after the Meeting.

## Contents

M. Rao and H. Wedel, Poisson processes as renewal processes	
invariant under translations	3
A.M. Kagan, Some analytical aspects of estimation theory	7
H.K. Ury, The behavior of some tests for ordered alternatives under interior slippage	15
D. Bierlein, The development of the concept of statistical decision theory	27
R.L. Brown and J. Durbin, Methods of investigating whether a regression relationship is constant over time	37
W.G. Cochran, Some effects of errors of measurement on multiple regression	47
J.F.C. Kingman, Some recent developments in the theory of Markov chains	71
B. Eichhorn, On sequential search	81

## POISSON PROCESSES AS RENEWAL PROCESSES INVARIANT UNDER TRANSLATIONS

## by Murali Rao and Hans Wedel (Sweden) University of Göteborg

INTRODUCTION: Thedéen has proved [3] that a renewal process whose renewals are "stationary under translation" is necessarily Poisson. In search for some sort of generalization of this interesting result we arrived at a very simple proof which we present.

Let  $\{X_n : n = \pm 1, \pm 2, \ldots\}$  be a sequence of random variables such that a.s.

$$x_{-2} < x_{-1} < 0 < x_1 < x_2 \dots$$

Put  $Y_0 = X_{-1}$ ,  $Y_1 = X_1$ ,  $Y_n = X_n - X_{n-1}$  for  $n \neq 0$ , 1. Assume that:

- i)  $\{(Y_0, Y_1), Y_n \ n \neq 0, 1\}$  is a set of independent random variables.
- ii)  $\{Y_n : n \neq 0, 1\}$  are independent, identically distributed positive
- random variables with  $P[Y_n \le y] = F(y)$ , F(0) = 0 and  $E[Y_n] = \frac{1}{m} < \infty$ . iii) E[N(I)] = m|I| where |I| denotes the Lebesgue measure of I,

and N(I) = number of  $X_n \in I$ . Then iii) is equivalent to

$$P(Y_i > u) = \int_u^{\infty} (1 - F(t)) dt \text{ for } i = 0, 1; \text{ see } [1, p. 354].$$

Let  $\{\xi_n: n = \pm 1, \pm 2, \ldots\}$  be a sequence of random variables which is independent of the sequence  $\{X_n\}$ . We shall assume that for all n,m n  $\neq$  m  $(\xi_n, \xi_m)$  have the same joint distribution G and that the support group of G, i.e. the group generated by the support of G, has an element of the form (0,d) with d > 0; if  $\xi_n$  and  $\xi_m$  are independent and have a nondegenerate distribution then this is certainly true.

Put 
$$Z_n = X_n + \xi_n$$
,  $n = \pm 1, \pm 2, \ldots$  and  $N(I) =$  number of  $Z_n \in I$ .

Theorem 1.

Let  $X_n$ ,  $\xi_n$ ,  $Z_n$  be as above. If  $E[\tilde{N}(I) \ \tilde{N}(J)] = E[N(I) \ N(J)]$  for all I, J then  $\{X_n\}$  is Poisson, i.e.  $F(y) = 1 - e^{-my}$ .

Proof. Put 
$$\phi(\mathbf{I}, \mathbf{J}) = \mathbf{E}[\mathbf{N}(\mathbf{I}) \ \mathbf{N}(\mathbf{J})] - \mathbf{E}[\mathbf{N}(\mathbf{I} \cap \mathbf{J})] = \sum_{n \neq m} \mathbf{P}[\mathbf{X}_n \in \mathbf{I}, \ \mathbf{X}_m \in \mathbf{J}].$$
  
Using independence of  $\{\xi_n\}$  and  $\{\mathbf{X}_n\}$  we get  
 $\mathbf{E}[\tilde{\mathbf{N}}(\mathbf{I}) \ \tilde{\mathbf{N}}(\mathbf{J})] - \mathbf{E}[\tilde{\mathbf{N}}(\mathbf{I} \cap \mathbf{J})] = \iint_{\mathbf{V}} \phi(\mathbf{I} - \mathbf{u}, \ \mathbf{J} - \mathbf{v}) d\mathbf{G}(\mathbf{u}, \mathbf{v}).$   
The condition  $\mathbf{E}[\mathbf{N}(\mathbf{I})] = m |\mathbf{I}|$  implies  $\mathbf{E}[\tilde{\mathbf{N}}(\mathbf{I})] = m |\mathbf{I}|$ . Thus  $\mathbf{E}[\tilde{\mathbf{N}}(\mathbf{I}) \ \tilde{\mathbf{N}}(\mathbf{J})] - \mathbf{E}[\tilde{\mathbf{N}}(\mathbf{I} \cap \mathbf{J})] = \mathbf{E}[\mathbf{N}(\mathbf{I} \cap \mathbf{J})] = \mathbf{E}[\mathbf{N}(\mathbf{I} \cap \mathbf{J})] = \mathbf{E}[\mathbf{N}(\mathbf{I} \cap \mathbf{J})] : \phi = \phi \star \mathbf{G}.$ 

A simple consequence of the renewal theorem is that for any finite intervals I, J, E[N(I + h) N(J + k)] is a bounded function of (h,k). The Choquet-Deny theorem [2, p. 152] applies and we deduce that every point of support of G is a period for  $\Phi$ . The set of periods for  $\Phi$  is a group and this group contains the element (0,d) and hence (0,kd) where k is any positive integer (indeed any integer). Thus for all I, J and all positive integers  $k, \Phi(I,J) = \Phi(I, J + kd)$ . Take I = (0,x], J = (0,x] with x < kd. Then  $I \cap (I + kd) = \emptyset$ .

Also 
$$\Phi(\mathbf{I}, \mathbf{I} + \mathbf{kd}) = \sum_{n \neq m} P[X_n \in \mathbf{I}, X_m \in \mathbf{I} + \mathbf{kd}] = \sum_{m,n \geq 1} P[X_n \in \mathbf{I}, X_m \in \mathbf{I} + \mathbf{kd}] =$$

$$= \sum_{n=1}^{\infty} \sum_{m>n} P[X_n \in I, X_m \in I + kd] = \sum_{n=1}^{\infty} \int_0^x H(I + kd - u) d(F_0 * F^{(n-1)*})(u) =$$
$$= m \int_0^x H(I + kd - u) du \text{ where } H(x) = \sum_{n=1}^{\infty} F^{k*}(x) \text{ and iii} \text{ implies}$$

= m 
$$\int_0^x H(I + kd - u) du$$
 where  $H(x) = \sum_{k=1}^{\infty} F^{k^*}(x)$  and iii) implies

$$mx = \sum_{k=0}^{\infty} F_0 \star F^{k \star}(x) \cdot x > 0$$

Similar calculations give  $\Phi(I,I) = 2m \int_0^x H(x - u) du = 2m \int_0^x H(u) du$ .

Thus 
$$2 \int_0^x H(u) du = \int_0^x H(I + kd - u) du = \int_0^x [H(x + kd - u) - H(kd - u)] du = \int_0^x [H(kd + u) - H(kd - u)] du.$$

This equality for all  $x \le kd$  implies 2H(u) = H(kd + u) - H(kd - u);  $u \le kd$ .

R&W 2

It is possible to show that the only solution of this functional equation is  $H(x) = \lambda x$ . However we take a short cut.

Suppose  $d_0$  is any positive number such that  $F(d_0) > 0$  and  $F(d_0^{-}) = 0$ . Then  $F^{n^{**}}$  has an atom at  $nd_0$  and thus H has a mass at every positive integral multiple of  $d_0$ , but H(u) = 0 for  $u < d_0$ . Choose k so that kd >  $d_0$ . For  $u < d_0$ , H(u) = 0 and the functional equation for H shows that H(kd - u) = H(kd + u);  $u < d_0$ . This is absurd since every interval of length larger than  $d_0$ , has a multiple of  $d_0$  and H has a mass at such a point. Thus F certainly cannot be arithmetic. As  $k \to \infty$  the renewal theorem shows that 2H(u) = 2um. This is equivalent to F being exponential. Q.E.D.

## Theorem 2.

If the support group of G is dense in the plane, condition iii) can be removed and a sufficient condition for the preceding theorem is

$$\mathbf{E}\left[\mathbf{N}(\mathbf{I}) \ \mathbf{N}(\mathbf{J}) - \mathbf{N}(\mathbf{I} \cap \mathbf{J})\right] = \mathbf{E}\left[\mathbf{\tilde{N}}(\mathbf{I}) \ \mathbf{\tilde{N}}(\mathbf{J}) - \mathbf{\tilde{N}}(\mathbf{I} \cap \mathbf{J})\right].$$

<u>Proof</u>. With the same notation as above the Choquet-Deny theorem shows that  $\phi(\mathbf{I}, \mathbf{J}) = \lambda |\mathbf{I}| |\mathbf{J}|$ , thus

$$\lambda hx = E[N[-h,0) N(0,x]] = \sum_{n \neq m} P[X_n \in [-h,0), X_m \in (0,x]] =$$

$$= \sum_{\substack{n < 1 \\ m \ge 1}} \int_0^x \int_0^h P[X_n \in [-h,0), X_m \in (0,x]] | y_0 = u, y_1 = v] dK(u,v) =$$

$$= \int_0^x \int_0^h \sum_{\substack{n < -1 \\ m \ge 1}} P[X_n \in [-h,0) | y_1 = u] \sum_{\substack{m \ge 1 \\ m \ge 1}} P[X_m \in (0,x]] | y_1 = v] dK(u,v)$$

$$= \int_0^h \int_0^x U(x - v) U(h - u) dK(u,v)$$

where U(I) =  $\sum_{h=0}^{\infty} F^{n^*}(I)$  and K is the joint distribution of  $y_0$  and  $y_1$ .

R&W 3

Since the left side is a product measure this implies that dK is a product measure:

Say  $dK(u) \cdot dK(v) = dK(u,v)$ .

Further 
$$\lambda \mathbf{x} = \int_0^{\mathbf{x}} U(\mathbf{x} - \mathbf{u}) \, dK(\mathbf{v}) \cdot \frac{1}{h} \int_0^h U(\mathbf{h} - \mathbf{v}) \, dK(\mathbf{v}).$$

Put x = 1, we see  $\frac{1}{h} \int_0^h U(h - v) dK(v) = \mu$  where  $\mu$  is a constant. This implies  $(K \neq U)(x) = \mu x$ , i.e.  $E[N(I)] = \mu |I|$ . Put K  $\neq$  U = V and I = J = (0,x] then  $\Phi(I,I) = \lambda x^2$ ,  $\Phi(I,I) =$ 

$$= \sum_{n \neq m} P[x_n \in I, x_m \in I] = 2V * H = \mu 2 \int_0^x H(t)dt, \text{ thus we have } H(x) = \frac{\lambda}{\mu}x.$$

<u>Acknowledgement</u>: It was H. Bergström who pointed out that our conditions are sufficient. We are grateful to him and P. Jagers for valuable discussions.

## References:

- 1. W. FELLER, An Introduction to Probability Theory and its Applications, Vol. II, Wiley, New York 1966.
- 2. P. MEYER, Probability and Potentials, Blaisdell, 1966.
- 3. T. THEDÉEN, On stochastic stationary of renewal processes, Arkiv för matematik, Band 7, Häfte 3, 1967, pp. 249-263.

R&W 4

## SOME ANALYTICAL ASPECTS OF ESTIMATION THEORY

by A.M. Kagan (USSR) Mathematical Institute of the Academy of Sciences, Leningrad

In the paper the close connection between certain problems of the statistical theory of estimation and analytical problems of the characterization of distributions is demonstrated.

More precisely let  $(x_1, \ldots, x_n)$  be a repeated sample from the population with distribution function (d.f.)  $F(x;\theta)$ , where  $\theta$  is a parameter. Suppose that  $g(x_1, \ldots, x_n)$  is an admissible or optimal (in a certain sense) estimator of the parametric function  $\gamma(\theta)$ . What are the conditions imposed on  $F(x;\theta)$  by admissibility or optimality of more or less simple estimators  $g(x_1, \ldots, x_n)$  - that is the question we discuss in the paper. We mention also a number of related results concerning sufficiency of statistics and Fisher's information.

Recently certain results have been obtained for the exponential families and for the families with group parameters - location and scale. We shall restrict ourselves with the families depending on the scale parameter because the principal results for the exponential families and for the families with the location parameter were reported in Linnik's paper [1] at the previous European Meeting of Statisticians (London, 1966).

Everywhere the quadratic loss function will be used; it means that the quality of an estimator  $g(x_1, \ldots, x_n)$  of  $\gamma(\theta)$  is measured by  $E_{\theta}(g - \gamma(\theta))^2$ . The agreement automatically defines the conceptions of admissibility and optimality.

## 1. ESTIMATION OF POLYNOMIALS OF SCALE PARAMETER.

Let  $(x_1, \ldots, x_n)$  be a repeated sample from the population with d.f.  $F(\frac{x}{\sigma})$  depending on the scale parameter  $\sigma \in (0, \infty)$ . Everywhere in the paper F(x) is supposed to be concentrated on  $(0, \infty)$ . Assume that

AMK 1

 $\int_0^\infty x \, \mathrm{d}F(x) = \alpha_1^{<\infty}, \quad \int_0^\infty x^2 \, \mathrm{d}F(x) = \alpha_2^{<\infty};$  $\alpha_1^{-1}\overline{x} = \frac{x_1^{+} \cdots + x_n}{\alpha_1^{n}}$ 

then

will be an unbiased estimator of  $\sigma$  with finite variance. It is easily to see that apart from the trivial case of degenerate F(x) the best estimator of  $\sigma$  of the form  $c\overline{x}$  (we shall denote it  $c_n^{0}\overline{x}$ ) which has a bias, is better than  $\alpha_1^{-1}\overline{x}$ , i.e. it satisfies the condition

$$\mathbb{E}_{\sigma}(c_{n}^{0}\overline{x} - \sigma)^{2} \leq \mathbb{E}_{\sigma}(\alpha_{1}^{-1}\overline{x} - \sigma)^{2} \text{ for all } \sigma \in (0, +\infty).$$

That is why it is natural to clear up the conditions of admissibility of  $\alpha_1^{-1}\overline{x}$  among all unbiased estimators of  $\sigma$  and the conditions of admissibility of  $c_n^{0-x}$  among all estimators of  $\sigma$ .

The next theorems were proved in [2]; there F(x) was a priori supposed to satisfy the condition

$$\int_{0}^{\infty} x^{k} dF(x) < \infty, \ k = 1, \ 2, \ \dots \ .$$
 (1)

Theorem 1.1. Let F(x) satisfy the condition (1). Then the necessary and sufficient condition for  $\alpha_1^{-1-x}$  to be admissible among unbiased estimators of  $\sigma \in (0,\infty)$  for two sample sizes  $n = n_1$ ,  $n = n_2$ ,  $n_2 > n_1 \ge 3$ , is that F(x) is either a degenerate d.f. or a d.f. of the gamma-distribution.

Theorem 1.2. Let F(x) satisfy the condition (1). The necessary and sufficient condition for  $c_n^{0-x}$  to be admissible among all estimators of  $\sigma \in (0,\infty)$ , for two sample sizes  $n = n_1$ ,  $n = n_2$ ,  $n_2 > n_1 \ge 3$ , is that F(x) is either a degenerate d.f. or a d.f. of the gamma-distribution. It should be noticed that using analytical results obtained recently by C.G. Khatri and C.R. Rao [3] one can avoid the condition (1). We shall now outline briefly the scheme of the proof of Theorems 1.1 and 1.2. Sufficiency is proved in the following manner. The case of AMK 2

degenerate F(x) is trivial. If F(x) is a function of the gamma-distribution then  $\overline{x}$  will be a complete sufficient statistic for the family  $\frac{x}{r(\frac{1}{\sigma})} \dots F(\frac{n}{\sigma})$ . Hence according to the Rao-Blackwell-Kolmogorov theorem it follows that in this case  $\alpha_1^{-1}\overline{x}$  is for all n not only admissible but the best unbiased estimator of  $\sigma \in (0,\infty)$ . The proof of the admissibility of  $c_n^{0}\overline{x}$  in this case is also based on sufficiency of  $\overline{x}$  and on the Cramér-Rao inequality (cf. [4]). Necessity of the conditions of Theorems 1.1 and 1.2 is proved almost in the same manner. We shall restrict ourselves to Theorem 1.2.

Let us consider the estimator  $S_n(x_1, \ldots, x_n) = S_n$ ,

$$S_{n} = c_{n}^{0} \frac{E_{1}(c_{n}^{0} | y)}{E_{1}((c_{n}^{0} x)^{2} | y)}, \qquad (2)$$

where  $y = (\frac{x_2}{x_1}, \dots, \frac{x_n}{x_1})$ . It can be proved that  $E_{\sigma}(S_n - \sigma)^2 \leq E_{\sigma}(c_n^{0-1} - \sigma)^2$ 

and the equality sign in (3) holds - simultaneously for all  $\sigma \in (0,\infty)$  - if and only if

$$E_{1}(c_{n}^{0-}|y) = E_{1}((c_{n}^{0-})^{2}|y)$$
(4)

(3)

with probability 1.

Analytically the condition (4) is convenient enough. Denoting

$$x_{i} = e^{\xi_{i}},$$
  
G(u) = P{ $\xi_{i} < u; \sigma = 1$ }

we can rewrite the condition (4) as

$$E\{\frac{1}{n}\sum_{1}^{n} e^{\xi_{1}} | \xi_{2} - \xi_{1}, \dots, \xi_{n} - \xi_{1}\} =$$

$$= c_{n}^{0} E\{(\frac{1}{n}\sum_{1}^{n} e^{\xi_{1}})^{2} | \xi_{2} - \xi_{1}, \dots, \xi_{n} - \xi_{1}\}.$$
(5)

Multiplying both sides of (5) by  $\exp(t_2(\xi_2 - \xi_1) + \ldots + t_n(\xi_n - \xi_1))$ and taking the expectations of both sides we obtain the following functional equation for the Laplace transform

$$P(z) = \int_{-\infty}^{+\infty} e^{zu} dG(u):$$

$$B_{n} \{ P(1 - \sum_{i=2}^{n} t_{i}) \prod_{i=2}^{n} P(t_{i}) + P(-\sum_{i=2}^{n} t_{i}) \sum_{k=2}^{n} [P(1+t_{k}) \prod_{\substack{i=2 \ i \neq k}}^{n} P(t_{i})] \} =$$

$$= P(2 - \sum_{i=2}^{n} t_{i}) \prod_{i=2}^{n} P(t_{i}) + P(-\sum_{i=2}^{n} t_{i}) \sum_{k=2}^{n} [P(2+t_{k}) \prod_{\substack{i=2 \ i \neq k}}^{n} P(t_{i})] +$$

$$+ 2\{ P(1 - \sum_{i=2}^{n} t_{i}) \sum_{k=2}^{n} [P(1+t_{k}) \prod_{\substack{i=2 \ i \neq k}}^{n} P(t_{i}) + (6)$$

$$+ P(-\sum_{i=2}^{n} t_{i}) \sum_{j>k\geq 2}^{n} [P(1+t_{i}) P(1+t_{k}) \prod_{\substack{i=2 \ i \neq k}}^{n} P(t_{i})] \},$$

$$= P(2 - \sum_{i=2}^{n} t_{i}) \sum_{j>k\geq 2}^{n} [P(1+t_{i}) P(1+t_{k}) \prod_{\substack{i=2 \ i \neq j \ i \neq j}}^{n} P(t_{i})] \},$$

where  $B_n = constant$  and  $Re t_i = 0$ .

.

From (6) the desired result follows after certain analytic transformations.

Note that if F'(x) = f(x) exists then the estimator (2) takes the form obtained originally by E. Pitman [5]:

$$S_{n} = \frac{\int_{0}^{\infty} u^{n} \prod_{i=1}^{n} f(ux_{i}) du}{\int_{0}^{\infty} u^{n+1} \prod_{i=1}^{n} f(ux_{i}) du}$$

A development of the method used in [2] allowed to obtain the following theorem (see [6]):

AMK 4

Theorem 1.3. Suppose that F(x) satisfies the condition (1). The necessary and sufficient condition for a polynomial  $g(\overline{x}) = a_0 \overline{x}^k + \ldots$  $\ldots + a_k, a_0 \neq 0$ , of degree  $k \geq 1$  to be optimal for all  $\sigma \in (0,\infty)$  among unbiased estimators of  $\pi_k(\sigma) = E_{\sigma}q(\overline{x})$  for k sample sizes  $n = m, m+1, \ldots$  $\ldots, m+k-1, m \geq 3$ , is that F(x) is either a degenerate d.f. or a d.f. of the gamma-distribution.

The proof of Theorem 1.3 appears to be equivalent to solving the following equation for P(z):

$$\sum_{j=0}^{k} \left[ P(j - \sum_{2}^{n} t_{i}) \sum_{\substack{m_{2} + \dots + m_{n} = k - j \\ m_{i} \ge 0}} \prod_{2}^{n} P(m_{k} + t_{k}) \right] =$$

$$= b_{n} P(-\sum_{2}^{n} t_{i}) \prod_{2}^{n} P(t_{i}), n = m, m+1, \dots, m+k-1,$$
(7)

 $b_n = const., Re t_i = 0.$ 

## 2. SUFFICIENCY AND PARTIAL SUFFICIENCY

Under different conditions on F(x) there was proved in the papers [7,8,9] that sufficiency of the statistic  $\overline{x}$  for the family

$$F(\frac{x_1}{\sigma}) \ldots F(\frac{x_n}{\sigma}), \sigma \in (0,\infty)$$
 (8)

is equivalent to the fact that F(x) is a d.f. of the gamma-distribution. It appears [9] that the independence of  $\sigma$  of the conditional expectation  $E_{\sigma}(Q|\overline{x})$  for a separate polynomial Q of general form imposes strong conditions upon F(x). In particular the following theorem holds.

Theorem 2.1. If  $1^{\circ}$ .  $\int_{0}^{\infty} x^{2} dF(x) < \infty$ ,  $2^{\circ}$ . the n-th convolution  $F^{\mathbf{x}n}(x)$  is absolutely continuous,  $3^{\circ}$ .  $\mathbf{E}_{\sigma}(x_{1}^{2}|\overline{x})$  is independent of  $\sigma$ , then F(x) is a d.f. of the gamma-distribution.

Now we are going to generalize the conception of sufficiency (cf. [10]). Assume that for some integer  $k \ge 1$ 

AMK 5

$$\int_{0}^{\infty} x^{2k} dF(x) < \infty.$$
 (9)

Under this condition the set of all polynomials  $Q(x_1, \ldots, x_n)$  of degree  $\leq k$  forms a Hilbert space if one defines the scalar product of elements  $Q_1$  and  $Q_2$  as

$$(Q_1, Q_2)_{\sigma} = E_{\sigma}(Q_1 Q_2).$$
 (10)

We shall denote this space by  $L_k^{(2)}$  and its subspace formed by all polynomials  $q(\overline{x}) = a_0 \overline{x}^k + \ldots + a_k$  of the sample mean  $\overline{x}$  will be denoted by  $\mathcal{T}_k$ .

The subspace  $\mathcal{T}_k$  is said to be  $L_k^{(2)}$ -sufficient for the family (8) if for any  $Q \in L_k^{(2)}$  there exists an element  $q \in \mathcal{T}_k$  independent of  $\sigma \in (0, \infty)$  such that

$$\hat{\mathbf{E}}_{\sigma}(\mathbf{Q}|\mathcal{T}_{\mathbf{k}}) = q$$

where  $\hat{E}_{\sigma}(\cdot | \mathcal{T}_{k})$  denotes the projection into  $\mathcal{T}_{k}$  when the scalar product in  $L_{k}^{(2)}$  is defined by the formula (10) with given  $\sigma$ .

Theorem 2.2 (cf. [9]). The necessary and sufficient condition for  $\mathcal{T}_k$  to be  $L_k^{(2)}$ -sufficient for the family (8) is that either F(x) is a degenerate d.f. or the first 2k moments of F(x) coincide with the corresponding moments of the gamma-distribution.

From Theorem 2.2 it follows that if the first 2k moments of F(x) coincide with the corresponding moments of the gamma-distribution then any polynomial  $Q \in L_k^{(2)} \setminus \mathcal{T}_k$  will be inadmissible among unbiased estimators of  $E_{\sigma}Q$ .

3. EXTREMAL ROLE OF THE GAMMA-DISTRIBUTION IN INFORMATIONAL SENSE Suppose that F(x) has the density f(x). The integral

$$\mathcal{I}_{f}(\sigma) = \int_{f(\frac{x}{\sigma}) > 0} \left[ \frac{\partial \log \frac{1}{\sigma} f(\frac{x}{\sigma})}{\partial \sigma} \right]^{2} \frac{1}{\sigma} f(\frac{x}{\sigma}) dx \qquad (11)$$

is well known as Fisher information for the family of densities  $\frac{1}{\sigma} f(\frac{x}{\sigma})$ .

Suppose that the following conditions are satsfied.

- 1. f(x) is continuously differentiable,
- $2. \int_0^\infty x^2 f(x) dx < \infty,$
- 3.  $\lim_{x \to 0} xf(x) = 0$ ,  $\lim_{x \to \infty} x^2 f(x) = 0$ .

Theorem 3.1 (see [11]). Within the class of densities with given moments  $\alpha_1$ ,  $\alpha_2$  satisfying the conditions 1-3, min  $\mathcal{I}_f(\sigma)$  is attained simultaneously for all  $\sigma \in (0, \infty)$  - by the gamma-distribution. The comparison of the results of this paper with the results of [10, 12,13] shows that for problems concerning the scale parameter the gamma-distribution plays the same role as the normal law does for problems concerning the location parameter.

REFERENCES

- [1] Yu.V. Linnik, Some new investigations on parametric statistics, Communication at the European Meeting of Statisticians, London (1966).
- [2] A.M. Kagan, A.L. Ruhin, On the estimation theory of scale parameter (in Russian), Teoriya veroyatn. i primeneniya, XII, 4 (1967).
- [3] C.G. Khatri and C.R. Rao, On some characterizations of gamma-distribution, Sankhya (1968).
- [4] J. Hodges, E. Lehmann, Some applications of Cramér-Rao inequality, Proc. II Berkeley Symposium on Prob. Statist., Vol. II (1952).
- [5] E. Pitman, The estimation of location and scale parameters of a continuous population of any given form, Biometrika, 30, III-IV (1938).

13

AMK 7

	14
[6] F.M. Kagan,	Optimality condition of certain estimators for families with scale parameter (in Russian), Doklady Akademii Nauk Uzbek. SSR, 6 (1968).
[7] E.B. Dynkin,	Necessary and sufficient statistics for a family of probability distributions (in Russian), Uspehi Matem. Nauk, VI, I (1951).
[8] T. Ferguson,	Location and scale parameters in exponential families of distributions, Ann. Math. Stat., 33, 3 (1962).
[9] F.M. Kagan,	Characterization of gamma-distribution by statis- tical properties of sample mean (in Russian), Izvestiya Akademii Nauk Uzbek. SSR, 4 (1967).
[10] A.M. Kagan,	Partial sufficiency and unbiased estimation of polynomials of location parameter (in Russian), Doklady Akademii Nauk SSSR, 174, 6 (1968).
[11] F.M. Kagan,	An informational property of gamma-distribution (in Russian), Izvestiya Akademii Nauk Uzbek. SSR, 5 (1967).
[12] A.M. Kagan, O.V.	Shalaevsky, Characterization of the normal law by the property of partial sufficiency (in Russian), Teoriya veroyatn.i primeneniya, XII, 3 (1967).
[13] A.M. Kagan,	On the estimation theory of location parameter, Sankhya, A 28, 3-4 (1966).

AMK 8

# THE BEHAVIOR OF SOME TESTS FOR ORDERED ALTERNATIVES UNDER INTERIOR SLIPPAGE

by Hans K. Ury (U.S.A) San Francisco Medical Center, University of California and California State Department of Public Health

#### 0. INTRODUCTION AND SUMMARY

For testing the standard null hypothesis against ordered alternatives in a one-way analysis of variance with k samples, Bartholomew [1], [3] proposed some test statistics for the case of underlying normal distributions, and Chacko gave the corresponding nonparametric test in [8]. In this note it is shown that under "interior slippage" (i.e., when one population other than the first or  $k^{th}$  is larger or smaller than the others), the probability that these tests will reject the null hypothesis simultaneously in favor of *both* the alternatives of upward and downward ordering goes to 1 in the limit, as the sample sizes grow sufficiently large, regardless of the significance level.

Since a nonparametric test against trend by Terpstra [13] and Jonckheere [10] is shown to behave somewhat more "normally" under interior slippage, it is suggested that this test might well be preferable at least for small k, particularly for k = 3.

## 1. NOTATION AND THE TESTS

Since the underlying models for the tests of (a) Bartholomew, (b) Chacko, and (c) Terpstra are different, they will be given separately.

(a) Here we have k independent normal random variables,  $X_1, \ldots, X_k$  with unknown means  $\mu_1, \ldots, \mu_k$  and a common but unknown variance  $\sigma^2$ . Let  $x_{ij}$  (i = 1, ..., k; j = 1, ..., n<sub>i</sub>) be independent observations on the k variables, with  $x_{ij}$  the j<sup>th</sup> observation from the i<sup>th</sup> variable. Let  $\overline{x}_i = \sum_{j=1}^{n_i} x_{ij}/n_i$  and  $s_i^2 = \sum_{j=1}^{n_i} (x_{ij} - \overline{x}_i)^2/n_i$  denote the sample mean and HKU 1

variance for the i<sup>th</sup> variable. The null hypothesis  $H_0: \mu_1 = \ldots = \mu_k$  is tested against either the alternative of upward ordering,  $H_1: \mu_1 \leq \ldots \leq \mu_k$  or against that of downward ordering,  $H_2: \mu_1 \geq \ldots \geq \mu_k$ , with at least one inequality strong in each case and with  $\sigma^2$  unspecified for all three hypotheses. Denote by

(1) 
$$\overline{x}[t,s] = (n_t \overline{x}_t + n_{t+1} \overline{x}_{t+1} + \dots + n_s \overline{x}_s)/(n_t + n_{t+1} + \dots + n_s)$$

the pooled sample mean of  $\overline{x}_t$ ,  $\overline{x}_{t+1}$ , ...,  $\overline{x}_s$ , where s and t are positive integers with  $1 \leq t \leq s \leq k$ .

The MLE's (maximum likelihood estimates) of the µ's under  $H_0$  are  $\hat{\mu}_1 = \ldots = \hat{\mu}_k = \bar{x}_{[1,k]}$ . Under  $H_1$  or  $H_2$  the MLE's are obtained by pooling successive sample means which violate the restriction specified by the alternative, continuing this procedure until no violations exist among the remaining pooled or unpooled means. If there are m distinct estimates obtained by pooling, respectively, the first  $t_1$  means, the next  $t_2$  means, ..., and the last  $t_m$  means,  $t_j > 0$ ,  $\sum_{j=1}^m t_j = k$ , and if we set  $\tau_0 = 0$ , (2)  $\tau_i = t_1 + t_2 + \ldots + t_i$  (i = 1, 2, ..., m)

 $\tau_{m} = k$ ,

then

(3) 
$$\hat{\mu}_{\tau_{i}+1} = \hat{\mu}_{\tau_{i}+2} = \dots = \hat{\mu}_{\tau_{i+1}} = \overline{x}[\tau_{i}+1,\tau_{i+1}]$$
 (i = 0, 1, ..., m-1).

Denote the m distinct estimates by  $\overline{x}_{t_j}$  (j = 1, ..., m), where  $\overline{x}_{t_j}$  = =  $\overline{x}_{[\tau_{j-1}+1,\tau_j]}$ . Let the sum of the sample sizes pooled into  $\overline{x}_{t_j}$  be denoted by N<sub>t\_j</sub>. [It follows from Brunk [5], [6] and van Eeden [9] that these MLE's are unique and can be formally represented as  $\hat{\mu}_i = \max_{\substack{1 \le r \le s \\ 1 \le r \le i}} \min_{\substack{i \le s \le k \\ i \le s \le k}} \overline{x}_{[r,s]}$  for the case of H<sub>1</sub>, and as  $\hat{\mu}_i = \min_{\substack{1 \le r \le i \\ 1 \le r \le i \le s \le k}} \overline{x}_{[r,s]}$ 

Finally, let  $p_{m,k}$  stand for the probability of obtaining exactly m distinct estimates out of a possible k. Under  $H_0$ , and for equal sample sizes, Chacko [7] and Miles [11] showed that

HKU 2

$$p_{m,k} = |S_k^m| / k!$$

where  $|S_k^m|$  is the coefficient of  $z^m$  in  $z(z+1) \dots (z+k-1)$  (i.e., is the modulus of a Stirling's number of the First Kind). They also showed that (4) holds not only for  $X_1, \dots, X_k$  normally distributed, but whenever their joint distribution is a symmetric function of them. Bartholomew's likelihood ratio test at significance level  $\alpha$  [3] calls for rejection of  $H_0$  when  $T_k^1 \ge C_1$ , where the test statistic

(5) 
$$T_{k}^{1} = \sum_{i=1}^{k} n_{i} [\hat{\mu}_{i} - \overline{x}_{[1,k]}]^{2} / s_{0}^{2} = \sum_{j=1}^{m} N_{t_{j}} [\overline{x}_{t_{j}} - \overline{x}_{[1,k]}]^{2} / s_{0}^{2},$$

with  $s_0^2 = \sum_{i=1}^k n_i [\overline{x}_i - \overline{x}_{[1,k]}]^2 + \sum_{i=1}^k n_i s_i^2$ , and where  $C_1$  is determined by

(6) 
$$\alpha = \sum_{m=2}^{k} p_{m,k} P[\beta_{(m-1)/2, (N-m)/2}] \ge C_1].$$

Here  $\beta [(m-1)/2, (N-m)/2]$  is a random variable having the Beta distribution with parameters (m-1)/2 and (N-m)/2, and  $N = \sum_{i=1}^{k} n_i = \sum_{j=1}^{m} N_t$ .

(a') When  $\sigma^2$  is known, the test statistic is

(7) 
$$T_{k}^{2} = \sum_{i=1}^{k} n_{i} \left[ \hat{\mu}_{i} - \overline{x}_{[1,k]} \right]^{2} / \sigma^{2} = \sum_{j=1}^{m} N_{t_{j}} \left[ \overline{x}_{t_{j}} - \overline{x}_{[1,k]} \right]^{2} / \sigma^{2},$$

and  ${\rm H}_0$  is rejected at level  ${\alpha}$  if  ${\rm T}_k^2 \geq {\rm C}_2,$  with  ${\rm C}_2$  determined by

(8) 
$$\alpha = \sum_{m=2}^{k} p_{m,k} P[\chi_{m-1}^{2} \ge C_{2}],$$

where  $\chi^2_{m-1}$  is a random variable having the Chi square distribution with m-1 degrees of freedom.

(a") The computation of the  $p_{m,k}$  imposes limitations on the use of  $T_k^1$  and  $T_k^2$ . For unequal sample sizes, these probabilities have been determined only for k = 3, 4, and 5 [1]. For equal  $n_i$ , they were tabulated HKU 3

for k  $\leq$  12 in [11], and tables of the Stirling Numbers of the First Kind can of course be used for k  $\leq$  50. Barton and Mallows [4] give some approximations which should prove useful for k quite large. But in the general case of unequal samples (or for moderately large k) one could use a conditional test, suggested by Tukey and mentioned in [2], which does not require a knowledge of the  $p_{m,k}$ . To apply this test, one simply computes the test statistic as in (5) or (7) and determines significance by referring to the percentage points of  ${}^{\beta}[(m_{0}-1)/2,(N-m_{0})/2]$  or  $\chi^{2}_{m_{0}}-1$ , respectively, where  $m_{0}$  is the observed value of m. The corresponding test statistics will be denoted  $T_{k}^{1c}$  and  $T_{k}^{2c}$ .

(b) Since (4) holds only for symmetrically dependent random variables, Chacko's rank test [8] calls for equal sample sizes. Thus, let k independent random samples be drawn from populations with unknown continuous (to avoid ties) cumulative distributions  $F_i$  (i = 1, ..., k) respectively. We now have  $H_0$ :  $F_1 = \ldots = F_k$ ;  $H_1$ :  $F_1 \geq \ldots \geq F_k$ ;  $H_2$ :  $F_1 \leq \ldots \leq F_k$  with at least one inequality strong in each case. Chacko's test procedure consists of replacing each  $x_{ij}$  by its rank  $R_{ij}$  in the overall sample and, letting  $\overline{R_i} = (1/n) \sum_{j=1}^{n} R_{ij}$ , formally operating on the  $\overline{R_i}$  as one previously did on the  $x_i$ , pooling when necessary, to obtain a final distinct set of m quantities  $\overline{R_{t_1}}, \ldots, \overline{R_{t_m}}$ . Let N = nk. Then the test statistic is

(9) 
$$T_{k}^{3} = \frac{12n}{N(N+1)} \sum_{j=1}^{m} t_{j} \left[ \overline{R}_{t_{j}} - \frac{N+1}{2} \right]^{2}$$

As shown in [8], H<sub>0</sub> is for large n rejected at level approximately  $\alpha$  if  $T_k^3 \ge C_2$ , where  $C_2$  is determined by (8).

(b') For unequal sample sizes one could again use a conditional test analogous to those mentioned in (a"). This test will be denoted  $T_k^{3c}$ . (c) The underlying model for Terpstra's [13] or Jonckheere's [10] test is

HKU 4

$$x_{ij} = \alpha + \beta i + \varepsilon_{ij}, i = 1, ..., k; j = 1, ..., n_i,$$

where  $\alpha$  and  $\beta$  are unknown constants and where the  $\varepsilon_{ij}$  have a common continuous cdf F. H<sub>0</sub>, H<sub>1</sub> and H<sub>2</sub> now correspond to  $\beta = 0$ ,  $\beta \ge 0$  and  $\beta \le 0$ , respectively. The test statistic will be stated as follows:

(10) 
$$T_{k}^{4} = \sum_{i=1}^{k-1} \sum_{j=i+1}^{k} h_{ij},$$
where  $h_{ij} = \sum_{r=1}^{n_{i}} \sum_{c=1}^{n_{j}} h_{ir,js}$ 

and 
$$h_{ir,js} = \begin{cases} 1 \text{ if } x_{ir} < x_{js} \\ -1 \text{ if } x_{ir} > x_{js} & (i < j). \end{cases}$$

It follows from [13] or [10] that  $T_k^4$  is under  $H_0$  asymptotically normally distributed with zero mean and variance

$$D^{2} = \frac{1}{18} \left\{ N(N+1)(2N+1) - \sum_{i=1}^{k} n_{i}(n_{i}+1)(2n_{i}+1) \right\},$$

where N =  $\sum_{i=1}^{k} n_i$ . Therefore if all  $n_i$  are large, H<sub>0</sub> is rejected at level approximately  $\alpha$  in favor of H<sub>1</sub> if  $T_k^4 \ge t_\alpha D$ , and in favor of H<sub>2</sub> if  $T_k^4 \le -t D$ , where  $t_\alpha$  is the 100(1- $\alpha$ )%-point of a standard normal distribution.

## 2. BEHAVIOR OF THE TESTS UNDER INTERIOR SLIPPAGE

We define interior slippage as follows for model (a):  $\mu_1 = \ldots = \mu_{m-1} = \mu_{m+1} = \ldots = \mu_k = \mu; \ \mu_m = \mu + \Delta$ , with  $\Delta \neq 0$ . The result for Bartholomew's test will be shown in detail for the case of known  $\sigma^2$ .

THEOREM 1. Under interior slippage, when testing H<sub>0</sub> against either H<sub>1</sub> or H<sub>2</sub>, lim  $P[T_k^2 \ge C_2] = 1$  as  $N \to \infty$  with  $n_i/N \ge a > 0$  for i = 1, ..., k. PROOF. It suffices to show that lim  $P[T_k^2 < \infty] = 0$  as  $N \to \infty$  with  $n_i/N \ge a > 0$ . Let us assume the slippage is upward, i.e.  $\Delta > 0$ . HKU 5 (The case  $\Delta < 0$  is similar.)

(i) Testing  $H_0$  against  $H_1$ : To begin with, the probability of complete amalgamation, m = 1, is zero in the limit. This follows immediately from Theorem 1 of [8], which states that a necessary (and sufficient) condition for complete pooling is

(11) 
$$\overline{x}_{[1,j]} > \overline{x}_{[1,k]}$$
 for  $j = 1, ..., k-1$ .

By the consistency of  $\overline{x}_i$  as estimator of  $\mu_i$  we have  $\lim P[\overline{x}_{[1,m-1]} < \overline{x}_{[1,k]}] = 1$ . Therefore there will be, with limiting probability 1, a contribution to  $T_k^2$  from at least the first sample mean. (We can ignore the contribution from  $\overline{x}_{[m,k]}$ .)

Again from the consistency of  $\overline{x}_i$  and of  $\overline{x}_{[1,k]}$  we know that for any  $\varepsilon_1 > 0$  there exists an  $N_1$  such that for  $N \ge N_1$ ,

(12) 
$$\mathbb{P}[|\overline{x}_1 - \mu| > \varepsilon_1] \leq \varepsilon_1,$$

and for any  $\boldsymbol{\epsilon}_{2}$  > 0 there exists an  $\boldsymbol{N}_{2}$  such that for  $\boldsymbol{N}$   $\geq$   $\boldsymbol{N}_{2},$ 

(13) 
$$P[|\overline{\mathbf{x}}_{[1,k]} - \mu| < \Delta a - \varepsilon_2] \leq \varepsilon_2.$$

Hence for  $N \ge max(N_1, N_2)$ ,

(14) 
$$P[|\overline{x}_{1} - \overline{x}_{[1,k]}| < \Delta a - \varepsilon_{1} - \varepsilon_{2}] \leq \varepsilon_{1} + \varepsilon_{2},$$

and therefore lim  $P[T_k^2 > n_1(\Delta a - \varepsilon_1 - \varepsilon_2)^2 / \sigma^2] = 1$ . Since  $\Delta$ , a,  $\varepsilon_1$ ,  $\varepsilon_2$ , and  $\sigma^2$  are constant, this proves the theorem for (i).

(ii) Testing  $H_0$  against  $H_2$ : An analogous proof holds in this case. The necessary (and sufficient) condition for complete analgamation is now

(15) 
$$\overline{x}_{[1,j]} < \overline{x}_{[1,k]}$$
 for  $j = 1, ..., k-1$ ,

and we now have  $\lim P[\bar{x}_{[1,m]} > \bar{x}_{[1,k]}] = 1$ . We eventually obtain  $\lim P[\bar{T}_k^2 > n_k(\Delta a - \varepsilon_1 - \varepsilon_2)^2 / \sigma^2] = 1$ , which proves it for (ii).

COROLLARY 1. Under interior slippage, when testing  $H_0$  against either  $H_1$  or  $H_2$ , lim  $P[T_k^1 \ge C_1] = 1$  as  $N \to \infty$  with  $n_i/N \ge a > 0$  for i = 1, ..., k. HKU 6 The proof is similar to that of Theorem 1. One shows that lim  $P[T_{b_r}^1 \le a/g] = 0$ , where  $g = \Delta^2/a^2$ .

COROLLARY 2. Under interior slippage, when testing  $H_0$  against either  $H_1$  or  $H_2$ , lim  $P[T_k^{jc} \ge C_j] = 1$  for j = 1, 2 as  $N \rightarrow \infty$  with  $n_i/N \ge a > 0$ ,  $i = 1, \ldots, k$ .

For model (b), interior slippage means  $F_1(x) = \ldots = F_{m-1}(x) = F_{m+1}(x) = \ldots = F_k(x) = F(x)$ ,  $F_m(x) = F(x+\theta)$ ,  $\theta > 0$  (for upward slippage). The equivalent result for Chacko's test follows.

THEOREM 2. Under interior slippage, when testing H<sub>0</sub> against H<sub>1</sub> or H<sub>2</sub>, lim  $P[T_k^3 \ge C_2] = 1$  as N  $\rightarrow \infty$ .

PROOF. Because of the equal sample sizes, the m<sup>th</sup> population will be entitled to a mean rank of [(N+1)/2] (1+ $\Delta$ ) and the others, to mean ranks of  $[(N+1)/2][1 - \Delta/(k-1)]$  in an overall sample of N, for some  $0 < \Delta < 1$ . For testing H<sub>0</sub> against H<sub>1</sub>, a proof similar to that of Theorem 1 establishes that for any  $\varepsilon > 0$  there exists an N<sub>0</sub> such that for N  $\geq$  N<sub>0</sub>,

(16) 
$$P\left[\left|\overline{R}_{1} - \frac{N+1}{2}\right| < \frac{N+1}{2} \frac{\Delta}{k-1} - \varepsilon\right] \leq \varepsilon,$$

and hence

(17) 
$$\lim P\left\{T_k^3 > \frac{12n}{N(N+1)} \left[\frac{N+1}{2} \frac{\Delta}{k-1} - \varepsilon\right]^2\right\} = 1,$$

or, in effect,  $\lim P[T_k^3 > dn] = 1$  for some constant d > 0. This proves Theorem 2 for the case of H<sub>0</sub> against H<sub>1</sub>. The proof for the case H<sub>0</sub> against H<sub>2</sub> is similar.

COROLLARY 1. Under interior slippage, when testing H<sub>O</sub> against H<sub>1</sub> or H<sub>2</sub>, lim  $P[T_k^{3c} \ge C_2] = 1$  as N  $\rightarrow \infty$ .

For model (c), interior slippage can be defined by  $x_{ij} = \alpha + \varepsilon_{ij}$ , i = 1, ..., m-1, m+1, ..., k, and  $x_{mj} = \alpha + \Delta + \varepsilon_{ij}$ , with  $\Delta > 0$  for upward slippage.

It is easiest here to consider the expected value of the test statistic, HKU 7

(18) 
$$ET_{k}^{4} = n_{m} \left( \sum_{i=1}^{m-1} n_{i} - \sum_{i=m+1}^{k} n_{i} \right).$$

In particular, for the case of equal sample sizes,

(19) 
$$ET_k^4 = n(2m - k - 1).$$

We therefore see at once that  $\mathrm{ET}_{\mathbf{k}}^4 = 0$  if  $\mathbf{m} = (\mathbf{k}+1)/2$  for  $\mathbf{k}$  odd. For equal sample sizes it can be shown, analogously to the earlier proofs, that if  $\mathbf{m} < (\mathbf{k}+1)/2$ , the limit of the probability of rejecting  $\mathrm{H}_0$  in favor of  $\mathrm{H}_2$  is 1, and hence the limit of the probability of rejecting  $\mathrm{H}_0$  in favor of  $\mathrm{H}_1$  is zero. The opposite holds for  $\mathbf{m} > (\mathbf{k}+1)/2$ . These situations are, in turn, reversed for  $\Delta < 0$ .

For unequal sample sizes, no such statement can be made. The probabilities depend on whether  $\sum_{i=1}^{m-1} n_i < \sum_{i=m+1}^{k} n_i$  or vice versa.

#### 3. COMMENTS

In Section 2 we dealt with what Mosteller [12] calls "the error of the third kind": rejecting  $H_0$  correctly, but for the wrong reason. It is of course a matter of opinion whether it is worse in general to accept or reject  $H_0$  under such circumstances. However, in particular cases it does seem reasonable "to prefer one wrong decision over the other".

For example, if a population has slipped upward for m > (k+1)/2, particularly for m = k-1, it seems worse to reject  $H_0$  in favor of  $H_2$  than to reject it in favor of  $H_1$ ; the opposite situation would hold for m < (k+1)/2, especially m = 2. The Terpstra test will only make the less poor decision in these cases, at least for reasonably equal sample sizes. The other tests can make both.

If the central population has slipped, m = (k+1)/2, it seems preferable to accept H<sub>0</sub>, since no "trend" of any sort can possibly be claimed. Terpstra's test will accept H<sub>0</sub> here, not only for equal sample sizes but whenever  $\sum_{i=1}^{m-1} n_i = \sum_{i=m+1}^{k} n_i$ .

HKU 8

 $T_k^4$  is of course a test designed for a more fully specified model than  $T_k^j$ , j = 1, 2, 3. As shown in [3], it should have better power than the tests against ordering when there is a (reasonably) linear trend, and poorer power when there is considerable variation in the differences between successive means. (Asymptotic efficiency comparisons have not been possible because the tests have different limiting distributions.) However, for k = 3 or 4, the difference in power between  $T_k^4$  and the other tests appears to be quite small.

 $T_k^4$  cannot be used unless a complete *a priori* ranking of the µ's is feasible. But if this can be done, then in view of its comparable power and the protection which it offers against really bad slippage decisions,  $T_k^4$  would seem to be preferable for k = 3 and 4, and possibly 5. This is particularly so for k = 3, since here any interior slippage must be that of the central population.

It may also be noted that  $T_k^4$ , unlike  $T_k^1$  and  $T_k^2$ , does not presuppose underlying normal distributions and, unlike  $T_k^3$ , can be used with unequal  $n_i$ 's. (The use of the three conditional tests mentioned may be hard to justify, since the value of m obtained is clearly not irrelevant for deciding between  $H_0$  and  $H_1$  or  $H_2$ ; also, some power studies in [2] indicate that these tests have distinctly lower power than their unconditional counterparts.)

Incidentally, the two tests given by Whitney [14] for k = 3 show the same behavior under interior slippage as  $T_{k}^{4}$ .

Finally, when several interior populations have slipped in the same direction but by arbitrary amounts, the preceding results will of course hold true. When some have slipped upwards and some downwards, it is possible to have a "weighted mean slippage" of zero; in this case only one of the alternative hypotheses would be rejected. In general, however, both H<sub>1</sub> and H<sub>2</sub> would again be rejected with probability 1 in the limit.

ACKNOWLEDGEMENTS. I wish to express my thanks to Dr. Alvin Wiggins and to Dr. Robert Elashoff for their helpful comments.

HKU 9

REFERENCES

- [1] BARTHOLOMEW, D.J. (1959). A test of homogeneity for ordered alternatives. Biometrika <u>46</u>, 36-48.
- [2] BARTHOLOMEW, D.J. (1961). A test of homogeneity of means under restricted alternatives. J. Roy. Statist. Soc. Ser. B 23, 239-281.
- [3] BARTHOLOMEW, D.J. (1961). Ordered tests in the analysis of variance. Biometrika 48, 325-332.
- [4] BARTON, D.E. and MALLOWS, C.L. (1961). The randomization bases of the problem of the amalgamation of weighted means. J. Roy. Statist. Soc. Ser. B <u>23</u>, 423-433.
- 5 BRUNK, H.D. (1955). Maximum likelihood estimates of monotone parameters. Ann. Math. Statist. <u>26</u>, 607-616.
- [6] BRUNK, H.D. (1958). On the estimation of parameters restricted by inequalities. Ann. Math. Statist. 29, 437-454.
- [7] CHACKO, V.J. (1959). Testing homogeneity against ordered alternatives. Ph.D. thesis, University of California, Berkeley.
- [8] CHACKO, V.J. (1963). Testing homogeneity against ordered alternatives. Ann. Math. Statist. 34, 945-956.
- [9] VAN EEDEN, C. (1957). Maximum likelihood estimation of partially or completely ordered parameters. Indag. Math. <u>19</u>, 128-211.

[10] JONCKHEERE, A.R. (1954). A distribution free <u>k</u> sample test against ordered alternatives. Biometrika <u>41</u>, 135-145.

[11] MILES, R.E. (1959). The complete amalgamation into blocks by weighted means of a finite set of real numbers. Biometrika <u>46</u>, 317-327.

HKU 10

[12] MOSTELLER, F. (1948). A k-sample slippage test for an extreme population. Ann. Math. Statist. <u>19</u>, 58-65.

.

[13] TERPSTRA, T.J. (1952). The asymptotic normality and consistency of Kendall's test against trend, when ties are present in one ranking. Indag. Math. <u>14</u>, 327-333.

[14] WHITNEY, D.R. (1951). A bivariate extension of the U statistic. Ann. Math. Statist. <u>22</u>, 272-284.

HKU 11

# THE DEVELOPMENT OF THE CONCEPT OF STATISTICAL DECISION THEORY By D. Bierlein, Germany Technische Hochschule, Karlsruhe

1. A. WALD'S MODEL AND ITS PRACTICAL APPLICATION

What is a "good" statistical procedure? - Up to the middle of the 20<sup>th</sup> century the answer to this question generally was given by an ad-hoc definition of an "optimum" property for a special case. Dissatisfaction with that divergent tendency has led A. Wald (about 1939) to the outline of statistical decision theory. Its concern is a uniform model for quite different statistical decision problems, with an accent on the consideration of the consequences of every decision. The model for a general statistical decision situation is - as is well known - characterized by the following data:

- sample space M (= set of potential results of the planned observation of a random vector *Nt*)
- 2) set  $\int$  of distribution functions which have to be taken into consideration for  $\mathcal{M}$ .
- 3) set D of possible decisions
- loss function W by which for every element of the consequence of a decision is evaluated.

A statistical procedure  $\delta$  for a decision problem formalized in this way is a map from M into D and may be interpreted as a strategy which attaches a decision to every potential information. If, concerning  $\delta$ , you make the additional assumption that for every  $F \in \mathcal{F}$  the expectation of the loss function exists - and therefore the risk function  $r_{\delta} | \mathcal{F}$ you get to the definition of the *statistical decision function*. Let  $\Delta$ be the set of all decision functions of the decision problem (M, $\mathcal{F}$ ,D,W). The question which statistical procedure from  $\Delta$  should be used is reduced to the choice of a suitable principle of optimality. The main principles in statistical decision theory are the dominance principle, related to classes of unbiased or invariant procedures, and the

DB 1

minimax principle. The task of mathematical statistics is to find explicitly, for as extensive classes of statistical decision problems as possible, the optimal procedures based on the most important principles of optimality or at least to give a practicable algorithm for finding them. For the practical statistician or for the consumer of statistics in economics, engineering and empirical sciences there remains the task 'to give. in an adequate manner, a concrete value to the parameters M, f, D and W of the statistical decision problem. This task meets some difficulties: for example it strikes against the old custom of many experimental scientists to use rigorously  $\alpha$ -level tests with a value  $\alpha$  which depends only on the scientist's special branch. But, there are even men in statistical practice who increasingly regret that an  $\alpha$ -level test does not give any concrete recommendation how to act. This is just recently pointed out again by J. Wolfowitz<sup>1</sup>). My own experience as statistical consultant mainly to engineers encourages me to hope that in a large area of application the consumers of statistics will rather soon accept the framework of decision theory. As soon as mathematical statistics will present a way with or without the help of a computer - to calculate the optimal procedure, in practice most of the classical standard methods will be replaced by more adequate procedures of statistical decision theory. However, many of the classical procedures, which do not fit (or fit only in an insatisfactory manner) into the scheme of statistical decision theory, will remain practically important. In the first place that seems to be true for the methods of correlation and regression theory and, generally, for descriptive statistics. But in the model of a statistical decision problem they have their legitimate place too: In addition to intuition and special scientific investigations, they are useful for making precise the datum f (in several cases also D) of the statistical decision problem. That takes place in a preceding analysis of a separated sample with the help of descriptive statistics. In general the preceding inspections and considerations supply more information than may be inferred for the formal decision problem by a suitable choice of f, but not enough information to be able to determine a precise a-priori

J. Wolfowitz: Remarks on the theory of testing hypotheses. The New York Statistician 18 (1967).

probability à la Bayes. This situation gave rise to several modifications of statistical decision theory: Hodges and Lehmann, for example, offer the "Restricted Bayes Solution" as an alternative to the minimax solution <sup>2</sup>). In the following we will discuss a generalization of Wald's decision model where **f** is replaced by the set  $\phi$  of all such probability measures on **f** which are compatible with the "pre-statistical information" of the decision maker (R. Bartoszýnski, D. Bierlein, O. Bunke, H. Richter and others).

## 2. EXAMPLES

How important a pre-statistical information may be in practice, will be demonstrated by means of some examples of point estimation of a probability:

Consider the situation that for a river there are constructed new installations regulating the stream, for instance a new dam. Let A be the event that a flood of well-defined power will occur at least once during a time-interval of ten years and put x = Pr(A). It is required to estimate x as soon as possible, say after 10, 20 or 30 years; i.e. it is not possible to make a sample of size greater than 1, 2 or 3. But beyond that there may be given a number of flood-observations which were made during many decades before the construction of the new installations. Add to this the hydrologically established knowledge, that Pr(A) is not increased by the new construction, perhaps is even decreased. Both these facts - the knowledge of results of former observations made under similar conditions and the use of theoretical considerations allow to estimate a subjective probability  $\phi$  for that x does not exceed a certain bound x.

$$\phi(\mathbf{x} \leq \mathbf{x}_{0}) = \phi_{0} \text{ or } \phi(\mathbf{x} \leq \mathbf{x}_{0}) \geq \phi_{0}.$$

Already these crude specifications of  $\phi$  can represent a pre-statistical information diminishing the minimax risk.

A quite similar situation occurs in medium range or long-term weather-

<sup>&</sup>lt;sup>2</sup>) J.L. Hodges and E.L. Lehmann: The use of previous experience in reaching statistical decisions. Ann. Math. Stat. 23 (1952).

*forecast*. To be able to make observations in order to test a forecasting rule one has to wait for the - in general rare - years where the conditional meteorological situation, to which the rule attaches a forecast, is obtained. For that reason the pre-statistical information which results from theoretical meteorology and from observations made under similar situations should not be disregarded.

Further actual examples can be found in *techniques of astronautics*. Here the sample size is bounded drastically by high costs of experiments and by time limitation. The information resulting from former experiments with older construction units and a knowledge about the tendency of the effects of technical innovations should - compared with the information resulting from sampling - be of some importance for the decision. In this example too, it should not be difficult to express the pre-statistical information as a system of subjective probabilities for a certain number of subsets  $A_i$  of the parameter interval [0, 1].

3. **<b>\phi-OPTIMAL** PROCEDURES

. .

To develop a general model for these examples, let us assume that for certain subsets  $A_i$ ,  $B_k$  of f the subjective probabilities them-selves or bounds for them are known:

$$\phi(A_i) = \phi_i \text{ for } i \in J_1, \quad \phi(B_k) \leq \psi_k \text{ for } k \in J_2.$$

By these data  $\phi$  is, in general, not made suffciently precise in order to form the risk expectation  $\int_{\mathcal{F}} r_{\delta}(F) d\phi$ . The set of all precise apriori probabilities is the set P of all probability measures p, for which

$$\mathbf{r}_{\delta}(\mathbf{p}) := \int_{\mathbf{f}} \mathbf{r}_{\delta}(\mathbf{F}) d\mathbf{p}$$

exists for all  $\delta \in \Delta$ . Under these measures p the elements of the subset

$$\phi: = \{ p \in P: p(A_i) = \phi_i \text{ for } i \in J_1, p(B_k) < \psi_k \text{ for } k \in J_2 \}$$

DB 4

are compatible with  $\phi$ . For that  $\phi$  is not empty, the data for  $\phi$  must, of course, not be self-contradictory.

In the generalized decision theory other non-empty subsets of P are also admitted as pre-statistical informations.

For a generalized statistical decision situation, characterized by the data M,  $\phi$ ,D and W, a statistical decisions function  $\delta^{\ddagger}$  is called  $\phi$ -optimal, if

$$\sup_{\delta} \mathbf{r}_{\delta}^{\star}(\mathbf{p}) = \min_{\delta} \sup_{\delta} \mathbf{r}_{\delta}(\mathbf{p}),$$

$$p \boldsymbol{\epsilon} \Phi \qquad \delta \boldsymbol{\epsilon} \Delta \mathbf{p} \boldsymbol{\epsilon} \Phi$$

. . . . . .

i.e.  $\delta^{\pm}$  is minimax strategy of player 2 in the zero-sum two-person game  $(\phi, \Delta, r_{\delta}(p))$ . The term

$$\mathbf{v}^{\bigstar}(\Phi, \Delta)^{\cdot} = \inf \sup_{\delta \in \Delta} \mathbf{r}_{\delta}(\mathbf{p}),$$

i.e. the upper value of the game  $(\Phi, \Delta, r_{\delta}(\mathbf{p}))$ , is called  $\Phi$ -minimax risk. A  $\Phi$ -optimal procedure guarantees that the expected loss  $r_{\delta}(\mathbf{p})$  does not exceed the bound  $v^{*}(\Phi, \Delta)$ .

#### 4. EFFECTIVITY OF AN INFORMATION $\boldsymbol{\varphi}$

At first there is the question: Under what conditions is it possible to take advantage of a pre-statistical information, i.e. under what conditions is the minimax risk  $v^{\bigstar}(\Phi, \Delta)$  less than  $v^{\bigstar}(P, \Delta)$ , which is equal to the minimax risk of the decision problem without using a pre-statistical information (or - more precisely - with respect to the trivial information  $\phi(\varsigma) = 1$ )?

We say  $\Phi$  is *effective*, if

 $v^{\bigstar}(\Phi, \Delta) < v^{\bigstar}(P, \Delta),$ 

and define as *effectivity* of  $\Phi$  the term

Eff 
$$(\phi)$$
: =  $\frac{\mathbf{v}^{\mathbf{x}}(\mathbf{P},\Delta) - \mathbf{v}^{\mathbf{x}}(\phi,\Delta)}{\mathbf{v}^{\mathbf{x}}(\mathbf{P},\Delta) - \mathbf{r}}$ 

where  $r_{\underline{\star}}$  is the infimum of the risk function  $r_{\delta}(p)$  on  ${}_{\Delta} \not x$  P.  $^3). It is easy to see that$ 

$$0 = Eff(P) \leq Eff(\Phi) \leq Eff(\{F_0\}) = 1$$

for all  $F \in \mathcal{L}$  with

$$\inf_{\mathbf{d}\in D} W(\mathbf{F}_{\mathbf{0}}, \mathbf{d}) = \mathbf{r}_{\mathbf{x}}.$$

Necessary and sufficient conditions for  $\Phi$  to be effective are formulated <sup>4</sup>) making use of Wald's <sup>5</sup>) intrinsic metric on the space P and of the condition that the games (P, $\Delta$ ,r) and ( $\Phi$ , $\Delta$ ,r), respectively, are strictly determined.

Because of this aspect - and in another connection - criteria for that  $(\Phi, \Delta, \mathbf{r})$  is strictly determined are interesting. With the help of a minimax theorem of Ky Fan<sup>6</sup>) such criteria are offered in a form which allowes applications to important special cases of point estimation<sup>4</sup>).

One may ask about the relation between effective pre-statistical information and Bayesian a-priori probabilities. The guess, that every precise a-priori probability is an effective pre-statistical information, is wrong.

Indeed, if  $(P, \Delta, r)$  is strictly determined, then for every minimax strategy  $p^{\bigstar}$  of "player 1"

$$v^{\bigstar}(\{p^{\bigstar}\}, \Delta) = v^{\bigstar}(P, \Delta)$$

and consequently

 $Eff(\{p^{\ddagger}\}) = 0;$ 

thus the information  $p^*$  is not effective.

<sup>&</sup>lt;sup>3</sup>)  $r_{\pm}$  is zero, if the loss function W is reduced, that is, if inf W(F,d)=0 for all F  $\epsilon_{f}$ . d $\epsilon_{D}$ 

On the other side, just that a-priori probability which is not seldom taken to be equivalent with absolute ignorance, namely the rectangular distribution L, is effective by all means: In the example of point estimation of a probability with quadratic loss function and sample size 1 you have

Eff ({L}) = 
$$\frac{1}{9}$$
.

5. HOW TO GAIN A &-OPTIMAL PROCEDURE

The statistical decision theory often meets the reproof that the gap between theory and application is rather large. Is not this gap increased further by including an additional parameter into the model of a decision situation? - Just in the practically important case of a quite vague pre-statistical information a  $\Phi$ -optimal procedure can be found in a relatively convenient manner. That may be demonstrated on the example of point estimation of a probability x: Let W be a quadratic loss function, and a pre-statistical information be given in the form

 $\phi(\mathbf{x} \leq \tau) = \lambda.$ 

Then in the first instance you can use tables or diagrams of those pairs  $(\tau, \lambda)$  for which the information

 $\phi(\mathbf{x} \leq \tau) = \lambda$ 

is effective. If the pre-statistical information is effective, you can find the  $\Phi$ -optimal estimator in a table, which till now has been computed for sample sizes 1 and 2 and is being prepared for sample sizes 3 and 4.

33

<sup>&</sup>lt;sup>4</sup>) D. Bierlein: Zur Einbeziehung der Erfahrung in spieltheoretische Modelle. Op. Res. Verf. III(1967).

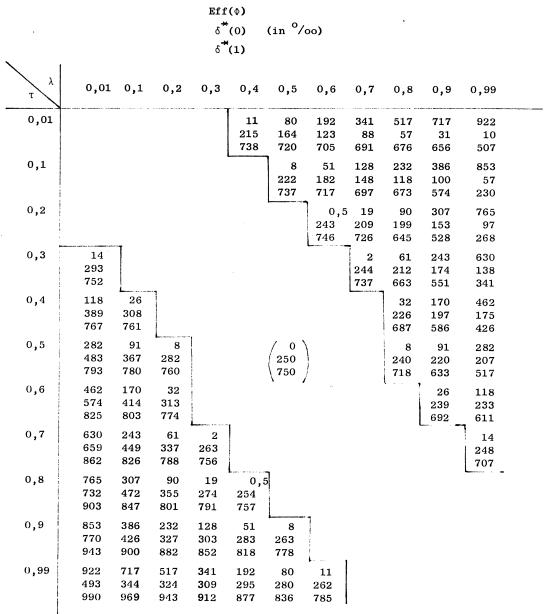
<sup>&</sup>lt;sup>5</sup>) A. Wald: Statistical decision functions. Wiley (1950).

<sup>&</sup>lt;sup>6</sup>) Ky Fan: Minimax theorems. Proc. Nat. Ac. Sc. 39 (1953).

The effectivity of  $\Phi$  and the  $\Phi$ -optimal estimator  $\delta^*|M$  as a function of  $\tau$  and  $\lambda$ .

Sample size 1

. .



۰

DB 8

Sample size 2

.

- Eff(φ)  $δ^{*}(0) (in <sup>0</sup>/00)$   $δ^{*}(1)$   $δ^{*}(2)$

1				د <b>* (2</b> )	)							
				• (=)								
τλ	0,01	0,1	0,2	0,3	0,4	0,5	0,6	0,7	0,8	0,9	0,99	
0,01				22	91	190	310	445	593	752	950	
				162	122	92	69	49	33	19	11	
				493	484	475	467	458	446	424	131	
				780	768	758	750	742	736 `	727	711	
0,1					5	33	83	164	317	550	870	
0,1					186	155	131	117	110	85	58	
					493	480	461	406	319	237	113	
					787	776	765	756	738	704	541	
		٦										
0,2	2						21	92	219	414	714	
	221						185	162	141	122	96	
	497						454	401	343	274	200	
	796						778	757	730	680	440	
0,3	63	2	1					36	117	259	552	
	307	225						185	169	155	126	
	497	501						438	387	333	300	
	808	797						765	735	684	411	
0,4	168	25							35	119	404	
-, -	390	266	1						192	179	151	
	516	511							441	400	400	
	819	805							756	701	455	
0,5	285	58	0,2	5		( 0)			0,2	2 58	285	
0,5	470	298	212	1		0 207 500	١		206	187	169	
	546	529	502			500			498	471	454	
	831	813	794			793	/		788	702	530	
0.0						· · · /						
0,6	404	119	35							25	168	
*.	545	299	244							195	181	
	600 849	600 821	559 808							489 734	484 610	
	045											
0,7	552	259	117	36						2	63	
	589	316	265	235						203	192	
	700	667	613	562						499	503	
	874	845	831	815						775	693	
0,8	714	414	219	92	21						2	
•	560	320	270	243	222						204	
	800	726	657	599	546						503	
	904	878	859	838	815	1					779	
0,9	870	550	317	164	83	33	5	-			1 <u></u>	
- ,-	459	296	262	244	235	224	213					
	887	763	681	594	539	520	507					
	942	915	890	883	869	845	814					
0,99												
0,99	950	752	593 264	445	310	190	91 222	22				
	289	273 576	264 554	258 542	250	242	232	220				
	869 989	576 981	554 967	542 951	533 931	525 908	516 878	507 838				
	505	201	507	901	221	200	010	030				
	•										DB 9	

\*

35

.

# METHODS OF INVESTIGATING WHETHER A REGRESSION RELATIONSHIP IS CONSTANT OVER TIME

by R.L. Brown and J. Durbin (England) Central Statistical Office, London, and London School of Economics and Political Science

### 1. INTRODUCTION

Regression analysis of time-series data is usually based on the assumption that the regression relationship is constant over time. In some applications, particularly in the social and economic field, the validity of this assumption is open to question and it is important that methods of detecting and allowing for changes should be included in the analysis.

In this paper we consider a number of techniques for detecting departures from constancy. Although we shall present several formal tests of significance our approach is essentially that of Tukey's data analysis (Tukey, 1962), that is we try to develop techniques which bring out in a graphic way whatever departures from constancy are present in the data rather than parametrise in advance particular types of departure and develop formal significance tests which have high power against these particular alternatives.

The present paper should be regarded as a preliminary report on and summary of our work on this subject, a full acount of which will be published later (Brown and Durbin, 1969). The later paper will contain proofs of theoretical results and applications to real and artificial examples.

The basic regression model we are concerned with is

$$y_{+} = x_{+}^{\dagger}\beta_{+} + u_{+}, \quad t = 1, \dots, T$$
 (1)

where  $x_t$  is the column vector of observations at time t on each of k regressors,  $\beta_t$  is the vector of regression coefficients, where we have attached the suffix t to indicate that  $\beta_t$  may not be constant, and  $u_1$ , ...,  $u_T$  are independent normal variables with zero means and variances

 $\sigma_1^2$ , ...,  $\sigma_T^2$ . The first element in each  $x_t$  is unity, representing the constant term in the model, and the remaining elements are assumed to be non-stochastic. Thus autoregressive and other models containing lagged y's are excluded from consideration. The hypothesis  $H_0$  we wish to investigate is  $\beta_1 = \ldots = \beta_T = \beta$  and  $\sigma_1^2 = \ldots = \sigma_T^2 = \sigma^2$ ; we are, however, more concerned about departures from equality among the  $\beta$ 's than among the  $\sigma$ 's.

## 2. METHODS BASED ON LEAST-SQUARES RESIDUALS

Assuming H<sub>0</sub> to be true, let b denote the least-squares estimate of  $\beta$ , i.e.  $b = \begin{bmatrix} T \\ 1 \\ 1 \\ 1 \end{bmatrix}^{-1} \begin{bmatrix} T \\ 1 \\ 1 \\ 1 \end{bmatrix}^{-1} \begin{bmatrix} T \\ 1 \end{bmatrix}^{-1} \begin{bmatrix} T \\ 1 \end{bmatrix}^{-1} \begin{bmatrix} T \\ 1 \\ 1 \end{bmatrix}^{-1} \begin{bmatrix} T \\ 1$ 

In this respect the problem resembles that of detecting changes in the mean in industrial quality control for which the cumulative sum or cusum technique, introduced by Page (1954) and discussed further by Barnard (1959) and by Woodward and Goldsmith (1964), has been found to be a more effective tool for detecting small changes than the ordinary control chart. This suggests that instead of plotting out the individual  $z_t$  the cusums  $Z_r = \frac{1}{s} \int_{1}^{r} z_t$ ,  $r = 1, \ldots, T$  should be plotted, where we have divided by the estimated standard deviation s to eliminate the irrelevant scale factor. The difficulty about this suggestion is that there seems to be no way of assessing the significance of the departure of the observed graph of  $Z_r$  against r from the mean-value line  $E(Z_r) = 0$ . The intractability of the problem arises from the fact that in general the covariance function  $E(Z_r Z_s)$  does not reduce to a

form that is manageable by standard Gaussian-process techniques (c.f. Mehr and McFadden, 1965). For instance, for the simple case of regression through the origin on a linear time trend the covariance function is asymptotically  $r - 3r^2s^2/4T^3$  (r < s) which is an unmanageable form.

An alternative is to consider the cusum of squares  $\frac{1}{s^2} \sum_{1}^{r} z_t^2$ . Although more tractable, this is still fairly difficult to deal with and is hard to interpret. Instead of considering this we prefer to make the transformation given in the following section which enables us to treat the problem in terms of cusums and cusums of squares of independent  $N(0,\sigma^2)$  variables.

### 3. METHODS BASED ON RECURSIVE RESIDUALS

Let  $b_{\mbox{\bf r}}$  be the least-squares estimate of  $\beta$  based on the first r observations and let

$$w_{r} = \frac{y_{r} - x_{r}'b_{r-1}}{\sqrt{1 + x_{r}'(x_{r-1}' - x_{r-1})^{-1} x_{r}}}, r = k+1, ..., T$$
(2)

where  $X'_{r-1} = [x_1, \ldots, x_{r-1}]$ . It can be shown that  $w_{k+1}, \ldots, w_T$  are independent  $N(0, \sigma^2)$ . These quantities are easy to obtain recursively on a modern computer without the necessity of repeated matrix inversions in virtue of the relations

$$b_{r} = b_{r-1} + (X'_{r}X_{r})^{-1}x'_{r}(y_{r} - x'_{r}b_{r-1})$$
(3)

and

$$(X'_{r}X_{r})^{-1} = (X'_{r-1}X_{r-1})^{-1} - \frac{(X'_{r-1}X_{r-1})^{-1} x_{r}x'_{r}(X'_{r-1}X_{r-1})^{-1}}{1 + x'_{r}(X'_{r-1}X_{r-1})^{-1}x_{r}} .$$
(4)

Denoting by  $S_r$  the residual sum of squares after fitting the model from the first r observations, we have the further relation

$$S_r = S_{r-1} + w_r^2$$
 (5)

To avoid difficulties due to ill-conditioning of the matrices  $X'_rx_r$ , it is recommended that all elements of  $x_t$  except the value unity in the leading position should be replaced by their differences from the overall sample mean.

(2) is a generalisation of the regression model of the Helmert transformation. Interesting applications of (3) and (4) to the fitting of regression models in the frequency domain and to the fitting of nonlinear models are given by Duncan and Jones (1966) and by Walker and Duncan (1967). The basic relation (4) which enables repeated matrix inversions to be avoided is due to Bartlett (1951).

If  $\beta_t$  is constant up to time  $t = t_0$  and differs from this constant value from then on, the w<sub>r</sub>'s will have mean zero up to  $r = t_0$  but in general will have non-zero means subsequently. This suggests that plotting the cusum quantity

$$W_{r} = \frac{1}{s} \sum_{k=1}^{r} W_{j}, r = k+1, ..., T$$
 (6)

against r should be a useful technique for detecting changes in  $\beta_t$ . As previously, s denotes the estimated standard deviation determined by  $s^2 = S_{\pi}/(T - k)$ .

We require a method of testing the significance of the departure of the sample path of  $W_r$  from its mean-value line  $W_r = 0$ . A suitable procedure is to find a pair of lines lying symmetrically above and below the line  $W_r = 0$  such that the probability of crossing one or both lines is  $\alpha$ , the required significance level. Since the variance of  $W_r$  increases with r, it is clearly desirable to choose the slopes of the lines so that they diverge with increasing r. We suggest taking slopes such that the probability that the point  $(r, W_r)$  lies outside the lines is a maximum for r half-way between r = k and r = T. This leads to the lines joining the points  $(k, \pm a\sqrt{T-k})$  and  $(T, \pm 3a\sqrt{T-k})$  where a is determined by the equation

(3a) + 
$$e^{-4a^2}$$
 {1 -  $\Phi(a)$ } =  $\frac{1}{2}\alpha$  (7)

in which  $\Phi(z) = \frac{1}{\sqrt{2\pi}} \int_{z}^{\infty} e^{-\frac{1}{2}u^{2}} du$ .

These results are obtained by approximating  $W_r$  by the continuous Gaussian process having the same mean and covariance functions. We have assumed that the probability that a particular sample path of  $W_r$  crosses both lines is negligible which will be justifiable for values of  $\alpha$  normally used for significance testing, say 0.1 or less. Useful values of a are:  $\alpha = 0.05$ , a = 0.948;  $\alpha = 0.01$ , a = 1.143.

We believe that the proper function of these lines is to provide a yardstick against which to assess the observed pattern of the sample path, though of course they can be used to provide a formal test of significance by rejecting if the sample path travels outside the region between the lines.

Another useful plot is that of the cusum of squares

$$s_{r} = \frac{\sum_{k=1}^{r} w_{t}^{2}}{\sum_{k=1}^{T} w_{t}^{2}} = \frac{s_{r}}{s_{T}}$$
,  $r = k+1, ..., T$  (8)

where we have standardised by dividing by the overall residual sum of squares  $S_T$ . On  $H_0$ ,  $E(s_r) = \frac{r-k}{T-k}$ . This suggests drawing a pair of lines  $s_r = \pm c_0 + \frac{r-k}{T-k}$  on the diagram parallel to the mean-value line such that the probability that the sample path crosses one or both lines is  $\alpha$ , the significance level.

We are able to obtain the required significance values  $c_0$  from the theory of Kolmogorov-Smirnov statistics in the study of the sample distribution function. This possibility arises because when T - k is even the joint distribution of  $s_{k+2}$ ,  $s_{k+4}$ , ...,  $s_{T+k-2}$  is the same as that of an ordered sample of  $\frac{1}{2}$  (T - k) - 1 independent observations from the uniform (0,1) distribution. Let

$$c^{+} = \max_{j=1,...,m-1} (s_{k+2j} - \frac{j}{m})$$

$$c^{-} = \max_{j=1,...,m-1} (\frac{j}{m} - s_{k+2j})$$

where  $m = \frac{1}{2} (T - k)$ . Then  $c^+$  and  $c^-$  have the same distribution and are distributed as Pyke's modified Kolmogorov-Smirnov statistic  $C_n^+$  with n = m - 1; significance values have been computed by C.E. Rogers and are tabulated in Table 1 of Durbin (1969). We suggest that these values are used as approximations to the significance values of

$$\max_{\substack{i=1,...,T-k-1}} (s_{k+i} - \frac{i}{T-k}) \text{ and}$$

$$\max_{\substack{i=1,...,T-k-1}} (\frac{i}{T-k} - s_{k+i}),$$

entering the table at  $m = \frac{1}{2} (T - k)$  when T - k is even and interpolating linearly between the values for  $m = \frac{1}{2} (T - k) - \frac{1}{2}$  and  $m = \frac{1}{2} (T - k) + \frac{1}{2}$ when T - k is odd. Our expectation is that the approximation should be good unless T - k is small.

Let  $c_0$  be the significance value obtained in this way corresponding to significance level  $\frac{1}{2} \alpha$ . The pair of lines  $s_r = \frac{1}{2} c_0 + \frac{r-k}{T-k}$  are then drawn on the diagram plotting  $s_r$  against r. Since for the values of  $\alpha$ normally used, say 0.1 or less, the probability of crossing both lines is negligible, we may take  $\alpha$  as the probability of crossing either line.

It is sometimes appropriate to consider a one-sided test. For example, if it is assumed that  $\beta_t = \beta^*$  for  $t \leq r$  and  $\beta_t = \beta^{**} \neq \beta^*$  for t > r while  $\sigma_t^2 = \sigma^2$  for all t, then  $E(w_t^2) = \sigma^2$  for  $t \leq r$  and  $E(w_t^2) > \sigma^2$ for t > r. One would therefore expect the departure from the null hypothesis to be indicated by a tendency for the sample path of  $s_r$  to lie below the mean-value line, and would therefore use a one-sided test. For this purpose, one would take the significance value of  $c_0$  corresponding to significance level  $\alpha$ , not  $\frac{1}{2} \alpha$ .

But whether the two-sided or one-sided situations are envisaged we ourselves prefer to regard the lines constructed in this way as yardsticks against which to assess the observed sample path rather than as providing formal tests of significance.

These procedures based on the plot of  $s_r$  represent a development of a test of constancy proposed by Durbin (1960).

Further useful plots are obtained by graphing the components of  $b_r$  against r. If the regression relationship does indeed vary over time, these plots may serve to identify the source of the variation.

Finally, we remark that an alternative set of plots can be obtained by running the analysis backwards through time instead of forwards through time. The pictures provided by the two plots will differ according to where along the time scale any variation in the regression relationship takes place. Since both analyses are informative, we suggest that both should be carried out.

### 4. MOVING REGRESSIONS

Another useful way of investigating the time-variation of  $\beta_t$  and  $\sigma_t^2$  is to plot out the estimated regression coefficients and residual variance obtained from a segment of  $\ell$  successive observations, this segment being moved along the time scale. The significance of differences over time can then be assessed by a variant of the ordinary analysis-of-variance test for non-overlapping groups. This technique will be considered in Brown and Durbin (1969).

### 5. ILLUSTRATIONS

Some examples will be presented at the Conference to illustrate the procedures desfribed in section 3.

This work was done in the Research and Special Studies Division of the Central Statistical Office. Durbin's part of the work was done in the capacity of consultant to the Central Statistical Office.

43

### REFERENCES

- Anscombe, F.J. (1961). Examination of residuals. Proc. 4th Berkeley Symp.,  $\underline{1}$ , 1.
- Anscombe, F.J. and Tukey, J.W. (1963). The examination and analysis of residuals. Technometrics, 5, 141.
- Barnard, G.A. (1959). Control charts and stochastic processes. J. Roy. Statist. Soc. B, <u>21</u>, 239.
- Bartlett, M.S. (1951). An inverse matrix adjustment arising in discriminant analysis. Ann. Math. Statist., <u>22</u>, 107.
- Brown, R.L. and Durbin, J. (1969) Time-varying regression methods. (To be published.)
- Duncan, D.B. and Jones, R.H. (1966). Multiple regression with stationary errors. J. Amer. Statist. Assoc., <u>61</u>, 917.
- Durbin, J. (1960). Testing the hypothesis that a regression surface remains constant through time. Technical Report No. 5, September 1969, Applied Mathematics and Statistics Laboratories, Stanford University. (Unpublished.)
- Durbin, J. (1969). Tests for serial correlation in regression analysis based on the periodogram of least-squares residuals. (Submitted for publication.)
- Mehr, C.B. and McFadden, J.A. (1965). Certain properties of Gaussian processes and their first-passage times. J. Roy. Statist. Soc. B, <u>27</u>, 505.

Page, E.S. (1954). Continuous inspection schemes. Biometrika, <u>41</u>, 100. Tukey, J.W. (1962). The future of data analysis. Ann. Math. Statist., <u>33</u>, 1.

Walker, S.H. and Duncan, D.B. (1967). Estimation of the probability of an event as a function of several independent variables. Biometrika, 54, 167.

B&D 8

44

Woodward, R.H. and Goldsmith, P.L. (1964). Cumulative Sum Techniques, Monograph No. 3, I.C.I. Series on Mathematical and Statistical Techniques for Industry, Oliver and Boyd Ltd., Edinburgh.

B&D 9

45

. •

SOME EFFECTS OF ERRORS OF MEASUREMENT IN MULTIPLE REGRESSION 1)

by W.G. Cochran (USA) Department of Statistics. Harvard University

### 1. INTRODUCTION

In recent years there has been an increase in multiple regression studies on problems in which some of the independent variables represent quantities that are obviously difficult to measure, and are presumably measured with substantial errors. In the social sciences, for example, these variables may include measures of a person's skills at certain tasks or his attitudes and psychological characteristics, the data being obtained from a questionnaire, plus perhaps some kind of examination. Such studies raise the question: to what extent do these errors of measurement vitiate the uses to which the multiple regression is put? In examining this question, my results are less general than is desirable. The only tractable model is simpler than is needed for many applications. Even with this model, the effects of the errors are complex. I have, however, tried to indicate approximately what happens in situations representative of at least a substantial number of applications.

Frequent uses of multiple regression are: (1) to predict a variable y. The relevant quantity here is the residual variance  $\sigma_y^2(1-R^2)$ , where R is the population multiple correlation coefficient between y and the x's. (2) To study and try to interpret the values of the individual regression coefficients. The nature of the effects of errors of measurement on the values of the  $\beta_i$  has been indicated in a previous paper, Cochran, 1968. Consequently, this paper deals mainly with  $R^2$ , although the effects on the  $\beta_i$  will be discussed briefly in section 8.

### 2. MATHEMATICAL MODEL

Using capital letters to denote correctly measured values, we suppose that in the population the variate  $Y_u$  has a linear regression  $\alpha + \sum_{i} \beta_i X_{iu} + d_u$  on the <u>k</u> X's, where <u>d</u> is the random residual from the

<sup>1)</sup> This work was assisted by a Contract between the Office of Naval Research and the Department of Statistics, Harvard University.

regression. Owing to errors of measurement, the variates actually recorded for Y and for the  $\underline{i}$  X-variate are

(2.1) 
$$y_{u} = Y_{u} + a + e_{u}, \quad x_{iu} = X_{iu} + a_{i} + e_{iu},$$

where  $\underline{a}$  and the  $\underline{a}_i$  represent overall constant biases of measurement, while  $\underline{e}_u$  and the  $\underline{e}_u$  are fluctuating components which follow frequency distributions with means zero.

For this type of model, Lindley, 1947, gave the necessary and sufficient relations that must hold between the joint frequency function of the  $X_{iu}$  and that of the  $e_{iu}$  in order that the regression of  $y_{u}$  on the  $x_{iu}$  remain linear. In particular, if  $y_{u}$  and the  $X_{iu}$  follow a multivariate normal distribution, the  $e_{iu}$  must also follow a multivariate normal. This case is assumed here. Clearly, if  $Y_{u}$ ,  $e_{u}$ ,  $X_{iu}$  and the  $e_{iu}$  jointly follow a multivariate normal, it follows from relation (2.1) that  $y_{u}$  and the  $x_{iu}$  also follow a multivariate normal and hence that the regression of  $y_{u}$  on the  $x_{iu}$  is linear.

For the present it is assumed further that  $e_u$  is independent of  $Y_u$  and that any  $e_{iu}$  is independent of  $X_{iu}$  or any  $X_{ju}$  ( $j \neq i$ ) and of any other  $e_{ju}$ . These last assumptions are not essential to ensure linearity of the regression of  $y_u$  on the  $x_{iu}$ , and some remarks about the non-independent case will be made in section 7. For many applications in which both the  $X_{iu}$  and the  $e_{iu}$  appear non-normal, it would be desirable to bypass the normality assumptions, but I have no results for this situation.

The bias terms <u>a</u> and  $a_i$  in (2.1) affect the constant term in the regression of y on the  $x_i$ , but do not affect the multiple correlation coefficient between y and the  $x_i$ , and hence do not enter into the following sections on  $\mathbb{R}^2$ .

# 3. EFFECT ON R<sup>2</sup> WHEN X'S ARE INDEPENDENT

With  $\underline{k}$  X-variates, the following notation will be used for the relevant population parameters.

48

$$\sigma_{i}^{2} = \text{variance of the correct } X_{iu};$$

$$\epsilon^{2}, \epsilon_{i}^{2} = \text{variance of } e_{u}, e_{iu};$$

$$\rho_{ij} = \text{correlation coefficient between } X_{iu} \text{ and } X_{ju};$$

$$\delta_{i} = \text{correlation coefficient between } X_{iu} \text{ and } Y_{u}.$$

The symbol  $\delta_i$  is used instead of the more natural  $\rho_{iy}$  because this helps to avoid confusion between different kinds of correlation in later discussion. The sign attached to each  $X_{iu}$  is assumed chosen so that  $\delta_i \geq 0$ .

The value of  $R^2$ , the population squared multiple correlation between Y and the  $X_i$ , is completely determined by the  $\rho_{ij}$  and the  $\delta_i$ . Primes will be used to denote the corresponding correlations  $R'^2$ ,  $\rho'_{ij}$ ,  $\delta'_i$  between the observed y and the  $x_i$ . From the assumptions we have

(3.1) 
$$\rho'_{ij} = \frac{\operatorname{Cov}(X_i + e_i)(X_j + e_j)}{\sqrt{(\sigma_i^2 + \varepsilon_i^2)(\sigma_j^2 + \varepsilon_j^2)}} = \frac{\rho_{ij}\sigma_i\sigma_j}{\sqrt{(\sigma_i^2 + \varepsilon_i^2)(\sigma_j^2 + \varepsilon_j^2)}} = \frac{\rho_{ij}\sigma_i\sigma_j}{\sqrt{(\sigma_i^2 + \varepsilon_i^2)(\sigma_j^2 + \varepsilon_j^2)}}$$

$$(3.2) \qquad \delta'_{i} = \frac{\operatorname{Cov}(X_{i} + e_{i})(y + e)}{\sqrt{(\sigma_{i}^{2} + \varepsilon_{i}^{2})(\sigma_{Y}^{2} + \varepsilon^{2})}} = \frac{\delta_{i}\sigma_{i}\sigma_{Y}}{\sqrt{(\sigma_{i}^{2} + \varepsilon_{i}^{2})(\sigma_{Y}^{2} + \varepsilon^{2})}}$$

In psychometric writings the quantity  $\sigma_i^2/(\sigma_i^2+\epsilon_i^2)$  is often called the *coefficient of reliability* of  $x_i$ . We shall follow this terminology and define

.

$$g_i = \sigma_i^2 / (\sigma_i^2 + \epsilon_i^2) = \text{coefficient of reliability of } x_i$$
.

Similarly,

$$g_y = \sigma_y^2 / (\sigma_y^2 + \epsilon^2) = coefficient of reliability of y.$$

Hence, from (3.1) and (3.2),

(3.3) 
$$\rho'_{ij} = \rho_{ij} \sqrt{g_i g_j}, \qquad \delta'_i = \delta_i \sqrt{g_i g_y}$$

If the X's are mutually independent, it is well known that

(3.4) 
$$R^2 = \sum_{i=1}^{k} \delta_i^2$$
.

Since our assumptions guarantee that the x's are also independent in this case,

(3.5) 
$$R'^{2} = \sum_{i=1}^{k} \delta'^{2}_{i} = g_{y} \sum_{i=1}^{k} \delta^{2}_{i} g_{i}$$
.

Hence,

(3.6) 
$$R'^{2} = R^{2}g_{y}\sum_{i=1}^{k} \delta_{i}^{2}g_{i} / \sum_{i=1}^{k} \delta_{i}^{2} = R^{2}g_{y}\bar{g}_{w}$$

where  $\bar{g}_{w}$  is a weighted mean of the coefficients of reliability of the  $x_{i}$ .

Consider now the residual variance from the regression. With the correct measurements this is  $\sigma_Y^2(1-R^2)$ . With the fallible measurements it becomes

(3.7) 
$$\sigma_{y}^{2}(1-R'^{2}) = \sigma_{y}^{2} + \varepsilon^{2} - \sigma_{y}^{2} g_{y}^{-} g_{w}^{-} R^{2} = \sigma_{y}^{2}(1-R'^{2} g_{w}^{-}) + \varepsilon^{2}$$

since  $\sigma_{ygy}^2 = \sigma_y^2$ . Equation (3.7) contains the well-known result that under this model the effect of errors of measurement of y with variance  $\epsilon^2$  is simply to increase the residual variance by  $\epsilon^2$ , the variance of these errors.

As regards errors of measurement of the  $x_i$ , two points are worth noting in relation to applications. For a given reliability of measurement, i.e. a given value of  $\bar{g}_w$ , the deleterious effect on the residual variance increases as  $R^2$  increases, being greater when the prediction formula is very good than when it is mediocre. For example, suppose that  $\bar{g}_w = 0.5$ , representing a poor reliability in measurement of the  $x_i$ . If  $R^2 = 0.9$ ,

50

the residual variance is increased from  $0.1\sigma_Y^2$  to  $0.55\sigma_Y^2$ , over a five-fold increase. With  $R^2 = 0.4$ , the increase is only from  $0.6\sigma_Y^2$  to  $0.8\sigma_Y^2$ , a 33% jump.

Secondly, as would be expected, the quality of measurement of those  $X_i$  that are individually good predictors is much more important than that of poorer predictors. With k = 2,  $\delta_1 = 0.9$ ,  $\delta_2 = 0.3$ , we have  $R^2 = 0.90$ ,  $(1-R^2) = 0.1$ . If  $g_1 = 0.5$ ,  $g_2 = 1$ , this gives  $(1-R^{'2}) = 0.505$ , but with  $g_1 = 1$ ,  $g_2 = 0.5$ ,  $(1-R^{'2}) = 0.145$ , a much smaller increase.

### 4. EFFECT OF CORRELATION BETWEEN X's: TWO VARIATES

After working several numerical examples, my approach was to try to construct an approximation of the form  $R'^2 = R^2 g_y \bar{g}_w f$ , where <u>f</u> is a correction factor to allow for the effect of correlations among the X's, being equal to 1 when the X's are independent. But with numerous X variables, all intercorrelated, it soon appeared that no simple correction factor was likely to be generally applicable. However, we will continue to study the relation of  $R'^2$  to  $R^2 g_y \bar{g}_w$ . Further, since the effect of errors in y under this present model is always just to introduce the factor  $g_y$ , this factor will be omitted in what follows so as to concentrate attention on correlations among the X's.

With 2 X-variates having a correlation  $\rho$ , the values of  $R^2$  and  $R'^2$  work out as follows:

(4.1) 
$$R^{2} = (\delta_{1}^{2} + \delta_{2}^{2} - 2\rho \delta_{1} \delta_{2})/(1-\rho^{2}) ,$$

(4.2) 
$$R'^{2} = (g_{1}\delta_{1}^{2}+g_{2}\delta_{2}^{2} - 2g_{1}g_{2}\delta_{1}\delta_{2})/(1-g_{1}g_{2}\delta_{2}^{2}).$$

For given  $\delta_1$ ,  $\delta_2$ , the correlation  $\rho$  lies within the limits  $\delta_1 \delta_2 + \sqrt{(1-\delta_1^2)(1-\delta_2^2)}$ , otherwise R<sup>2</sup> would exceed 1. Within these limits,

(4.3) 
$$\mathbf{R'^{2}} = \mathbf{R}^{2} \frac{(\mathbf{g}_{1}\delta_{1}^{2} + \mathbf{g}_{2}\delta_{2}^{2} - 2\mathbf{g}_{1}\mathbf{g}_{2}\rho\delta_{1}\delta_{2})}{(\delta_{1}^{2} + \delta_{2}^{2} - 2\rho\delta_{1}\delta_{2})} \cdot \frac{(1-\rho^{2})}{(1-\mathbf{g}_{1}\mathbf{g}_{2}\rho^{2})}$$

51

Taking out  $\overline{g}_{W} = (g_1 \delta_1^2 + g_2 \delta_2^2) / (\delta_1^2 + \delta_2^2)$  as a factor, we may write

(4.4) 
$$R^{2} = R^{2} \overline{g}_{w}(A)(B)$$

where B is the term

(4.5) 
$$B = (1-\rho^2)/(1-\rho^2 g_1 g_2) .$$

For  $g_1g_2 < 1$ ,  $\rho \neq 0$ , this term is always < 1. For fixed  $g_1g_2$  it decreases monotonically towards 0 as  $\rho$  moves from 0 towards either +1 or -1.

Factor A takes the form

(4.6) 
$$A = \left(1 - \frac{2\rho\delta_{1}\delta_{2}}{\frac{\delta_{1}^{2} + \delta_{2}^{2}}{\frac{1}{g_{2}} + \frac{\delta_{2}^{2}}{g_{1}}}}\right) / \left(1 - \frac{2\rho\delta_{1}\delta_{2}}{\delta_{1}^{2} + \delta_{2}^{2}}\right)$$

For  $0 < g_1 g_2 < 1$ , it follows that A > 1 if  $\rho$  is positive while A < 1 if  $\rho$  is negative, provided that  $\delta_1, \delta_2$  are both > 0.

Hence, if  $\rho$  is negative the factor f = AB is always < 1, decreasing towards zero as  $\rho$  approaches -1. If  $\rho$  is positive the situation is not so clear, since A > 1 and B < 1. However, when  $\rho$  is small the factor A, which contains only linear terms in  $\rho$ , tends to dominate B which is quadratic in  $\rho$ . Thus when  $\rho$  is positive, f = AB increases and is greater than 1 for a time, but then decreases as  $\rho$  increases further, becoming less than 1 if  $\rho$  is high enough. The only exception is the case  $\delta_1 = \delta_2$ ,  $g_1 = g_2 = g$ : f then reduces to  $(1+\rho)/(1+g\rho)$ , which increases from f = 1 at  $\rho = 0$  to f = 2/(1+g) at  $\rho = 1$ . Incidentally, when  $\delta_1 = \delta_2$ , the range of  $\rho$  is from  $(-1 + 2\delta_1^2)$  to +1.

The size of the product  $g_1g_2$  is also relevant to f. For given  $\rho$ , both A and B tend to approach 1 as  $g_1g_2$  increases towards 1. Thus the formula  $R'^2 = R^2g_w$  is closer to the truth when  $g_1$  and  $g_2$  are high.

In a previous paper, Cochran, 1961, the **ef**fect of  $\rho$  on the value of  $R^2$  was studied in connection with applications to the discriminant

52

function. From (4.1) it is clear that negative values of  $\rho$  are helpful to prediction, since when  $\rho$  is negative,  $\mathbb{R}^2$  always exceeds the value  $(\delta_1^2 + \delta_2^2)$  that it would have if  $\rho$  were 0. With  $\rho$  positive,  $\mathbb{R}^2$  decreases at first but has a minimum at  $\rho = \delta_2/\delta_1$ , where  $\delta_2 < \delta_1$ , and thereafter increases. It does not reach  $(\delta_1^2 + \delta_2^2)$  until  $\rho = 2\delta_1\delta_2/(\delta_1^2 + \delta_2^2)$ , which is high if  $\delta_1$  and  $\delta_2$  are not too different. Thus, positive correlations are harmful to prediction unless  $\rho$  is high enough.

As an illustration, table 4.1 shows the values of  $R^2$  and f for  $\rho = -.0.5(0.1) + 0.9$ , for six examples. In the first three,  $\delta_1 = 0.6$ ,  $\delta_2 = 0.4$ , and in the second,  $\delta_1 = 0.7$ ,  $\delta_2 = 0.2$ . The three pairs  $g_1, g_2 = (0.9, 0.7)$ , (0.8, 0.6), and (0.7, 0.5) are given. The behavior of  $R^2$  and f as described above may be noted, as well as the increasing departure of f from 1 as the product  $g_1g_2$  decreases. The principal difference between the cases  $\delta_1 = 0.6$ ,  $\delta_2 = 0.4$  and  $\delta_1 = 0.7$ ,  $\delta_2 = 0.2$  is as follows. When  $\delta_1$  and  $\delta_2$  differ greatly and  $\rho$  is positive,  $R^2$  begins to increase and f to decrease for quite moderate values of  $\rho$  (around 0.3 for  $\delta_1 = .7$ ,  $\delta_2 = .2$ ), while when  $\delta_1$  and  $\delta_2$  are more nearly equal,  $R^2$  decreases and f increases until  $\rho$  is closer to 1. The turning value of f is a complicated expression, but is usually close to that of  $R^2$ .

The complementary sets  $g_1, g_2 = (0.7, 0.9)$ , (0.6,0.8), (0.5,0.7), not shown here, exhibit the same behavior with f lying a little nearer 1, except for high, positive  $\rho$  when f becomes less than 1.

	$\delta_1 = .6, \delta_2 = .4$					$\delta_1 = .7, \ \delta_2 = .2$				
	g <sub>i</sub> =	.9,.7	.8,.6	.7,.5		g <sub>i</sub> =	.9,.7	<b>.8,.</b> 6	.7,.5	
ρ	R <sup>2</sup>	f	f	f		$\mathbf{R}^{2}$	f	f,	f	
5	_ <b>x</b>	_	-	_		.893	0.84	0.78	0.74	
4	.848	0.87	0.82	0,78		.764	0.89	0.85	0.81	
3	.730	0.91	0.88	0.85		.675	0.93	0.90	0.88	
2	.642	0.95	0.92	0.90		.610	0.96	0,94	0.93	
1	.574	0,98	0.97	0.96		.564	0.98	0.98	0.97	
0	.520	1.00	1.00	1.00		.530	1.00	1.00	1.00	
.1	.477	1.02	1.03	1.04		.507	1.01	1.02	1,02	
.2	.442	1.04	1.06	1.07		.494	1.02	1.02	1.03	
.3	.413	1.06	1.08	1.10		.490	1.02	1.02	1.03	
.4	.390	1.07	1.10	1.12		.498	1.00	1.00	1.01	
.5	.373	1.08	1.11	1.14		.520	0.98	0.97	0.97	
.6	.362	1.08	1.11	1.14		.566	0.94	0,91	0.90	
.7	.361	1.07	1.09	1.12		.655	0.86	0.82	0.79	
.8	.378	1.03	1.03	1.06		.850	0.73	0.67	0.63	
.9	.463	0.86	0.85	0.85		_*	-	-	-	
	ē <sub>w</sub> =	.838	.738	.638		Ē <sub>w</sub> =	.885	.785	.685	

Table 4.1 Values of  $f = R'^2/R^2 \overline{g}$  for six example.

\* Impossible because  $R^2 > 1$ .

.

In these examples  $\bar{g}_{w}$  lies between 0.638 and 0.885. As regards the crude approximation  $R^{'2} \stackrel{\vee}{=} R^2 \bar{g}_{w}$ , in these examples this is correct to within  $\pm$  15% for  $\rho$  lying between -0.3 and +0.5, being much closer than this throughout most of table 4.1.

To summarize for two independent variates: when  $\rho$  is negative, f < 1 because the negative correlation  $\rho' = \rho \sqrt{g_1 g_2}$  is less helpful to R<sup>'2</sup> than the negative correlation  $\rho$  is to R<sup>2</sup>. When  $\rho$  is positive and small or modest, f exceeds 1, because the harmful positive correlation is decreased by the errors of measurement. If  $\rho$  becomes high enough, however, positive correlation becomes helpful and f drops below 1. For given  $\rho$ , f departs further from 1 as the product  $g_1g_2$  decreases.

With 3 X-variables denoted by the subscripts i, j, and k, the value of  $R^2$  may be expressed as

(4.7) 
$$R^{2} = \frac{\sum_{i} \delta_{i}^{2} (1-\rho_{jk}^{2}) - 2 \sum_{j>i} (\rho_{ij}-\rho_{ik}\rho_{jk})\delta_{i}\delta_{j}}{1 - \sum_{j>i} \rho_{ij}^{2} + 2\rho_{12}\rho_{13}\rho_{23}}$$

while  $R'^2$  has the corresponding value found by substituting  $\delta'_i = \delta_i \sqrt{g_i}$ ,  $\rho'_{ij} = \rho_{ij} \sqrt{g_i g_j}$ . These expressions are discouraging to the prospect of finding an approximation for f that would be valid over a wide range of values of the  $g_i$  and the  $\rho_{ij}$ . With regard to  $R^2$  itself, (4.6) suggests that with all  $\delta_i > 0$ , negative values of  $\rho_{ij}$  are likely to be helpful, since the only linear term in the  $\rho_{ij}$  is  $-2\rho_{ij}\delta_i\delta_j$  in the numerator.

Before proceeding further, we digress to consider the values of the  $\rho_{i,j}$  and the  $g_i$  likely to occur in practice.

# 5. SOME VALUES OF r<sub>i</sub> IN PRACTICAL APPLICATIONS

When the sign attached to each  $x_i$  is chosen so that  $\delta_i \geq 0$ , these decisions determine the sign attached to every  $\rho_{ij}$ . In studying the estimates  $r_{ij}$  of the  $\rho_{ij}$  found in 12 well-known examples of the discriminant function, Cochran, 1961, I noted that most of the  $r_{ij}$ 

are positive and modest in size, while those that are negative are usually small. The same situation appears to hold in many applications of multiple regression. Table 5.1 shows the distributions of the  $r_{ij}$  in (i) the discriminant function examples, (ii) the numerical examples of a multiple regression given in 12 standard statistical texts, (iii) a single large example--the prediction of verbal ability scores of 12th grade white students in the north of the U.S. from 20 variables representing data on the student, the quality of the school, and the student's home environment, Coleman et al., 1966.

#### Table 5.1

	Number of Cases				Number of Cases		
r ij	D.F.	Texts	Verbal	r ij	D.F.	Texts	Verbal
<5	1	1	0	0 to .1	15	5	58
5 to4	2	0	0	.1 to .2	22	8	41
4 to3	1	0	1	.2 to .3	25	7	<b>2</b> 5
3 to2	4	2	2	.3 to .4	18	9	7
2 to1	4	2	10	.4 to .5	6	7	3
1 to 0	9	5	36	.5 to .6	10	6	6
				.6 to .7	4	5	1
				.7 to .8	1	4	0
				> .8	0	3	0
Total	21	10	49	Total	101	54	141
r	-0.19	-0.17	-0.09	r	+0.30	+0.41	+0.16

Distributions of estimated correlations between x's

The percentages of r's that are positive are 83%, 84%, and 74% in the three sets. The negative r's average to between -0.2 and 0, the averages of the positive r's being a little higher. While some allowance is needed for the sampling errors of the  $r_{ij}$ , since our interest is in the unknown  $\rho_{ij}$ , my impression is that most of the degrees of freedom are large enough so that the effect of sampling errors on the average

56

r's should be small. In the discriminant function and text examples there may have been some selection towards more interesting examples, but this would probably affect the sizes of  $\delta_i$  rather than the  $\rho_{ii}$ .

In calculations for 3 or more x's, these results led me to concentrate on  $\rho_{ij}$  less than 0.5, and mainly on two cases: (1) all  $\rho_{ij}$  positive, (2) only a minority negative.

6. THE PROBLEM OF ESTIMATING RELIABILITY

With variables that are hard to measure, the problem of estimating the reliability of measurement is also formidable, and I have not come across any set of g values that might be regarded as representative. Direct estimation of g is possible only when it is feasible to measure both the correct value X and the fallible value x for a sample of items. This situation is likely to be confined mainly to applications in which (1) g is high and (2) the fallible measurement is enough cheaper or more convenient to make it preferable to the correct measurement. Occasionally, an opportunity to measure X may present itself even though X is not usually available. Thus, the reliability of appraiser's estimates of the values of homes may be estimated by finding the actual sellingprices for these homes that happen to have been sold recently: data of Kish and Lansing, 1954, indicate a g of around 0.83 in this situation.

When X cannot be measured, assume first that for the uth item the correct measurement  $X_u$  is constant (i.e. not varying with time). If two independent measurements  $x_{u1}$ ,  $x_{u2}$  of each item by the fallible instrument can be made, the quantity  $Cov(x_{u1}, x_{u2})$  estimates  $\sigma_X^2$ , so that  $\hat{g} = Cov(x_{u1}, x_{u2})/s_x^2$  is an estimate of g. This method is widely used in appraising the reliability of examinations, the two measurements being either random halves or alternative forms of an examination. Naturally, values of g over 0.9 are sought, though values between 0.7 and 0.9 may be considered acceptable if the skill in question is difficult to measure by examination. The assumption of independence is crucial in this approach. With a positive correlation between the errors  $e_{u1}$  and  $e_{u2}$ ,  $Cov(x_{u1}, x_{u2})$  overestimates  $\sigma_x^2$  so that g is overestimated.

57

Alternatively, the same measurement may be made on the specimens at two different times. Examples are examinations given a week apart, or questions repeated to a respondent on a later occasion. If these questions refer to memory of a definite past event, the errors of measurement may be smaller on the first than on the second occasion. Fortunately, still assuming independence and constant  $X_u$  for given  $\underline{u}$ , all three quantities  $\sigma_X^2$ ,  $\sigma_X^2$ , and  $\sigma_Z^2$  can be estimated, as can  $g = \sigma_X^2/\sigma_{X_1}^2$ , the reliability of the answers on the first occasion. With questions involving memory, however, positive correlation between errors is a constant danger, since the respondent may recall the same wrong answer on both occasions.

When the correct measurement varies with time, interpretation becomes more complex. For the uth item on the jth occasion, the simplest model is to write the correct value as  $X_u + t_{uj}$ , where  $X_u$  now represents an average value over time for the uth item and  $t_{uj}$ represents the fluctuation over time for the true value. The observed value on the jth occasions is  $x_{uj} = X_u + t_u + e_{uj}$ . For simplicity, assume  $t_{uj}$  and  $e_{uj}$  independent from item to item and from occasion to occasion. Then if the objective is to measure the correct value of X on a specific occasion, i.e. to measure  $X_u + t_{uj}$ , the reliability of our measuring process is

$$g = (\sigma_X^2 + \sigma_t^2) / \sigma_x^2 .$$

This quantity is estimated by  $\hat{g}$  if our data are a sample of two independent measurements of each specimen on the <u>j</u>th occasion. But if our sample consists of independent measurements on two different occasions,  $\hat{g} = \text{Cov}(x, x_{uj})/s_x^2$  estimates  $\sigma_X^2/\sigma_x^2$ , which can be a serious underestimate if  $\sigma_t^2$  is large. For instance, Guilford, 1959, from measurements one day apart on the same subjects, reports estimated g values of 0.22 for respiration period, 0.36 for white blood cell count, 0.65 for blood sugar content, and 0.74 for systolic blood pressure. Presumably, these relatively low values are in part caused by real day to day variation in the values of the items.

Further, when X varies through time, the relevant quantity for

WGC 12

58

prediction of y may not be the value of X on the first occasion, but some function of its levels over time, as for instance in the prediction of death rate from cigarette smoking history. This point may be put more generally. In difficult problems of measurement we may be attempting, through ignorance, to measure the wrong quantity. At its simplest, suppose that the relevant true measurement for the <u>u</u>th item in the population is  $X_u$ , which does not vary with time. The "true" value which we attempt to measure is  $X_u + a_u$ , where  $a_u$  is a random variable representing the extent to which we are measuring the wrong quantity. Our observed values for two independent fallible measurements are  $X_u + a_u + e_u$ . Hence,  $\hat{g}$  estimates  $(\sigma_X^2 + \sigma_a^2)/\sigma_X^2$ , whereas the relevant g is  $\sigma_X^2/\sigma_X^2$ .

For these reasons I am unable to name any narrow range of values of g which can represent practical experience in difficult measurement problems. My calculations have been done for the range  $g \ge 0.5$ : they should perhaps have been extended to lower g's. Ignorance of the values of the  $g_i$  for a specific application of interest is, of course, a considerable detriment to the use of any results of this paper. There is, however, increased interest in studying errors of measurement, as evidenced by the work of the U.S. Census Bureau, e.g. Hansen, Hurwitz, and Bershad, 1961, and later papers, by Kish's study, 1962, of interviewer variance, and by Mandel's study, 1959, of errors of measurement by different laboratories.

It should not be forgotten that g is a measure of precision of measurement relative to the true variation in the population. High values of g may be found with what appears quite sloppy and imprecise measurements, because the population is highly variable. A low value, such as  $\hat{g} = 0.41$  reported by Kinsey, 1948, for "Age at first knowledge of venereal disease", (by repeating the question on a later occasion) may in part reflect the fact that the correct ages have a small standard deviation.

7. EFFECT OF ERRORS WITH MORE THAN TWO X VARIATES

Returning to the relation between  $R^{2}$  and  $R^{2}$  with k X-variates

WGC 13

59

(k>2), the only case which will be discussed algebraically is that in which  $\rho_{ij} = \rho > 0$ ,  $g_i = g$ . In this case R<sup>2</sup> and R<sup>'2</sup> have the simple formulas

(7.1) 
$$R^{2} = \frac{\sum \delta_{i}^{2}}{1 + (k-1)\rho} \left[1 + \frac{k\rho\sum (\delta_{i} - \overline{\delta})^{2}}{(1-\rho)\sum \delta_{i}^{2}}\right]$$

(7.2) 
$$R'^{2} = \frac{g\sum \delta_{i}^{2}}{1 + (k-1)g\rho} \left[ 1 + \frac{gk\rho\sum (\delta_{i} - \overline{\delta})^{2}}{(1-g\rho)\sum \delta_{i}^{2}} \right] \cdot$$

If the  $\delta_i$  are approximately equal, i.e. the  $X_i$  are individually about equally good, the first terms in (7.1) and (7.2) dominate. Then we have

(7.3) 
$$f = \frac{R^2}{gR^2} \cong \frac{1 + (k-1)\rho}{1 + (k-1)g\rho}$$

For  $\rho$  positive, this f exceeds 1 and increase steadily as  $\rho$  goes from 0 to 1.

The case  $\delta_{i} = \delta$ ,  $\rho = 1$  is of some interest. In measuring a trait of a subject, such as aggressiveness, a common practice is to ask <u>k</u> questions, the answer to each being a measure of aggressiveness. If all the questions measured aggressiveness correctly, we would have  $\rho = 1$ ,  $\delta_{i} = \delta$ , making  $R^{2} = \delta^{2}$ , reflecting the fact that in this event any one question contains all the information. If the questions have reliability g and independent errors,  $\delta_{i} = \sqrt{g}\delta$ ,  $\rho' = g\rho$  and  $R'^{2} =$  $kg\delta^{2}/\{kg + (1-g)\}$  from (7.2). For example, with g = .6 and k = 5 questions,  $R'^{2} = 3\delta^{2}/3.4 = 0.88\delta^{2}$ , considerably better than the value  $R'^{2} = 0.6\delta^{2}$ that we would get by asking only one question. The argument here is the same as that used in the well-known correction for attenuation. If the errors of measurement in the different questions are positively correlated with one another though still uncorrelated with X, we do not do quite so well. With a correlation r between these errors,  $R'^{2}$  works out as  $kg\delta^{2}/\{kg + (1-g)(1+kr-r)\}$ . Thus with g = .6, k = 5, r = .5,  $R'^{2} = 0.71\delta^{2}$ , as against  $0.88\delta^{2}$ .

When the  $\delta_i$  vary in (7.1) and (7.2) the ratio of the second terms inside the brackets is  $g(1-\rho)/(1-g\rho)$ . This ratio is less than g, and therefore less than 1, and decreases as  $\rho$  increases. Thus this term slows down the increase in f. In applications the two terms in R<sup>2</sup> are often of the same order of magnitude, so that f shows only a small rise above 1 for positive and moderate values of  $\rho$ .

Table 7.1 shows the values of f for seven examples for 3, 5, and 10 x's, selected from those worked. In these examples, the reliability g is the same for all x's, f being given for g = .9(.1).5. For each example the number of x-variates k, and the values for the  $\delta_i$  and of the  $\rho$  are given. At the foot of the table are the values of  $R^2$  and of  $R_{ind}^{2ij} = \sum_{i}^{\delta_{i}^2}$ , the value that  $R^2$  would have if the X's were independent.

Table	7.	1
-------	----	---

3 5,.5,.4	5 .5,(.4) <sup>2</sup> .3,.2	10 .5,.4,(.3) <sup>2</sup> (.2) <sup>3</sup> ,(.1) <sup>3</sup>	3 .6,.5,.4	3 .6,.5,.4	5 .5,.4,.3	10 5 4 (3) <sup>2</sup>
5,.5,.4	.5,(.4) <sup>2</sup> .3,.2	$(.2)^3, (.1)^2$	.6,.5,.4	.6,.5,.4	.5,.4,.3	$54(3)^2$
					.2,.4	$(.2)^2, (.1)^3,$
.3	.3	.3	.2,2,	.3,.3,	.3(j≠5)	.3(j≠10)
			2	2	<b></b> 3(j=5)	<b>2(</b> j=10)
f	f	f	f	f	f	f
1.03	1.04	1.01	0.98	0.99	0.97	1.00
1.07	1.08	1.02	0.97	1.00	0.94	1.00
1.11	1.13	1.04	0.95	1.01	0.92	1.00
1.16	1.19	1.08	0.94	1.04	0.90	1.01
1.21	1.26	1.12	0.93	1.04	0.89	1.03
.497	.369	.390	.875	.630	.805	.511
	f 1.03 1.07 1.11 1.16 1.21	f       f         1.03       1.04         1.07       1.08         1.11       1.13         1.16       1.19         1.21       1.26         .497       .369	f       f       f         1.03       1.04       1.01         1.07       1.08       1.02         1.11       1.13       1.04         1.16       1.19       1.08         1.21       1.26       1.12         .497       .369       .390	$\begin{array}{c ccccccccccccccccccccccccccccccccccc$	$\begin{array}{c ccccccccccccccccccccccccccccccccccc$	$\begin{array}{c ccccccccccccccccccccccccccccccccccc$

## Values of f for seven selected examples

In general, the values of f in table 7.1 behave as would be expected from the results for k = 2 in section 4. In the three cases with  $\rho_{ij} = \frac{1}{10}$  +0.3, this correlation produces a marked decrease in R<sup>2</sup> as compared with R<sup>2</sup><sub>ind</sub>. Consequently, f rises above 1 because errors of measurement reduce the detrimental correlation to +0.3g.

When some  $\rho_{ij}$  are positive and some negative, we may in general regard the positive  $\rho_{ij}$  as harmful to prediction and the negative  $\rho_{ij}$  as helpful, though this is an oversimplification of what can happen with more than 2 x's. Two examples with substantial proportions of negative correlations have been included in table 7.1. The first example with k = 3 has 2 of the 3  $\rho_{ij}$  negative. In the example with k = 5, all correlations between  $x_5$  and the other x's are negative, making 4 negative  $\rho_{ij}$  out of 10. In both examples the net effect of the correlations is distinctly helpful,  $R^2$  exceeding  $R^2_{ind}$ . As would be anticipated, f < 1 in both examples. In the two remaining examples, with k = 3 and k = 10, one  $\rho_{ij}$  out of 3 and 9 out of 45 are negative. The net effect of the correlations is a decrease in  $R^2$  versus  $R^2_{ind}$ , though less marked than when all  $\rho_{ij}$  are positive. In both examples f rises only very slightly above 1.

When some  $\rho_{ij}$  are negative and the  $g_i$  are very unequal, f is more erratic. Its behavior still follows the lines indicated above. As an illustration, table 7.2 gives f for some unequal  $g_i$  for the example in table 7.1 with k = 5 and  $\rho_{i5}$  negative.

### Table 7.2

			-	I J	
<sup>g</sup> 1g4	=	1.0	.8	.6	,5
g <sub>5</sub>	=	.5	.6	.8	1.0
f	Ш	0.76	0.85	0.99	1.09
- g <sub>w</sub>	=	.886	.754	.646	.614

Values of f with unequal  $g_i = k=5$ ,  $\rho_{ij} = .3$ ,  $(j \neq 5)$ , = -.3(j=5)

In the first case in table 7.2,  $g_1 \dots g_4$  are without error, while  $g_5$  has reliability only 0.5. The harmful correlations  $\rho$  = +.3 are unaltered:

62

the helpful ones are reduced to  $-0.3\sqrt{.5} \cong -0.21$ . Consequently, f is well below 1 although  $\overline{g}_{w} = .886$  is quite high. The opposite case, with  $g_1 \dots g_4 = 0.5$ ,  $g_5 = 1$ , gives f = 1.09, the harmful correlations being reduced more than the helpful one. The two middle examples illustrate less extreme situations of the same type.

Since with electronic computers, multiple regression calculations having as many as 50 independent variates are becoming commoner, a summary statement could be made with greater confidence if examples of this size, with all  $\rho_{ij}$ ,  $\delta_i$  and  $g_i$  different, had been worked, and if more were known about the values of the  $g_i$  likely to occur in such problems and about the cruciality of the assumption of a multivariate normal. As a rough guide to the effects of errors of measurement, the relation  $R'^2 \stackrel{\sim}{=} R^2 g_y \bar{g}_w$  may serve, the value of  $R'^2$  being perhaps 10-20% higher than this if most correlations among the X's are positive and harmful, and perhaps 10-20% lower if we are lucky enough to have mainly helpful correlations. This rule assumes that the errors of measurement are independent of one another and of the correct X's.

If the  $e_i$  and  $e_j$  for two different X's have a correlation  $c_{ij}$  but are still independent of the true  $X_i$  or  $X_j$ , we have, in the notation of section 3,

$$\rho'_{\mathbf{i}\mathbf{j}} = \frac{\rho_{\mathbf{i}\mathbf{j}^{\sigma}\mathbf{i}^{\sigma}\mathbf{j}}^{\mathbf{j} + c}\mathbf{i}\mathbf{j}^{\varepsilon}\mathbf{i}^{\varepsilon}\mathbf{j}}{\sqrt{(\sigma_{\mathbf{i}}^{2} + \varepsilon_{\mathbf{i}}^{2})(\sigma_{\mathbf{j}}^{2} + \varepsilon_{\mathbf{j}}^{2})}} : \quad \dot{\delta}'_{\mathbf{i}} = \frac{\delta_{\mathbf{i}^{\sigma}\mathbf{i}^{\sigma}\mathbf{Y}}}{\sqrt{(\sigma_{\mathbf{i}}^{2} + \varepsilon_{\mathbf{i}}^{2})(\sigma_{\mathbf{Y}}^{2} + \varepsilon_{\mathbf{j}}^{2})}} \cdot$$

Thus  $\delta_i = \delta_i \sqrt{g_i g_y}$  as before, but if  $c_{ij} > 0$ ,  $\rho_{ij}$  now exceeds  $\rho_i \sqrt{g_i g_j}$ . With most  $\rho_{ij}$  positive and harmful, it looks as if the effect of a positive  $c_{ij}$  will be that f will lie closer to unity.

Suppose now that  $e_i$  and  $X_i$  are correlated, with  $Cov(e_iX_i) = c_i$ . To take the simplest case,  $e_i$  is assumed uncorrelated with any other  $X_j$  or  $e_j$ , though in applications it may happen that  $e_i$  is correlated with some other X 's. The variance of  $x_i$  now becomes  $(\sigma_i^2 + \varepsilon_i^2 + 2c_i)$ , so that  $g_i$  becomes  $\sigma_i^2/(\sigma_i^2 + \varepsilon_i^2 + 2c_i)$ . With the above assumptions, the equation  $\rho'_{ij} = \rho_{ij}\sqrt{g_ig_j}$  still holds. However,  $e_i$  becomes correlated with y through

63

its correlation with  $X_i$ . We have

(7.4) 
$$\operatorname{Cov}(\mathbf{y},\mathbf{x}_{i}) = \operatorname{Cov}(\mathbf{y},\mathbf{X}_{i}) + \operatorname{Cov}\left\{(\alpha + \sum_{j} \beta_{j} \mathbf{X}_{j} + d), (e_{i})\right\}$$
$$= \delta_{i} \sigma_{j} \sigma_{i} + \beta_{i} c_{i}.$$

Hence,

(7.5) 
$$\delta'_{i} = \rho_{yx_{i}} = \delta_{i}\sqrt{g_{i}} + \frac{\beta_{i}c_{i}}{\sigma_{y}\sigma_{x_{i}}} = \delta_{i}\sqrt{g_{i}} + \frac{\beta_{i}c_{i}\sqrt{g_{i}}}{\sigma_{y}\sigma_{y}} + \frac{\beta_{i}c$$

Equation (7.5) suggests that if  $\beta_i$  has the same sign as  $\delta_i$ , as should happen with the predominant regression coefficients, a positive correlation between  $e_i$  and  $X_i$  will increase  $\delta'_i$  and hence tend to increase  $R'^2$ . This might be expected since the  $e_i$ , as it were, are doing some of the work of the  $X_i$ . The most interesting case of a *negative* correlation is that studied by Berkson, 1950, and Box, 1961, in controlled experiments in which the fallible  $x_i$  are set at predetermined values, the errors of measurement  $e_i$  being therefore uncorrelated with the fallible  $x_i$ . Hence,  $Cov(e_i x_i) = Cov(e_{i'i} x_i) + \epsilon_i^2 = 0$ , making  $c_i = -\epsilon_i^2$ . In this case the effect on  $R'^2$  is almost easily seen by considering the values of the regression coefficients in the next section.

### 8. EFFECTS ON REGRESSION COEFFICIENTS

The assumptions and notation are the same as in section 2, except that for the moment we assume that  $e_i$  and  $X_i$  have a covariance  $c_i$ . The symbols  $\sigma_{ij}$ ,  $\sigma_{ij}$  denote  $Cov(X_iX_j)$  and  $Cov(x_iX_j)$ , where  $\sigma_{ii} = \sigma_i^2$ ,  $\sigma_{ii} = \sigma_i^2 + \varepsilon_i^2 + 2c_i$ . The assumption of multivariate normality guarantees that y has a linear regression both on the  $X_i$  and on the  $x_i$ . These relations provide two expressions for  $Cov(yx_i)$ .

(8.1) 
$$\operatorname{Cov}(yx_{i}) = \operatorname{Cov}\{(\alpha' + \sum_{j} \beta'_{j}x_{j} + d'), (x_{i})\} = \sum_{j} \beta'_{j}\sigma'_{ij},$$

(8.2) 
$$\operatorname{Cov}(y_{i}) = \operatorname{Cov}\{(\alpha + \sum_{j} \beta_{j} X_{j} + d), (X_{i} + e_{i})\} = \sum_{j} \beta_{j} \sigma_{ij} + \beta_{i} c_{i}$$

since the e are assumed uncorrelated with <u>d</u>. Hence the relations connecting the  $\beta_i$  and the  $\beta_i$  are

(8.3) 
$$\sum_{j=1}^{\sigma_{ij}\beta_{i}} = \sum_{j=1}^{\sigma_{ij}\beta_{j}} + \beta_{i}c_{i}.$$

. •

Since  $\sigma_{ii} = \sigma_{ii} - \epsilon_i^2 - 2c_i$ , these relations may be written

(8.4) 
$$\sum_{j}^{\sigma'} i j (\beta'_{j} - \beta_{j}) = -\beta_{j} (c_{i} + \varepsilon_{i}^{2}) .$$

Assuming  $\sigma'_{ij}$  non-singular, let its inverse be  $\sigma^{ij'}$ . Then

(8.5) 
$$\beta'_{i} = \beta_{i} - \sum_{j} \sigma^{ij'} \beta_{j} (c_{j} + \varepsilon_{j}^{2}) .$$

Thus the effect of errors of measurement is that  $\beta_i$  is a linear combination of  $\beta_i$  and of all the other  $\beta$ 's. The only case in which  $\beta_i \equiv \beta_i$  occurs when the values of all the  $x_j$  have been set at predetermined levels, making  $(c_j + \varepsilon_j^2) = 0$  for all j (the Berkson case).

With a single x-variate,  $\sigma^{11'} = 1/(\sigma^2 + \epsilon^2 + 2c)$  so that

(8.6) 
$$\beta' = \beta(\sigma^2 + c) / (\sigma^2 + c^2 + 2c)$$

Since with one x-variate  $R^2 \sigma_y^2 = \beta^2 \sigma^2$  and  $R'^2 \sigma_y^2 = \beta'^2 (\sigma^2 + \epsilon^2 + 2c)$ , this gives

(8.7) 
$$\frac{\frac{R'^2}{R}}{R^2} = \frac{(\sigma^2 + c)^2}{\sigma^2(\sigma^2 + \epsilon^2 + 2c)} .$$

For c positive, this ratio increases steadily as suggested in section 7, reaching the value 1 if X and e have correlation 1, making  $c = \sigma \epsilon$ . In the Berkson case, with  $c = -\epsilon^2$ ,

(8.8) 
$$\frac{R'^2}{R^2} = 1 - \frac{\varepsilon^2}{\sigma^2}$$
,

this being a reminder that although  $\beta$  is unchanged, the residual variance is increased by the errors in x. In the Berkson case with <u>k</u> variates, it is easily shown that

(8.9) 
$$\frac{\mathbf{R}^{2}}{\mathbf{R}^{2}} = 1 - \frac{\sum_{i=1}^{n} \beta_{i}^{2} \varepsilon^{2}}{\sum_{i=1}^{n} \beta_{i}^{2} \varepsilon^{2}}$$

We now assume the  $e_i$  and  $X_i$  uncorrelated, and briefly consider the  $\beta_i$ . From (8.5) it is evident that the effects on a specific  $\beta_i$  are complicated, depending on the signs and sizes of the other  $\beta_j$  and on the terms in the inverse matrix. As an approximation to applications in which most correlations among the X's are positive and modest, the following is the expression for  $\beta_i$  when  $\rho_{ij} = \rho$ ,  $g_j = g$ :

(8.10) 
$$\beta_{i} = \frac{g(1-\rho)\beta_{i}}{1-g\rho} + \frac{g(1-g)\rho(\sum_{i}\beta_{i})}{(1-g\rho)[1+(k-1)g\rho]}$$

The first term, which predominates when g is high, amounts to a reduction of  $\beta_i$  to a value somewhat less than  $g\beta_i$ . The second term is a common contribution to all the  $\beta_i$ , and is positive if  $(\sum \beta_i)$  is positive. In the examples that I have worked, the net effect is to make  $\beta_i'/\beta_i$  slightly greater than g for the larger  $\beta_i$  and substantially greater than g for the smaller  $\beta_i$ . The differences between the  $\beta_i$  are smaller than these between the  $\beta_i$  so that it becomes more difficult to distinguish the important from the unimportant regression coefficients.

A further consequence of the general relations (8.4), (8.5) is that if only one  $x_i$ , say  $x_1$ , is subject to error,

(8.11) 
$$\beta_{1}' = \frac{\beta_{1}}{1+\epsilon_{1}^{2}\sigma^{11}} : \beta_{i}' = \beta_{i} - \frac{\epsilon_{1}^{2}\sigma^{11}\beta_{1}}{(1+\epsilon_{1}^{2}\sigma^{11})} \quad (i > 1)$$

WGC 20

66

Since  $\sigma_1^{11} < 1/\sigma_1^2$  unless  $X_1$  is uncorrelated with any other  $X_i$ , it follows that  $\beta_1 < g_1\beta_1$  in this case. Also, every other  $\beta_i$  that is correlated with  $\beta_1$  is affected by errors in  $x_1$ .

In examples worked with unequal g's, the  $\beta'_i$  for those  $x_i$  having low  $g_i$  are very substantially reduced, while some  $\beta'_j$  with higher  $g_j$  may exceed  $\beta_j$  because of the contributions from  $\epsilon_i^2$  in (8.5). Consequently in this situation, in ignorance of the  $g_i$ , interpretations based on the relative sizes of the  $\beta'_j$  can become highly misleading.

In a more positive vein, equations (8.4) and (8.5) would also enable us to estimate the  $\beta_i$  from the  $\beta'_i$ , if we had good estimates of the  $g_i$ and if the model could be assumed to apply.

### SUMMARY

Multiple regression studies in which some or all of the variables are difficult to measure, and therefore presumably are measured with substantial errors, are increasingly common, particularly in the social sciences. This paper attempts to discuss the effects of such errors of measurement on the utility of multiple regression, both when the objective is prediction and when it is interpretation of the regression coefficients. Several different mathematical models are possible, since there may be correlations between the error of measurement of a variable and the true value of that variable and also between the errors for different variables.

A multivariate normal distribution of the correct Y's and the correct  $X_i$  is assumed. The errors of measurement are assumed normal with variances  $\varepsilon_i^2$  and at first independent of the correct values and of each other. Formulas available in the simplest cases and worked numerical examples indicate that the formula  $R'^2 = R^2 g_y \bar{g}_w$  approximates the relation between the squared multiple correlation coefficients  $R'^2$  and  $R^2$  in the presence and absence of errors. Here,  $g_y = \sigma_Y^2/(\sigma_Y^2 + \varepsilon^2)$  is the coefficient of reliability of y, while  $\bar{g}_w = \sum_i \delta_i^2 g_i/\sum_i \delta_i^2$  is a weighted mean of the coefficients of reliability of the  $x_i$ . The formula  $R'^2 = R^2 g_y \bar{g}_w$  slightly underestimates  $R'^2$  when the correlations among the  $X_i$  are positive, and may slightly overestimate  $R'^2$  when a minority of these correlations are negative, but appears correct to within + 20% in most cases. The effects of correlation

67

between the  ${\bf e}_i$  and  ${\bf e}_j$  for different x's and between  ${\bf e}_i$  and X are also indicated.

With errors of measurement in the  $x_i$ , any regression coefficient  $\beta_i$  becomes in general a linear function of all the  $\beta$ 's in the regression equation. Interpretation of the sizes of these  $\beta_i$  may become highly misleading. The problem of estimating the  $g_i$  in practice is also discussed.

WGC 22

.

## REFERENCES

Berkson, J., 1950. Are there two regressions? <u>Jour.Amer.Stat.Assoc</u>., 45, 164-180.

٠...

Box, G.E.P., 1961. The effects of errors in the factor levels and experimental design. <u>Bull.Int.Stat.Inst.</u>, 38, 3, 339-355.

Cochran, W.G., 1968. Errors of measurement in statistics. Technometrics, Vol. 10, (in press).

Cochran, W.G., 1961. On the performance of the linear discriminant function. Bull.Int.Inst.Stat., 39, 2, 435-447.

Coleman, J.B. et al., 1966. Equality of educational opportunity. U.S. Government Printing Office, Washington, D.C.

Guilford, J.P., 1959. Personality. McGraw Hill, New York.

- Hansen, M.H., Hurwitz, W.N., and Bershad, M., 1961. Measurement errors in censuses and surveys. <u>Bull.Int.Stat.Inst.</u>, 38, 2, 359-374.
- Kinsey, A.C., Pomeroy, W.B., and Martin, C.E., 1948. Sexual behavior in the human male. W.B. Saunders, Philadelphia.
- Kish, L. and Lansing, J.B., 1954. Response errors in estimating the value of homes. Jour.Amer.Stat.Assoc., 49, 520-538.
- Kish, L., 1962. Studies of interviewer variance for attitudinal variables. Jour.Amer.Stat.Assoc., 57, 92-115.
- Lindley, D.V., 1947. Regression lines and the linear functional relationship. <u>Jour.Roy.Stat.Soc.</u> B, 9, 218-224.

Mandel, J., 1959. The measuring process. Technometrics, 1, 251-267.

WGC 23

.

## SOME RECENT DEVELOPMENTS IN THE THEORY OF MARKOV CHAINS

#### by J.F.C. Kingman (England)

University of Sussex and Stanford University

For many years the analysis of stochastic processes involving some degree of Markovian behaviour has depended, implicitly or explicitly, on the discovery of "regeneration points" for the process. This concept was introduced by Palm, and a systematic account has been given by Smith [8]. Roughly speaking, a regeneration point for a process  $X_t$  is a random time  $\tau$  such that the process  $X_{\tau+t}$  (t > 0) is independent of  $X_u$  (u <  $\tau$ ) and has the same stochastic structure as  $X_t$  (t > 0). Thus at time  $\tau$  the process is "regenerated", and its random evolution from  $\tau$  follows the same laws as governed the process from t = 0.

If a regeneration point exists, it is not difficult to see that there is a whole sequence  $\tau_1, \tau_2, \cdots$  of these points, forming a renewal sequence in the sense that the positive random variables  $\tau_1, \tau_2 - \tau_1, \tau_3 - \tau_2, \cdots$ are independent and identically distributed. The process  $X_t$  can then be split up into independent segments  $X_t (\tau_{n-1} \leq t < \tau_n)$  which can be examined separately. In particular, many problems can be reduced to questions about the renewal theory of the sequence  $\{\tau_n\}$ , (cf. [8], [9]).

For some processes, however, regeneration points exist in much greater profusion than appears in the rather severe theory just mentioned, and to ignore all but a single sequence is to sacrifice a good deal of information. For example, in a queueing system fed by a Poisson arrival stream, every point of time at which the queue is empty is a regeneration point of the process, and the set of such instants is not a sequence, but a collection of intervals. More generally, if  $X_t$  is a Markov process with initial state  $X_0 = x$ , then any time with  $X_t = x$  regenerates the process.

In recent years a theory has been developed which generalizes classical renewal theory by allowing the set of regeneration points to be more substantial than a discrete sequence. This theory has application to the theory of continuous-time Markov chains (as developed for example in [2]), and in particular to the very difficult problem of characterizing Markov transition probabilities. The details may be found in [4] and [5] (and in other references cited in the latter paper). Let  $Z_t$  (t > 0) be a stochastic process taking only the values 0 and 1. Suppose that, for any T for which the event  $\{Z_T = 1\}$  has positive probability, the processes  $Z_t$ (t < T) and  $Z_t$  (t > T) are independent conditionally on  $\{Z_T = 1\}$ , and that moreover the conditional distributions of the latter process  $Z_{t+u}$  (u > 0) are the same as the distributions of the original process  $Z_u$  (u > 0). Then the process Z is said to define a *regenerative phenomenon*. The phenomenon is said to *occur* at time t if  $Z_t = 1$ . A typical example of a regenerative phenomenon is provided by a Markov process  $X_t$  with initial state x, by means of the formula

$$Z_{t} = f(X_{t}), \qquad (1)$$

where

f(x) = 1, $f(\xi) = 0, \quad (\xi \neq x).$ 

Denote the probability of occurrence at time t by

$$p(t) = P[Z_{+} = 1].$$
 (2)

The regenerative condition on Z implies that, whenever 0 < t\_1 < t\_2 <  $\ldots$  < ... < t\_k, then

$$P\{Z_{t_1} = Z_{t_2} = \dots = Z_{t_k} = 1\} = p(t_1)p(t_2-t_1)\dots p(t_k-t_{k-1}). (3)$$

It follows that the function p, called the p-function of the phenomenon, determines the finite-dimensional distributions of the process Z. The first problem is therefore to determine which functions p can arise in this way. This accomplished, a second problem is to describe the behaviour of the process Z in terms of properties of the p-function.

The first fact to notice about p-functions is that they have to satisfy certain functional inequalities. For instance, the probabilities

$$P\{Z_{s} = 0, Z_{s+t} = 1\} = P\{Z_{s+t} = 1\} - P\{Z_{s} = 1, Z_{s+t} = 1\}$$
  
= p(s+t) - p(s) p(t)

JK 2

and

$$P\{Z_{s} = 0, Z_{s+t} = 0\} = P\{Z_{s} = 0\} - P\{Z_{s} = 0, Z_{s+t} = 1\}$$
$$= 1 - p(s) - p(s+t) + p(s) p(t)$$

must be non-negative, so that p necessarily satisfies

$$p(s) p(t) \leq p(s+t) \leq 1 + p(s) p(t) - p(s).$$
 (4)

More complicated inequalities may similarly be derived from the fact that

$$\mathbb{P}\{\mathbb{Z}_{t_1} = \mathbb{Z}_{t_2} = \dots = \mathbb{Z}_{t_k} = 0\} \ge 0$$

and

$$P\{Z_{t_1} = Z_{t_2} = \dots = Z_{t_{k-1}} = 0, Z_{t_k} = 1\} \ge 0$$

for every k. Conversely, the Daniell-Kolmogorov theorem can be used to prove that this infinite family of functional inequalities is sufficient, as well as necessary, for a function p to be a p-function. The study of p-functions is therefore the study of the consequence of these inequalities.

Of particular importance among p-functions are those which satisfy

$$\lim_{t \to 0} p(t) = 1;$$
 (5)

these are called *standard*. Their significance has recently been stressed by the proof [7] that any measurable p-function is either

(i) of the form ap(t), where  $0 < a \leq 1$  and p is a standard p-function, or

(ii) equal to zero almost everywhere.

It is an easy consequence of (4) that a standard p-function is (uniformly) continuous on t  $\geq 0$ , and the corresponding process Z is continuous in probability.

The fundamental theorem in the theory of regenerative phenomena is an integral representation formula for the Laplace transform

JK 3

$$\mathbf{r}(\theta) = \int_0^\infty \mathbf{p}(t) e^{-\theta t} dt$$

of the standard p-function p. This can always be expressed in the form

$$\mathbf{r}(\boldsymbol{\theta}) = \left\{\boldsymbol{\theta} + \int (\mathbf{1} - \mathbf{e}^{-\boldsymbol{\theta}}\mathbf{x}) \boldsymbol{\mu}(\mathbf{d}\mathbf{x})\right\}^{-1}.$$
 (6)

where  $\mu$  is a positive measure on the interval  $(0,\infty]$ , uniquely determined by p. Conversely, if  $\mu$  is any positive measure on this interval with

$$\int (1-e^{-x}) \mu(dx) < \infty, \qquad (7)$$

then there is a unique continuous function with Laplace transform given by (6), and this is a standard p-function. The formula (6) therefore sets up a one-to-one correspondence between the standard p-functions and the positive measures satisfying (7).

Because  $(1-e^{-x})$  is small near x = 0, condition (7) does not necessarily imply that  $\mu$  is a finite measure. But if it has finite total mass q,  $\mu$  has a simple interpretation. In this case it can be written  $\mu = q\pi$ , where  $\pi$  is the probability measure of a positive, but possibly infinite random variable. Then it can be shown that  $Z_t$  is a step function, constant on intervals whose lengths are independent random variables. The lengths of intervals on which  $Z_t = 1$  have the negative exponential distribution with density  $qe^{-qt}$ , while the lengths of intervals on which  $Z_t = 0$  have distribution  $\pi$ . Following Bartlett [1], we say that  $\mu$  is a multiple of the recurrence time distribution of the phenomenon.

When  $\mu$  has infinite total mass, the structure of Z is much more complicated to describe. Although such phenomena are in a sense pathological, they do occur in models of practical importance, for instance in the theory of dams.

Using (6), a number of properties of standard p-functions can be established:

(a) p(t) is strictly positive,

(b) p(t) has finite right and left derivatives  $D^{\dagger}p(t)$  and  $\overline{D^{\phantom{\dagger}}p(t)}$  at every positive value of t, and  $\overline{D^{\dagger}p(t)} - \overline{D^{\phantom{\dagger}}p(t)}$  is equal to the atom (if any) of  $\mu$  at t, (in particular p is continuously differentiable in JK 4

t > 0 if and only if  $\mu$  has no atoms in (0, $\infty$ )).

. •

(c) the derivative  $D^{^+}p(0)$  exists, but may be  $-\infty,$  and  $-D^{^+}p(0)$  is the total mass of  $\mu,$ 

(d) p(t) tends to a limit  $p(\infty)$  as  $t \to \infty$ , and

$$p(\infty) = \{1 + \int x\mu(dx)\}^{-1}$$
. (8)

The form of equation (6) may suggest to some readers a connection with the theory of processes with non-negative independent increments, and such a connection does indeed exist. Writing

$$\Phi(\theta) = \int (1 - e^{-\theta x}) \mu(dx),$$

it is known that there exists such a process  ${\tt Y}_{\tt t}$  with

$$E(e^{-\theta Y}t) = e^{-t\phi(\theta)}.$$
 (9)

In terms of Y, define a process Z taking values 0 and 1 by

$$Z_{+} = 1 < => s + Y_{c} = t$$
 for some s

Then it has been shown by Kendall (in an as yet unpublished study of the sample functions of regenerative phenomena) that Z defines a regenerative phenomenon whose p-function satisfies

$$\int_{0}^{\infty} p(t) e^{-\theta t} dt = \left\{ \theta + \phi(\theta) \right\}^{-1}.$$
 (10)

The most important example of a regenerative phenomenon arises as in (1), where  $X_t$  is a Markov chain in the sense of Chung [2], a Markov process taking values in a countable set S. For any i  $\in$  S, we may take  $X_0 = i$ , and the phenomenon then occurs at time t if and only if the chain is in its initial state i. If the transition probabilities are written

$$p_{ij}(t) = P\{X_t = j | X_0 = i\}$$
 (11)

then the p-function of the phenomenon is just the diagonal transition probability  $p_{ii}(t)$ . It is usual to assume that

$$p_{ij}(t) \rightarrow \delta_{ij}$$
 (t  $\rightarrow$  0), JK 5

in which case p<sub>ii</sub>(t) is a standard p-function, to which all the general results about such functions may be at once applied.

This observation yields nearly all the known properties of the function  $p_{ii}(t)$  (as given for example in [2]). But there is one surprising exception, for it is known that every function of the form  $p_{ii}(t)$  is continuously differentiable in t > 0. In view of the remark made under (b) above, this is equivalent to the statement that the corresponding measure  $\mu$  has no atoms (except perhaps at  $\infty$ ). This can even be strengthened to show that  $\mu$  has a density in  $(0,\infty)$ .

It therefore follows that not every standard p-function is of the form  $p_{ii}$  in some Markov chain. On the other hand, it is not difficult to show that every standard p-function is the limit of a sequence of such functions  $p_{ii}$ . It is therefore a delicate, and at present unsolved, problem to determine which standard p-functions can arise from Markov chains. All that can be presented here is an account of the few known partial results.

Because (6) sets up a one-to-one correspondence, the problem of characterising the functions  $p_{ii}$  is equivalent to that of characterising the corresponding measures  $\mu$ . Call  $\mu$  a Markov measure if the standard p-function corresponding to it in (6) can be expressed in the form  $p_{ii}$  for some Markov chain. Then the following facts are known [6]:

(A) the value of the atom at  $\infty$  is irrelevant to deciding whether or not a measure is a Markov measure,

(B) any multiple of a Markov measure is a Markov measure,

(C) any finite or countable sum (subject to (7)) of Markov measures is a Markov measure,

(D) any Markov measure of infinite total mass admits a decomposition as a countable sum of totally finite Markov measures,

(E) the convolution of two totally finite Markov measures is a Markov measure,

(F) any Markov measure has a density on  $(0,\infty)$  which is lower-semicontinuous and strictly positive (unless it is identically zero),

(G) any measure on  $(0,\infty)$  satisfying (7) and having a density of the form

JK 6

$$\sum_{m=0}^{\infty} \sum_{n=0}^{\infty} a_{mn} t^{m} e^{-nt}, \qquad (12)$$

where  $a_{mn} \ge 0$ , is a Markov measure.

The difficulty of the characterization problem arises from the tension between the continuous time parameter and the essentially discontinuous nature of the stochastic process. If more general Markov processes are considered, the difficulty disappears. Indeed, let Z be any standard regenerative phenomenon, and let

$$X_{+} = t - \sup \{ u \le t; Z_{\mu} = 1 \}$$
 (13)

be the time elapsed at t since the last occurrence of the phenomenon Then X is a Markov process, and the phenomenon defined by (1) with x = 0 is a trivial modification of Z with the same p-function.

It should be remarked, however, that for this process the phenomenon defined by (1) with any non-zero value of x is not standard. If p arises from a Markov process in which all states determine standard phenomena, then [6, III] similar restrictions on p apply as in the Markov chain case.

The direct application of the theory of p-functions to Markov chains involves the functions  $p_{ij}(t)$  only for i = j, but the theory can be modified to deal also with the non-diagonal case  $i \neq j$ . The appropriate concept is that of a *quasi-Markov chain*, which is a stochastic process taking values 0, 1, 2, ..., N such that each of the states 1, 2, ..., N, though not the anomalous state 0, has the Markov property of regenerating the process. An example of such a process can be obtained from a Markov process X by taking distinct states  $x_1, x_2, ..., x_N$  and setting

$$Z_{t} = f(X_{t})$$
(14)

where

$$f(x_{r}) = r \qquad (r = 1, 2, ..., N)$$
  
$$f(\xi) = 0 \qquad (\xi \notin \{x_{1}, x_{2}, ..., x_{N}\}).$$

JK 7

The analogue of the p-function is a matrix-valued function, whose Laplace transform has a representation similar to (6). In particular, in the countable case, a typical p-matrix with N = 2 is

$$\begin{pmatrix} p_{ii}(t) & p_{ij}(t) \\ & & \\ p_{ji}(t) & p_{jj}(t) \end{pmatrix} .$$
(15)

The detailed consequences of the theory of quasi Markov chains for the characterization problem will be found in [6]; it suffices here to quote the main result. A function f(t) can be expressed in the form  $p_{ij}(t)$ , where i and j are distinct states of some Markov chain, if and only if f is expressible as a convolution

$$f = p_1 * d\mu * p_2, \qquad (16)$$

or more explicitly

$$f(t) = \int_0^t \int_0^{t-u} p_1(t-u-v) \mu(du) p_2(v) dv,$$

where

(i)  $p_1$  and  $p_2$  are diagonal Markov transition functions,

(ii)  $\mu$  is a totally finite measure on  $[0,\infty)$  which, apart from a possible atom at 0, is a Markov measure, and

(iii) 
$$\mu[0,\infty) \int_0^\infty p_1(t) dt \leq 1$$
.

It follows therefore that, if the class of Markov measures can once be determined, the characterization problem for the transition functions of Markov chains, both diagonal and non-diagonal, will have been solved, and the known theorems about the functions  $p_{ij}$  will fall into their natural perspective. This will not, however, dispose of all the outstanding problems even in the analytical part of Markov chain theory (leaving aside, that is, problems about sample function behavior). For example, Kendall and Speakman [3] have studied the function

$$g(t) = \inf_{i} p_{ii}(t), \qquad (17)$$

$$JK 8$$

and a systematic theory of such g-functions is urgently needed. Again, recent work of Williams (as yet unpublished) has added new interest to the problem of giving necessary and sufficient conditions for a set of numbers  $q_{ji}$  (i,  $j \in S$ ) to be expressible as the derivatives at the origin

$$q_{ij} = p'_{ij}(0)$$
 (18)

of the transition functions of some Markov chain. Indeed, it would seem that the deep problems of the analytical theory of Markov chains are to characterize the various functions and matrices arising, of which the problem described in this brief survey is the simplest, if not the least demanding.

## REFERENCES

- [1] M.S. Bartlett, Recurrence and first passage times, Proc. Camb. Phil. Soc. 49 (1953) 263-275.
- [2] K.L. Chung, Markov chains with stationary transition probabilities, Springer 1960.
- [3] D.G. Kendall and J.M.O. Speakman, On Markov groups, Proc. 5th Berkeley Symposium, II 2, 165-186.
- [4] J.F.C. Kingman, The stochastic theory of regenerative events, Z. Wahrscheinlichkeitstheorie <u>2</u> (1964) 180-224.
- [5] J.F.C. Kingman, An approach to the study of Markov processes, J. Roy. Stat. Soc. B 28 (1966) 417-447.
- [6] J.F.C. Kingman, Markov transition probabilities I IV, Z. Wahrscheinlichkeitstheorie <u>7</u> (1967) 248-270, <u>9</u> (1967) 1-9, and to appear.
- [7] J.F.C. Kingman, Measurable p-functions, Z. Wahrscheinlichkeitstheorie (to appear).
- [8] W.L. Smith, Regenerative stochastic processes, Proc. Roy. Soc. A <u>232</u> (1955) 6-31.
- [9] W.L. Smith, Renewal theory and its ramifications, J. Roy. Stat. Soc. B <u>20</u> (1958) 243-302.

JK 9

## ON SEQUENTIAL SEARCH

by B. Eichhorn (Israel) Dept. of Statistics, Tel-Aviv University

THE PROBLEM IN GENERAL.

We face an unknown function f, belonging to a family of real functions  $\mathcal{F}$ , say on the interval I= [0,1], with some given properties. For instance, the family of all monotone nonincreasing functions with exactly one zero on I.

We want to estimate a point (or points) of I where f assumes values of particular interest (for instance, the value zero or its maximum) by using some specified estimate for which a loss function giving "the loss due to estimation" L (estimate, f)  $\geq 0$  is defined. To help us make this estimate we are allowed to observe the values of the function f at points of the domain I which we can choose sequentially.

The problem varies according to the specific assumptions about the class  $\mathcal{F}$  and the way f is obtained from it, by the kind of estimate to be used (point or interval), the requirements imposed on it and the loss function attached, by the kind of estimating procedure allowed and kind of optimization sought, e.g. a minimax procedure.

We distinguish between two kinds of sequential estimating procedures:

(a) The fixed sample size sequential procedure, or n-observation sequential procedure  $T_n$ , where the fixed number of observations is n. The class of all  $T_n$  admissible with respect to the particular problem that is discussed will be denoted by  $\mathcal{T}_n$ . We shall call it briefly an n-seq. procedure. Here we are allowed exactly n observations whose places we can choose sequentially, making the place of the ith observation a function of the places and values of the first i-1 observations.

(b) The "true" sequential procedure Te  $\mathcal{T}_{m{r}}$  the class of all such admissible procedures, which we shall call just sequential procedures.

Here we do not fix ahead of time the number of observations. Instead we use a stopping rule  $\delta \simeq \{\delta_1, \delta_2, \dots\}$  which at each stage i tells us wether to stop or to continue and take the (i+1)st observation, or in a somewhat more general way tells us to stop with a certain probability  $\delta_i$ , each  $\delta_i$  being a function of the first i observations.

In this case, where the number of observations is not fixed, we assume a cost of ob**s**ervation in addition to the loss due to estimation. Usually we shall assume a constant cost c > 0 per observation. In this case we are concerned with the "total" cost which is defined as the sum of the cost of observation and the loss due to estimation

$$R(T,f) \simeq L(T,f) + cn(T,f)$$
.

R,L and n are functions of the procedure T and the function f. If T and f do not determine n and L completely but determine their distributions we define R(T,f) = E(L(T,f)) + cE(n(T,f)). We call a procedure nonrandomized if the places of observation are completely determined by the procedure T and the function f, each observation being a function of the information obtained so far. This is in contrast to randomized procedures where our next observation could be chosen randomly according to some probability distribution which is determined by the previous observations.

Let us state all this precisely.

DEFINITION, A nonrandomized n-seq. estimating procedure  $T_n$  is given by (1)  $T_n = [x_1, g_2, \dots, g_n, l]$ where  $x_1 \in I$  is the first place of observation, the other places of observation  $x_k$  are given by

$$g_k : I^{k-2} \times R^{k-1} \longrightarrow I, \qquad 2 \leq k \leq n,$$

which are functions of the former  $x_i$ 's and  $f(x_i)$ 's;  $\ell$  is the estimate, which is a function of  $x_2, \ldots, x_n$  and  $f(x_1), \ldots, f(x_n)$ , and its range depends on the particular problem and the kind of estimate we use. For instance if we use an interval estimate,  $\ell$  will take values [s,t] with s,t  $\epsilon$  I and  $s \leq t$ .

DEFINITION. A nonrandomized sequential procedure T is given by

(2) 
$$\mathbf{T} = \{\delta_0, \ell_0, \mathbf{x}_1, \delta_1, \ell_1, \mathbf{g}_2, \delta_2, \ell_2, \dots\}$$

where  $x_1$  and  $g_k$ ,  $(k \ge 2)$  are as in (1).  $\delta_k$ , k = 0, 1, ... is the probability of stopping with k observations given the values of the first k observations. So having reached stage k we stop with probability

$$\mathbf{P}_{\mathbf{k}} = \frac{\delta_{\mathbf{k}}}{1 - \sum_{i=0}^{\mathbf{k}-1} \delta_{i}} \quad \mathbf{a}$$

 $\delta_k : I^k \times R^k \longrightarrow [0,1]$  (usually  $\delta_k$  will be 0 or 1).

For any  $T \in \mathcal{T}$  we require that for each  $f \in \mathcal{F}$  together leading to the sequence  $x_1$ ,  $f(x_1)$ ,  $x_2$ ,  $f(x_2)$ ,... we shall have (3)  $\sum_{k=0}^{\infty} \delta_k$  (T, sequence) = 1

but each  $\delta_k$  depends only on the first 2k elements of the sequence. As each sequence is completely determined by T and f (T being nonrandomized) we can also write  $\delta_k$ (T,f) and (3) becomes

(3') 
$$\sum_{k=0}^{\infty} \delta_k(T, f) = 1,$$

for all T  $\epsilon \mathcal{T}$  and for all f  $\epsilon \mathcal{F}$ . Condition (3) assures us also that T will stop with probability 1. Finally the  $\ell_i$ 's are the estimates we would make if we stop after i observations; they are also functions of  $x_1, \ldots, x_i$  and  $f(x_1), \ldots, f(x_i)$ .

In a former paper [5] we found the following result concerning minimax procedures. The theorem is stated in slightly more general terminology than of a search problem.

THEOREM 1. Let  $\mathcal{F} \equiv \{f\}$  be a set of "states of nature," X be a set of possible places of observation on f, and  $d_i \equiv D_i(x_1, \dots, x_i, f(x_1), \dots, f(x_i))$ ,  $i = 1, 2, \dots$  be a set of admissible decisions, given the first i observations. Let there be a bounded loss function L(d,f), giving the loss for taking decision d for state f.

(4)  $0 \leq L(d, f) \leq 1$ . BE 3

The n-seq. nonrandomized procedure  $T_n$  and the "true" sequential decision procedure T are defined in analogy with the estimating procedures (1) and (2) respectively, as  $d_n$  replaces  $\ell$  and  $d_i$  replaces  $\ell_i$ ; a constant cost of observation c > 0 is assumed in the latter case. If for every integer  $n \ge 0$  (and  $\varepsilon > 0$ ) there exists a procedure  $T_n^* \varepsilon \mathcal{T}_n$  and a number  $L_n^*$  such that

s) 
$$\sup_{f \in \mathcal{F}} L(d_n(T_n^*, f), f) + (-\varepsilon) \leq L_n^* \leq \sup_{f \in \mathcal{F}} L(d_n(T_n, f), f)$$

for all  $\mathbf{T}_{\mathbf{n}} \in \mathcal{T}_{\mathbf{n}}$  and also

(A) the sequence L\* - L\* → 0 is strictly decreasing, and
 (B) for each "true" sequential procedure T and any given
 integer k > 0, there exists f\* = f\*(k,T) such that

(6) 
$$L(d_1(T, f^*), f^*) > L_2^*$$
, for  $i = 0, 1, 2, ..., k_2$ 

then if we define  $n_0$  such that

(7) 
$$L_{n_0}^* - L_{n_0}^* - 1 > c \ge L_{n_0}^* - L_{n_0}^*$$

or  $n_0 = 0$  if (7) does not hold for any n,  $T_n^*$  is ( $\varepsilon$ ) minimax among all nonrandomized sequential decision procedures  $T \in \mathcal{T}$ . That means, if R(T,f) = E(L(T,f)) + cE(n(T,f)),

(8) 
$$\sup_{f \in \mathcal{F}} R(T_{n_{o}}^{*}, f) = \sup_{f \in \mathcal{F}} EL(d_{n_{o}}(T_{n_{o}}^{*}, f), f) + cEn_{o}$$
$$= \sup_{f \in \mathcal{F}} L(d_{n_{o}}(T_{n_{o}}^{*}, f), f) + n_{o}c$$
$$\leq \sup_{f \in \mathcal{F}} \sum_{i=0}^{\infty} [L(d_{i}(T, f), f) + ic]\delta_{i}(T, f) + (c)$$

for all T  $\in \mathcal{T}$ .

This was applied to the following slach problem.

Let  $\mathcal{F}$  be the class of all unimodal functions on I, that means, for each f  $\varepsilon \mathcal{F}$  there exists  $x^{(f)} \varepsilon I$  such that

121

(9) f is strictly increasing for 
$$x \leq x^{(1)}$$
  
and strictly decreasing for  $x > x^{(f)}$ , or  
strictly increasing for  $x < x^{(f)}$  and  
strictly decreasing for  $x < x^{(f)}$ .

We want to estimate  $x^{(f)}$  by means of an interval estimate [s,t] which has to contain the true  $x^{(f)}$ . Our loss due to estimation will be the length of this interval L([s,t]) = t - s. J.Kiefer [1] found an  $\varepsilon$ -minimax solution for the n-seq. case. His procedure, called the Fibonacci method, which we shall denote  $T^*_n(\varepsilon)$ , is  $\varepsilon$ -minimax, namely, for any given  $\varepsilon > 0$ 

(10) 
$$\sup_{f \in \mathcal{F}} L(T^*(\varepsilon), f) - \varepsilon \leq L^* \leq \sup_{n} L(T, f), \text{ for all } T \in \mathcal{T}_n, f \in \mathcal{F}$$

where  $L(T_n, f)$  is the loss due to estimation resulting from using procedure  $T_n$  on the function f.  $L_n^*$  is known to be  $1/U_{n+1}^{-1}$ ,  $U_n^{-1}$  being the nth Fibonacci number defined as follows:

$$U_{0} = 0, U_{1} = 1, U_{n} = U_{n-2} + U_{n-1}$$
 for  $n \ge 2$ .

By showing that conditions A and B of Th. 1 hold for this problem we showed that the  $\varepsilon$ -minimax sequential procedure here is of a fixed size. Let us consider the simpler search problem of finding a root (zero) of a monotone function.

# PROBLEM 2.

Let  $\mathcal{F}$  be the set of all monotone nonincreasing functions on I with one zero at  $x^{(f)} \in I$ . We want to estimate  $x^{(f)}$  again using interval estimates that have to include the true  $x^{(f)}$  and where the loss due to estimation is the length of this interval.

(a) For each fixed n we have a minimax procedure  $T_n^*$  which is the Bolzano method of taking the next observation at the middle of the "interval of uncertainty". Thus  $T_n^*$  takes

(D.J. Wilde [4]).  $L_n^* = 1/2^n$  which is the length of the last interval of uncertainty.

Conditions (A) and (B) of Theorem 1 are easily seen to hold here, and using the theorem we have for part (b) of this problem we get:

COROLLARY. For problem 2, the Bolzano procedure  $T^*_n$  with n such that

$$\frac{1}{2^{n_0+1}} \leq c < \frac{1}{2^{n_0}}$$

is minimax among all "true" sequential procedures T  $\epsilon$   ${\cal T}$  .

**PROBLEM** 3. A PRIORI DISTRIBUTIONS ON  $x^{(f)}$  AND OPTIMAL PROCEDURES.

Let  $\mathcal{F}$  be the class of all monotone nonincreasing functions on I with one zero  $x^{(f)} \in I$ . This time we assume that f is picked in such a way that there is a continuous *a priori* distribution G of  $x^{(f)}$  on I. Again there is a constant cost c > 0 per observation and for each kind of estimate and loss function we want to find an *optimal* procedure, that is, the one which minimizes the *expected* cost.

For the fixed size procedure we aim to minimize  $E(L(T_n, f))$ , for the sequential procedure to minimize E(R(T, f)). We consider a few particular cases.

(a) We use a point estimate  $\hat{x}$  and assume loss due to estimation

(11) 
$$L(\hat{x}, x^{(f)}) = |\hat{x} - x^{(f)}|.$$

If we are given the a priori G with density g on I and have to estimate  $x^{(f)}$  without taking any observations, the best estimate, namely the one which minimizes E(L), is M the median of G.

From here we would like to proceed to the n-seq. procedure. It may seem that if we take an observation the best place to look for information is where we were going to make our estimate, namely at M. However, this is not true in general as we can see from an example. Consider the a priori distribution G with density g(x) = 2x. Let us take one obser-BE 6 vation and then estimate  $x^{(f)}$ . The median of G is  $M = \frac{1}{\sqrt{2}}$  and the expected loss due to estimation is

$$E(L(M)) \simeq 0.096$$

It is easy to calculate the optimal place of observation  $X_1 \approx 0.645$ with expected loss  $E(L(X_1)) \approx 0.02$ .

For this G and other "nice" a-priori distributions one could, in principle, find n-seq. optimal solutions by working backwards from stage n to 0. A particularly nice G is the uniform one.

(b) Assume a uniform apriori distribution G for the problem in (a).

THEOREM 2. The n-seq. procedure  $T_n^*$  defined by

(12) 
$$T_{n}^{*} = \left\{ \frac{1}{2}, x_{1}^{*} + \operatorname{sgn} f(x) \frac{1}{2^{2}}, \dots, x_{n-1}^{*} + \operatorname{sgn} f(x_{n-1}^{*}) \frac{1}{2^{n}}, \frac{1}{2^{n}} \right\}$$
$$\frac{1}{2^{n}} \max(0, x_{i}^{*} | f(x_{i}^{*}) \ge 0) + \min(1, x_{i}^{*} | f(x_{i}^{*}) \le 0)]$$

is optimal in  $\mathcal{T}_n$  , namely

(13) 
$$E(L(T_n^*)) \leq E(L(T_n))$$
 for all  $T_n \in \mathcal{J}_n$ 

**PROOF.** First we notice that after taking the n observations the best estimate is in the middle of the interval of uncertainty  $V_n$  (we shall use  $V_n$  ambiguously to denote also its length),

$$\mathbf{v}_{n} = \left[\max(\mathbf{o}, \mathbf{x}_{i} | \mathbf{f}(\mathbf{x}_{i}) \geq 0\right], \min(\mathbf{1}, \mathbf{x}_{i} | \mathbf{f}(\mathbf{x}_{i}) \leq 0)\right]$$

giving expected loss

$$E(L|V_n) = \frac{1}{4}V_n$$

This is clear since the aposteriori distribution is uniform on  ${\tt V}_{\tt n}$  . Consequently

(14)  $E(L(T_n)) \ge \frac{1}{4} E(V_n(T_n))$ ,

and we need only prove that  $E(V_n(T_n^*)) \leq E(V_n(T_n))$  for all  $T_n \in \mathcal{J}_n$ . From here on we consider interval estimates which are the intervals of uncertainty. This will only change the scale of the loss function, nultiplying it by 4, and all results will hold as well for the point estimate.

Let  $T_n$  be any procedure in  $\mathcal{T}_n$  and  $V_k$  be the interval of uncertainty after k observations. If  $T_n$  takes the next observation at  $x_{k+1}$  dividing  $V_k$  into two possible  $V_{k+1}$ 's of size  $y_{k+1}V_k$  and  $(1-y_{k+1})V_k$ ,  $0 \le y_k \le 1$ . At this stage we have a uniform apriori on  $V_k$  so we obtain

$$E(V_{k+1}(T_n) | V_k) = [y_{k+1}^2 + (1 - y_{k+1})^2]V_k \ge \frac{1}{2}V_k$$

This inequality holds for all  $y_{k+1}$  and any  $V_k$  . We can conclude that

$$E(V_{k+1}(T_n)) \geq \frac{1}{2} E(V_k(T_n))$$

and therefore

(15) 
$$E(V_n(T_n)) \geq \frac{1}{2^n}$$

for all T<sub>n</sub>.

Since  $E(V_n(T_n^*)) = 1/2^n$  we have established (13) and proved  $T_n^*$  to be optimal.

(c) We now introduce a cost c > 0 per observation, and try to find an optimal "true" sequential procedure, using the interval of uncertainty as estimate. We want to minimize  $E(R) = E(V_n) + cE(n)$ .

Let us consider the procedure  $T^*(c) = T^*_n$  for  $n_0 = n_0(c)$  satisfying

(16) 
$$\frac{1}{2^{n_0}+1} < c \le \frac{1}{2^{n_0}}$$
.

This procedure take<sup>S</sup> observations in the middle of the interval of uncertainty as long as one further observation will reduce the interval of uncertainty by at least c.

This however is not an optimal procedure. The following example will show this.

Let c = 31/120, then T\*(c) = T\_1^\*, and we get R = 1/2 + 31/120 = 91/120. Now consider instead the procedure T' which takes its first observation at  $x_1 < 1/2$ . If f(x)  $\leq 0$  it stops with  $V_1 = [0, x_1]$ , if f(x<sub>1</sub>) > 0 it takes a second observation at  $x_2 = x_1 + 1/2(1 - x_1)$  and stops with

 $V_2 = 1/2(1 - x_1)$ , then

$$E(R(T')) = x_1(x_1 + c) + (1-x_1)[\frac{1}{2}(1-x_1) + 2c]$$
  
=  $\frac{3x_1^2}{2x_1} - (1+c)x_1 + \frac{1}{2} + 2c$ ,  
 $\frac{dE(R(T'))}{dx_1} = 3x_1 - (1+c)$ .

This is minimized for  $x_1 = 1/3 + 1/3$  c = 151/360. Therefore, let T' take  $x_1 = 151/360$ . We obtain E(R(T')) = 65039/86400 < 91/120. It follows that T\* is not optimal. This is even true when c = 1/4,  $E(R(T^*)) = 3/4$ , while taking  $x_1 = 5/12$  in T' leads to

$$E(R(T')) = \frac{71}{96} < \frac{72}{96} = \frac{3}{4}$$
.

We can see, therefore, that for a small enough  $\varepsilon > 0$  we may put  $c = 1.4-\varepsilon$ and have T' take fewer observations than T\* and yet have a smaller expected total cost.

These examples show that T\* is not optimal, but we shall show that T\* is still valuable. While the optimal procedure is hard to calculate in each case,T\* is very simple and moreover:

THEOREM 3. T\* is c-optimal. That means the expected cost from T\* is less than c over the optimal expected cost. The proof of this will follow from the following theorem by D.Blackwell.

We introduce in our problem the following notation:

n = expected sample size.

 $\lambda$  = expected length of final interval.

We are interested in the question: Which pairs  $(n, \lambda)$  are attainable by sequential procedures?

THEOREM \* If  $(n, \lambda)$  is attainable,  $\lambda \ge 2^{-n}$ .

**PROOF.** First we shall notice that if we prove the theorem for procedures that at each stage decide whether to stop or to take another observation without randomization, then it will follow also for procedures that allow randomization. This is a result of the convexity of the function  $2^{-x}$ . So if we decide to stop with probability  $0 \le s \le 1$  and state  $(n_1, \lambda_1)$  and go on with probability 1 - s to state  $(n_2, \lambda_2)$ ,

$$\lambda_{i} \geq 2^{-n_{i}}$$

We are going to attain  $(sn_1 + (1 - s)n_2, s\lambda_1 + (1 - s)\lambda_2)$  and from the convexity it follows that

$$s\lambda_1 + (1-s)\lambda_2 \ge s2^{-n1} + (1-s)2^{-n2} \ge 2^{-sn}1^{-(1-s)n}2$$

So we have to prove the theorem only for the first kind of procedures.

We look at procedures truncated at k steps and use a proof by induction on k. Among procedures truncated at 0, the only point is (0,1) and the result holds. Suppose the result is true for procedures truncated at k. A procedure truncated at k+1 can either take no observations in which case we get (0,1) or it definitely takes one observation. This means it specifies an initial x, 0 < x < 1 and two procedures  $\pi_1$  and  $\pi_2$  truncated at k:

Use  $\pi_1$  when [0,x] occurs,  $\pi_2$  when [x,1] occurs. The resulting (n, $\lambda$  ) is then

$$x(1+n_1, x\lambda_1) + (1-x)(1+n_2, (1-x)\lambda_2)$$

where  $\pi_1$  yields  $(n_1, \lambda_1)$ ,  $\pi_2$  yields  $(n_2, \lambda_2)$ , therefore,  $\pi_1$  applied on an interval [0,x] yields  $(n_1, x\lambda_1)$ , etc..

So we must show that if  $\lambda_1 \ge f(n_1)$  and  $\lambda_2 \ge f(n_2)$ , then

$$x^{2}\lambda_{1} + (1-x)^{2}\lambda_{2} \ge f(1 + xn_{1} + (1-x)n_{2})$$

where  $f(t) = 2^{-t}$ . For this it is enough to show that

$$[x^{2}f(n_{1}) + (1-x)^{2}f(n_{2})] \ge f[xn_{1} + (1-x)n_{2}]/2$$

for  $n_1 \ge 0$ ,  $n_2 \ge 0$ ,  $0 \le x \le 1$ . Say  $n_1 \ge n_2$  and write  $n_1 - n_2 = t$ . Multiply both sides by  $f(-n_2)$ :

$$2[x^{2}f(t) + (1-x)^{2}] \ge f(xt),$$

put

$$f(t) = a$$
:  $2[x^{2}a + (1-x)^{2}] \ge a^{x}$ ,  $2^{-k} \le a \le 1$ ,  $0 \le x \le 1$ .

Fix x and maximize  $a^{X} - 2x^{2}a = \phi(a)$  over a. Let us maximize it over a larger range: 0 < a, which may, if anything, increase the maximum which we want to show to be less than or equal to  $2(1-x)^{2}$ .  $\phi^{*}(a) = xa^{X-1} - 2x^{2} + \phi' = 0$ :  $a^{X-1} = 2x$ :  $a = (2x)^{1/X-1}$ . The maximum of  $\phi(a)$  occurs at  $a = (2x)^{1/X-1}$  and is  $(1-x)^{X/X-1} - 2x^{2}(2x)^{1/X-1} = (1-x)(2x)^{X/X-1}$ . Is for  $0 \le x \le 1$ ,

$$(1-x)(2x)^{x/x-1} \leq 2(1-x)^2?$$

i.e., is 
$$(2x)^{x/x-1}(1-x)^{-1} \leq 2$$
?

i.e., is 
$$(2x)^{x}(1-x)^{1-x} \ge 2^{x-1}$$
?

i.e., is 
$$x^{x}(1-x)^{1-x} \ge \frac{1}{2}$$
?

Yes. Put  $\Psi(x) = x^{X}(1-x)^{1-x}$ ; then

 $\alpha = \log \Psi = x \log x + (1-x) \log(1-x)$ 

is convex and symmetric about 1/2, assuming its minimum at x = 1/2. Since  $\Psi(1/2) = 1/2$  this completes the proof of Theorem \*.

Now we come back to prove Theorem 3.

PROOF. Using T\* we get an expected cost

$$E(R(T^*,c)) = nc + 2^{-n}$$

as  $2^{-n-1} < c \le 2^{-n}$ .

For any other T we get an expected number of observations n(T) = x and  $E(L(T)) = \lambda \ge 2^{-x}$  due to Theorem \*. It follows that

$$E(R(T,c)) \geq \min (2^{-x} + xc) .$$
  
$$0 \leq x^{<\infty} \qquad BE 11$$

Put  $\rho(x) = 2^{-x} + xc$ , then  $E(R^{*},c) = \rho(n)$ ,

$$\rho'(x) = -2^{-x} \ln \beta + c$$
.

So the minimum occurs at c =  $2^{-x} \ln 2 \approx 2^{-x} \cdot 0.6931$ . It follows that

$$2^{n-1} \le 0.6931 \times 2^{-x} = c \le 2^{-n}$$

therefore,

$$2^{-n-1} < 2^{-x} < 2^{-n+1}$$

So  $\rho(x) \ge (n-1)c + 2^{-n-1} \ge \rho(n) - 2c$ . We can do better by considering separate cases: If x = n,

$$\rho(\mathbf{x}) = \rho(\mathbf{n}).$$

If x > n,  $\rho(x) \ge nc + 2^{-n-1} \ge \rho(n) - c$ .

If 
$$x < n$$
,  $\rho(x) \ge (n-1)c + 2^{-n} = \rho(n) - c$ .

Therefore, T\* is c-optimal.

Finally let us note that the last two results,  $T_n^*$  optimal in  $\mathcal{J}_n$  and  $T^*$  c-optimal, remain true when we allow also randomized procedures.

PROOF. Let  $V_i$  be the interval of uncertainty after i steps and let us take the next observation at  $y_{i+1}$  in such a way that  $V_i$  is divided into two intervals of lengths  $xV_i$  and  $(1-x)V_i$ , where  $y_{i+1}$  is a random variable and therefore x is a random variable on [0,1]. Then

$$E(V_{i+1} | V_i, x) = E(x^2 + (1-x)^2)V_i \ge {[[E(x)]^2 + [1-E(x)]^2]V_i}.$$

Therefore, we could have taken x = E(x) constant and done at least as well.

(d) Let us consider again a point estimate  $\hat{x}$  for  $x^{(f)}$  but with a square error loss function

$$L(X^{(f)}, \hat{X}) = (X^{(f)} - \hat{X})^2$$

For this case we obtained results exactly analogous with the results for  $$\operatorname{BE}\ 12$$ 

the loss function  $L = |X^{(f)} - \hat{X}|$  including theorems 2 and 3.

RANDOMIZED PROCEDURES.

Let us return to Problem 2. In view of the results obtained for uniform *a priori* distributions we may try to see what we can say about randomizations and minimax Procedures.

THEOREM 4. For Problem 2(a)  $T^*_n$  is minimax among all randomized n-seq. procedures.

PROOF. We know that  $T_n^*$  is the optimal procedure (Bayes solution) among all randomized procedures when a uniform a priori distribution is assumed.  $T_n^*$  has constant risk  $R = 1/2^n$  for all  $x^{(f)} \epsilon$  I, therefore following a known theorem and a lemma found in Lehmann [6] the uniform distribution is "least favorable" and  $T_n^*$  is minimax among all *randomized* procedures this time.

Let us go through the proof as a preparation for the next theorem. PROOF.  $T_n^*$  is constant over  $x^{(f)} \in I$ , therefore

 $\sup_{x \in I} R(T^*, x) = \int_{n} R(T^*, x) dx \leq \int_{n} R(T, x) dx$   $x \in I \qquad I \qquad I \qquad by optimality of T^*_{n}$ 

 $\stackrel{\leq \text{ sup } R(T_n, x)}{\underset{n \in I}{x \in I}}$  for all  $T_n \in \mathcal{J}_n$ , where  $\mathcal{J}_n$  is the class of all randomized n-seq. procedures. Q. E. D.

Now let us look at the problem 2(b). We know that T\*(c) is c-optimal among all randomized procedures, given a uniform a priori distribution. T\* has also a constant cost.

THEOREM 5. T\*(c) is c-minimax among all randomized procedures.

PROOF. 
$$\begin{aligned} \sup_{x \in I} R(T^*, x) &= \int R(T^*, x) dx &\leq \int R(T, x) dx - c \leq \sup_{x \in I} R(T, x) - c \\ x \in I & I & I & x \in I \\ & & & & & \\ & & & & \\ & & & & \\ & & & & & \\ & & & & \\ & & & & & \\ & & & & \\ & & & & & \\ & & & & & \\ & & & & & \\ & & & & & \\ & & & & & \\ & & & & & \\ & & & & & \\ & & & & & \\ & & & & & \\ & & & & & \\ & & & & & \\ & & & & & \\ & & & & & \\ & & & & & & \\ & & & & & \\ & & & & & \\ & & & & & \\ & & & & & \\ & & & & & \\ & & & & & \\ & & & & & \\ & & & & & \\ & & & & & \\ & & & & & \\ & & & & & \\ & & & & & \\ & & & & & \\ & & & & & \\ & & & & & \\ & & & & & \\ & & & & & & \\$$

## EXTENSIONS.

There are many possible extensions of the search problem some of which were explored by Kiefer [2] and Wilde [4] and several others. We did some further work considering restricting the class  $\mathcal{F}$  by putting a bound on the slope of f. This produced several new concepts and interesting examples.

J. Kiefer defines in [2] the *order* of the search problem as the minimum number of observations needed in order to be sure to obtain some information on the location of  $x^{(f)}$ .

Here we considered only first and second order problems. In [2] Kiefer solves the minimax n-seq. problem for the third order search, the search for an inflection point. It would be interesting to see the effect of randomization on the second and third order procedures, as well as the effect of various a priori distributions on these problems. It seems likely also that one would be able to find "almost optimal procedures" for nonuniform a priori distributions at least for the first order problem.

Another important extension would be the consideration of different loss functions; for instance, for the first order search problem with an a priori distribution given, define the loss function as the left end of the interval of uncertainty. Professor B. McGuire has some unpublished results on this problem which arose from a simplified practical problem.

The most important extension, in my opinion, is the extension to higher dimensions. This may also be the hardest extension. Kiefer reports on the difficulties involved already in the two dimensional search in [2]. We have considered the search for two zeros on an interval I, or equivalently, the search for an indicator function of a subinterval of I. This problem may be reduced to a search for a point of the subinterval. Under the assumption that we know a positive lower bound to the length of the sought interval, we could find a minimax search procedure, yet this procedure was not satisfactory, being inadmissible, actually dominated by many other procedures out of which we BE 14

could not find a best procedure.

. . .

Another extension would be considering different families of functions. This may not bring too many new results. Kiefer also reports on this in [2] and deals in particular with the problem of search on a lattice, where the domain of f is just a finite number of points. This is done for both one and two dimensions.

A last extension that is also mentioned by Kiefer is considering problems in which errors are involved in the observation. A paper on this subject was written by Kiefer and Wolfowitz [7].

#### REFERENCES.

 Kiefer, J., Sequential minimax search for a maximum, Proc. Amer. Math. Soc., (1953), Vol. 4, pp. 502-506.

1

- [2] Kiefer, J., Optimum sequential search and approximation methods under minimum regularity assumption, <u>J.Soc.Indust.Appl. Math</u>. (September 1956), Vol. 5, No. 3.
- [3] Johnson, S.M., Best exploration for maximum is Fibonaccian, <u>The Rand</u> Corporation, May 4, 1956.
- [4] Wilde, D.J., Optimum Seeking Methods, Englewood Cliffs, New Jersey, Prentice Hall, 1964.
- [5] Eichhorn, B.H., On Sequential Minimax, J.of Math. Analysis & Applications Vol. 14, No.1, April 1966.
- [6] Lehmann, E. L., Notes on the theory of estimation, University of
   [California, 1949-50.
   [7] Kiefer, J. and Wolfowitz, J., Stochastic estimation of the maximum of
- a regression function, <u>Ann.Math.Statist</u>.(1952),Vol.23,pp.162-466. BE 15