

A NOTE ON SOME METHODS FOR REGRESSION ANALYSIS WITH INCOMPLETE OBSERVATIONS

By R. D. GILL

*Centre for Mathematics and Computer Science,
Amsterdam, Netherlands*

SUMMARY. Some recent related proposals for estimating regression coefficients with incomplete observations are discussed. The proposals included approximate standard errors for the estimators. It is shown that the estimators of the regression coefficients are consistent under fairly weak conditions, but that only under rather strong ones can the usual (asymptotic) tests of significance be validly based on the estimated coefficients and the computed standard deviations. The results depend heavily on whether a random or a fixed specification for the covariates can be assumed.

I. INTRODUCTION

Very many procedures, both specific and general, have been suggested in the literature for dealing with the problem of incomplete observations in regression analysis; see the papers of Affi and Elashoff (1966, 1967, 1969a, b), Hartley and Hocking (1971) and Dempster, Laird and Rubin (1977). However few of the methods which are applicable in a general regression analysis situation give consistent estimators of the regression coefficients, and still fewer show how asymptotic standard deviations may be validly estimated (in order to carry out the usual *t*-tests, etc.). There are three very similar proposals which do at least give suggestions in this direction though little theoretical justification is given: these, the subject of this note, are Beale and Little's (1975) "method 5" and "method 6" and the method of Dagenais (1973).

Let us briefly sketch the kind of situation we are interested in. Each of N observations if complete would be a $(K+1)$ -vector of values taken by K independent or *predictor* variables and 1 dependent or *criterion* variable. However for some observations, the values taken by some of the predictor variables are missing. We suppose that the mechanism causing this works independently of that generating the values of both predictor and criterion variables. This assumption is only implicitly made in the sequel, but it is an assumption of major importance (as is usual in the literature on this subject).

Mathematics subject classification (1980) : 62J05.

Key words and phrases : Incomplete observations, Missing data, Regression analysis.

We shall work conditionally on the realized patterns of missing and non-missing values. For the sake of simplicity we assume that none of the values taken by the criterion variable are missing.

We wish to make as few assumptions as possible about the predictor variables. In particular we certainly want to avoid the assumption that each complete $(K+1)$ -vector observation is a drawing from a multivariate normal distribution. In practical applications of regression analysis it is usually possible to classify each predictor variable as being either a *fixed* "design variable" or a *random* "covariate". By design variables, we mean variables which are preset by the experimenter as part of the planned experimental design. Covariates on the other hand are variables which are generated by some stochastic mechanism jointly with the criterion variables, so that covariates and criterion variable together can be considered as drawn from some multivariate distribution (generally a different distribution for each set of values of the design variables).

Even if the natural specification of the predictor variables is random, the regression model may be specified in terms of the conditional distribution of the criterion variable given the predictor variables. So if one starts with a random specification of the predictor variables, one can step over to a fixed specification by conditioning. One also has the natural possibility, when some of the predictor variables are missing, of just conditioning on the non-missing predictor variables.

Our main aim is in fact to discover whether the methods for dealing with incomplete observations mentioned above are only appropriate with one of these possible specifications of the predictor variables. All three proposals work by filling in the missing values in each observation with least squares predictions based on the non-missing predictor variables in that observation; the coefficients needed for this are estimates based on all the present data. Then a standard weighted least squares regression analysis is carried out on the completed data set, supplying both estimates of the regression coefficients and standard errors for them. Weights are needed because the least squares prediction introduces an extra error of varying size in each incomplete observation. The proposals only differ in how the coefficients for the least squares predictions and how the weights for the final regression analysis are to be estimated (they all agree on what these coefficients and weights should be). Since it turns out that only consistency of the estimators of these quantities is needed, the differences between the proposals are not crucial.

We are interested in the problems of finding reasonable and non-technical conditions under which (i) the proposals yield consistent estimators of the regression coefficients, and (ii), suitably normed, these estimators are asymptotically normally distributed about the true regression coefficients with a covariance matrix which is *consistently estimated by that produced in the weighted least squares regression analysis*. Problem (i) turns out to have a satisfactory solution whatever the specification of the predictor variables. However in (ii), asymptotic normality is only easily proved with random predictor variables. Even when asymptotic normality can be proved, the asymptotic covariance matrix of the regression coefficient estimators only coincides with the limiting value of the covariance matrix estimator under conditions quite close to multivariate normality of the predictor variables.

In the next section we specify our general model, define the estimators, and prove consistency. Section 3 looks at asymptotic normality while in the final section we briefly discuss some implications of our results.

2. PROBLEM (i) : CONSISTENCY

First some notation. Random variables will be underlined, so that the same symbol not underlined represents a possible value of the corresponding variable. a^T denotes the transpose of the vector a . We specify a model for N observations for each $N = 1, 2, \dots$; all quantities (including the underlying sample space) may depend on N unless explicitly stated otherwise, though this dependence is generally suppressed in the notation. We write $\xrightarrow{\mathcal{P}}$ and $\xrightarrow{\mathcal{D}}$ for convergence in probability and in distribution respectively (always as $N \rightarrow \infty$) and denote a multivariate normal distribution with given mean vector and covariance matrix by $\mathcal{N}(\cdot, \cdot)$.

Let P and M (a *pattern* of observed predictor variables and its complement of *missing* ones) denote sets of indices such that $P \cup M = \{1, \dots, K\}$ (where K is the number of predictor variables), $P \cap M = \phi$, $P \neq \phi$ and if e.g. the *first* predictor variable is the constant 1, $1 \in P$. Vectors and matrices will often be partitioned according to P and M , e.g. if β is a $K \times 1$ vector and Σ a $K \times K$ matrix then

$$\beta = \begin{pmatrix} \beta_P \\ \beta_M \end{pmatrix},$$

$$\Sigma = \begin{pmatrix} \Sigma_{PP} & \Sigma_{PM} \\ \Sigma_{MP} & \Sigma_{MM} \end{pmatrix} = (\Sigma_P \ \Sigma_{\cdot M}). \quad \dots \quad (1)$$

Let $(\underline{x}^n, \underline{y}^n, \underline{e}^n)$, $n = 1, \dots, N$, denote the complete $K \times 1$ vector of predictor variables, the criterion variable, and the disturbance variable for the n -th observation, related by

$$\underline{y}^n = \beta^\top \underline{x}^n + \underline{e}^n \quad \dots (2)$$

for some fixed $K \times 1$ vector β of regression coefficients which we want to estimate. Let P^n and M^n , $n = 1, \dots, N$ be patterns of observed and missing predictor variables; the data consists of $(\underline{y}^n, \underline{x}_P^n, P^n)$, $n = 1, \dots, N$, where we have written \underline{x}_P^n for \underline{x}_{-P}^n . Similarly we often write \underline{x}_M^n for the unobserved \underline{x}_{-M}^n . To (2) we add the usual assumptions

$$\begin{aligned} \mathcal{E}(\underline{e}^n) &= 0 && \text{for all } n \\ \mathcal{E}(\underline{x}^n \underline{e}^{n'}) &= 0 && \text{for all } n, n' \quad \dots (3) \\ \mathcal{E}(\underline{e}^n \underline{e}^{n'}) &= \begin{cases} 0 & n \neq n' \\ \sigma^2 > 0 & n = n' \end{cases} \end{aligned}$$

where σ^2 like β does not depend on N . The second equality in (3) implies the first one if (2) includes a constant term, e.g.

$$\underline{x}_1^n = 1 \quad \text{almost surely for each } n. \quad \dots (4)$$

Note that our model does not yet assume independence or identical distributions for the N observations, nor have we excluded the predictor variables from being fixed design variables. We make the following further assumptions, which we shall illustrate with some important examples in a moment. For each pattern P let ρ_P be a non-negative number, and let Σ be a fixed $K \times K$ symmetric positive definite matrix. The symbol $\#$ denotes the number of elements in a set. Suppose that for each P , the following convergences hold as $N \rightarrow \infty$:

$$\begin{aligned} \text{A1} & \quad N^{-1} \# \{n : P^n = P\} \rightarrow \rho_P, \\ \text{A2} & \quad N^{-1} \sum_{n : P^n = P} \underline{x}_P^n \underline{x}_P^{n\top} \xrightarrow{\mathcal{P}} \rho_P \Sigma_{PP}, \\ \text{A3} & \quad N^{-1} \sum_{n : P^n = P} \underline{x}_P^n \underline{x}_M^{n\top} \xrightarrow{\mathcal{P}} \rho_P \Sigma_{PM}, \\ \text{A4} & \quad N^{-1} \sum_{n : P^n = P} \underline{x}_P^n \underline{e}^n \xrightarrow{\mathcal{P}} 0. \end{aligned}$$

Defining $\sigma_P^2 = \sigma^2 + \beta_M^\top (\Sigma_{MM} - \Sigma_{MP} \Sigma_{PP}^{-1} \Sigma_{PM}) \beta_M$ suppose we also have

$$\text{A5} \quad A = \sum_P \rho_P \sigma_P^{-2} \Sigma_{\cdot P} \Sigma_{PP}^{-1} \Sigma_P \text{ is non-singular.}$$

Interpreting Σ as a limiting average value of $\underline{x}^n \underline{x}^{n\tau}$, A2 and A3 together with A1 express the fact that the patterns of missing values are at least asymptotically not influenced by the predictor variables, while A4 expresses the same fact for the disturbance. The role of A5 will become clear later.

Example 1: Random predictor variables: $(\underline{x}^n, \underline{e}^n)$, $n = 1, \dots, N$, are independent over n and have the same distribution for all n and N , with

$$\mathfrak{E}(\underline{x}^n \underline{x}^{n\tau}) = \Sigma. \quad \dots \quad (5)$$

In this example A2, A3 and A4 follow from (3), (5) and A1 by the weak law of large numbers, taking special care for P such that $\rho_P = 0$. Here we use the fact that if $\underline{u}_k \xrightarrow{\mathfrak{P}} a$ as $k \rightarrow \infty$ for some sequence of random variables \underline{u}_k , and if k_N satisfies $k_N/N \rightarrow \rho$ as $N \rightarrow \infty$, then $(k_N/N) \underline{u}_{k_N} \xrightarrow{\mathfrak{P}} \rho a$ as $N \rightarrow \infty$.

This is trivial when $\rho > 0$; when $\rho = 0$ we still have that the sequence \underline{u}_{k_N} is bounded in probability and the result is then again trivial.

Example 2: Fixed predictor variables: For some vector x^n (depending also on N), $\underline{x}^n = x^n$ almost surely for each n and N ; these vectors satisfy A2 and A3 with convergence in probability replaced by ordinary convergence.

We now have that A4 follows from (3) and A2, since $N^{-1} \sum_{n:P^n=P} x_p^n \underline{e}^n$ has expectation zero and covariance matrix $(\sigma^2 N^{-1} \sum_{n:P^n=P} x_p^n x_p^{n\tau})/N$. This example can be obtained from Example 1 by conditioning. Suppose we are in the situation of Example 1 with $(\underline{x}^n, \underline{e}^n, P^n)$, $n = 1, \dots, N$ being the first N elements of a single infinite sequence. Then for almost all realizations $(\underline{x}^1, \underline{x}^2, \dots) = (x^1, x^2, \dots)$, conditions A2 and A3 continue to hold for the conditional distribution of the data given this realization of $(\underline{x}^1, \underline{x}^2, \dots)$. We need to assume separately that (3) also holds with the expectations replaced by conditional expectations given $\underline{x}^n = x^n$.

Example 3: Non-missing predictor variables fixed, missing predictor variables random: for some vectors x_p^n (depending also on N), $\underline{x}_p^n = x_p^n$ almost surely for each n and N ; these vectors satisfy A₂ with convergence in probability replaced by ordinary convergence.

Again A4 follows from (3) and A2, but A3 needs to be assumed separately. This example will generally arise from Example 1 by conditioning. Consider the situation of Example 1 with $(\underline{x}^n, \underline{e}^n, P^n)$, $n = 1, \dots, N$, being the first N elements of a single infinite sequence. Then for almost all realizations $(\underline{x}_p^1, \underline{x}_p^2, \dots) = (x_p^1, x_p^2, \dots)$, condition A2 continues to hold for the conditional distribution of the data given (x_p^1, x_p^2, \dots) . We need to assume A3 and the conditional form of (3) separately.

We shall not give explicit definitions of the estimators proposed in Dagenais (1973) and in Beale and Little's (1975) "method 5" and "method 6" but first work as if the parameters needed for the proposals (certain functions of σ^2 , β and Σ) were known. Define

$$\alpha_{MP} = \Sigma_{MP} \Sigma_{PP}^{-1} \quad (\text{where } \Sigma_{PP}^{-1} = (\Sigma_{PP})^{-1}.) \quad \dots (7)$$

$$\sigma_P^2 = \sigma^2 + \beta_M^T (\Sigma_{MM} - \Sigma_{MP} \Sigma_{PP}^{-1} \Sigma_{PM}) \beta_M \quad \dots (8)$$

$$\hat{\underline{x}}^n = \Sigma_{.P} \Sigma_{PP}^{-1} \hat{\underline{x}}_P^n = \begin{pmatrix} \alpha_P^n \\ \alpha_{MP} \hat{\underline{x}}_P^n \end{pmatrix} \quad \text{where } P = P_n \quad \dots (9)$$

$$\hat{\underline{x}} = (\hat{\underline{x}}^1, \dots, \hat{\underline{x}}^N)^T \quad \dots (10)$$

$$\hat{\Sigma} = \text{the } N \times N \text{ diagonal matrix with diagonal elements } \sigma_{Pn}^2, n = 1, \dots, N \quad (11)$$

$$\underline{Y} = (\underline{y}^1, \dots, \underline{y}^N)^T \quad \dots (12)$$

$$\hat{\beta} = (\hat{\underline{X}}^T \hat{\Sigma}^{-1} \underline{X})^{-1} (\hat{\underline{X}}^T \hat{\Sigma}^{-1} \underline{Y}) \quad \text{if } \hat{\underline{X}}^T \hat{\Sigma}^{-1} \hat{\underline{X}} \text{ is non-singular} \quad \dots (13)$$

If α_{MP} and σ_P^2 were known, $\hat{\beta}$ would be the proposed estimator of β and $(\hat{\underline{X}}^T \hat{\Sigma}^{-1} \hat{\underline{X}})^{-1}$ the proposed approximate covariance matrix for it. In fact in Example 1 $\hat{\underline{x}}_M^n$ is the best linear predictor of \underline{x}_M^n based on \underline{x}_P^n while σ_{Pn}^2 is the expanded variance of the error term in (2) if \underline{x}_M^n is replaced there by $\hat{\underline{x}}_M^n$. For defining

$$\hat{\underline{e}}^n = \underline{e}^n + \beta_M^T (\underline{x}_M^n - \alpha_{MP} \underline{x}_P^n) \quad (\text{where } P = P^n \text{ and } M = M^n) \quad \dots (14)$$

we rewrite (2) as

$$\underline{y}^n = \beta^T \hat{\underline{x}}^n + \hat{\underline{e}}^n \quad n = 1, \dots, N \quad \dots (15)$$

where in Example 1

$$\mathcal{E}(\hat{\underline{e}}^n) = 0$$

$$\mathcal{E}(\hat{\underline{x}}^n \hat{\underline{e}}^{n'}) = 0 \quad (\text{c.f. (3)}) \quad \dots (16)$$

$$\mathcal{E}(\hat{\underline{e}}^n \hat{\underline{e}}^{n'}) = (\hat{\Sigma})_{nn'}$$

After conditioning on $\underline{x}_P^n = x_P^n$, $n = 1, \dots, N$ (Example 3), (16) no longer necessarily holds, while in Example 2 it is generally false.

Theorem 1 : Under A1 to A5, $\hat{\beta}$ defined by (13) is a consistent estimator of β and $N(\hat{\underline{X}}^T \hat{\Sigma}^{-1} \hat{\underline{X}})^{-1}$ is a consistent estimator of A^{-1} . These statements are also true if in the definitions (7) to (13), α_{MP} and σ_P^2 are replaced by consistent estimators α_{MP} and σ_P^2 of the same quantities.

Proof: We first look at the estimation of A^{-1} .

$$\begin{aligned} N^{-1} \underline{\hat{X}}^{\tau} \underline{\hat{\Sigma}}^{-1} \underline{\hat{X}} &= \sum_P \sigma_P^{-2} \Sigma_{.P} \Sigma_{PP}^{-1} \left(N^{-1} \sum_{n: P^{n=P}} \underline{x}_P^n \underline{x}_{PP}^{n\tau} \right) \Sigma_{PP}^{-1} \Sigma_P. \\ &\xrightarrow{\mathcal{P}} \sum_P \sigma_P^{-2} \Sigma_{.P} \Sigma_{PP}^{-1} \rho_P \Sigma_{PP} \Sigma_{PP}^{-1} \Sigma_P = A. \end{aligned}$$

Because A is non-singular, the probability that $N^{-1} \underline{\hat{X}}^{\tau} \underline{\hat{\Sigma}}^{-1} \underline{\hat{X}}$ is non-singular converges to 1 as $N \rightarrow \infty$ and hence

$$N(\underline{\hat{X}}^{\tau} \underline{\hat{\Sigma}}^{-1} \underline{\hat{X}})^{-1} \xrightarrow{\mathcal{P}} A^{-1}. \quad \dots (17)$$

Next defining

$$\underline{\hat{E}} = (\hat{e}^1, \dots, \hat{e}^N)^{\tau} \quad \dots (18)$$

we can rewrite (15) as

$$\underline{Y} = \underline{\hat{X}} \beta + \underline{\hat{E}} \quad \dots (19)$$

and so by (13) and (17), with probability converging to 1,

$$\underline{\hat{\beta}} = \beta + N(\underline{\hat{X}}^{\tau} \underline{\hat{\Sigma}}^{-1} \underline{\hat{X}})^{-1} N^{-1} \underline{\hat{X}}^{\tau} \underline{\hat{\Sigma}}^{-1} \underline{\hat{E}}. \quad \dots (20)$$

So to prove

$$\underline{\hat{\beta}} \xrightarrow{\mathcal{P}} \beta \quad \dots (21)$$

it suffices to establish

$$N^{-1} \underline{\hat{X}}^{\tau} \underline{\hat{\Sigma}}^{-1} \underline{\hat{E}} \xrightarrow{\mathcal{P}} 0 \text{ as } N < \infty. \quad \dots (22)$$

Now,

$$\begin{aligned} N^{-1} \underline{\hat{X}}^{\tau} \underline{\hat{\Sigma}}^{-1} \underline{\hat{E}} &= \sum_P \sigma_P^{-2} \Sigma_{.P} \Sigma_{PP}^{-1} \left(N^{-1} \sum_{n: P^{n=P}} \underline{x}_P^n \hat{e}^n \right) \\ &= \sum_P \sigma_P^{-2} \Sigma_{.P} \Sigma_{PP}^{-1} \left(N^{-1} \sum_{n: P^{n=P}} (\underline{x}_P^n \underline{e}^n + (\underline{x}_P^n \underline{x}_M^{n\tau} - \underline{x}_P^n \underline{x}_P^{n\tau} \alpha_{MP}^{\tau}) \beta_M) \right) \\ &\xrightarrow{\mathcal{P}} 0 \quad \dots (23) \end{aligned}$$

by A2, A3, A4 and (7).

Finally even if α_{MP} and σ_P^2 are everywhere replaced by consistent estimators of the same quantities, all the above arguments remain valid. \square

Remark 1: The consistency of the estimators of α_{MP} and σ_P^2 in Beale and Little's (1975) "method 6" can be established by the same type of arguments as in Gill (1977) even though they derive their estimators from

considerations of maximum likelihood under multivariate normality of $(\underline{x}^n, \underline{e}^n)$, $n = 1, \dots, N$. Suitable conditions for consistency are A1, A2 and A4 supplemented with

A6 For P such that $\rho_P = 0$, $\# \{n: P^n = P\} = 0$ for sufficiently large N , and

A7
$$N^{-1} \sum_{n: P^n = P} (e^n)^2 \xrightarrow{\mathcal{P}} \rho_P \sigma^2 \text{ for all } P.$$

The (more complicated) other two proposals are harder to analyse though a similar approach is applicable; we have not worked out the full details. Since consistency of the estimators of α_{MP} and σ_P^2 is all that is needed (\sqrt{N} -consistency of the estimator of α_{MP} for asymptotic normality; see next section), less complicated proposals can easily be made which still have the required properties. The proof of Theorem 1 actually also shows consistency of Beale and Little's (1975) "method 4", where weights are not introduced.

Remark 2: If \underline{x}_M^n is predicted by regression on \underline{x}_P^n and y^n for each n , the resulting estimator of β is generally inconsistent. For instance in Example 1, if we let $\hat{\underline{x}}_M^n$ be the best linear predictor of \underline{x}_M^n based on \underline{y}^n and \underline{x}_P^n , and write

$$\underline{y}^n = \beta_P^\top \underline{x}_P^n + \beta_M^\top \hat{\underline{x}}_M^n + \hat{\underline{e}}^n,$$

then we find that in general $\mathcal{E} \underline{x}_P^n \hat{\underline{e}}^n \neq 0$ and so it does not hold that $N^{-1} \sum_{n: P^n = P} \underline{x}_P^n \hat{\underline{e}}^n \xrightarrow{\mathcal{P}} 0$ if $\rho_P > 0$. This fact makes another of Beale and Little's (1975) proposals (middle of their section 5) rather difficult to motivate.

3. PROBLEM (ii): ASYMPTOTIC NORMALITY WITH CORRECT COVARIANCE MATRIX

Reviewing the proof of Theorem 1, we see that under the conditions of that theorem,

$$N^{\frac{1}{2}}(\underline{\hat{\beta}} - \underline{\beta}) \xrightarrow{\mathcal{D}} \mathcal{N}(0, A^{-1}) \quad \dots \quad (24)$$

if and only if

$$N^{-\frac{1}{2}} \underline{\hat{X}}^\top \hat{\Sigma}^{-1} \underline{\hat{E}} \xrightarrow{\mathcal{D}} \mathcal{N}(0, A). \quad \dots \quad (25)$$

Since then (17) holds too we can indeed validly use $(\underline{\hat{X}}^\top \hat{\Sigma}^{-1} \underline{\hat{X}})^{-1}$ as an asymptotic covariance matrix for $\underline{\hat{\beta}}$ and carry out the usual tests of significance on regression coefficients. We shall prove a theorem giving conditions for (24) to hold in the special case of Example 1, but shall give some heuristic arguments that it cannot hold in Example 2, and only holds under rather special conditions in Example 3.

Theorem 2 : Consider the situation of Example 1 and suppose furthermore

$$\mathcal{E}((\hat{\epsilon}^n)^2 \underline{x}_P^n \underline{x}_P^{n\tau}) = \Psi_P \quad \dots \quad (26)$$

for some finite matrices Ψ_P . Then

$$N^{-\frac{1}{2}} \underline{\hat{X}}^\tau \hat{\Sigma}^{-1} \underline{\hat{E}} \xrightarrow{\mathcal{D}} \mathcal{N}(0, B) \quad \dots \quad (27)$$

where

$$B = \sum_P \rho_P \sigma_P^{-4} \Sigma_P \Sigma_{PP}^{-1} \Psi_P \Sigma_{PP}^{-1} \Sigma_P. \quad \dots \quad (28)$$

A sufficient condition for (24) to hold (i.e. for equality of A and B) is

$$\Psi_P = \sigma_P^2 \Sigma_{PP} \text{ for all } P. \quad \dots \quad (29)$$

These results also hold with α_{MP} and σ_P^2 replaced by consistent estimators provided the estimator of α_{MP} is actually \sqrt{N} -consistent, i.e. $N^{\frac{1}{2}}(\hat{\alpha}_{MP} - \alpha_{MP})$ is bounded in probability as $N \rightarrow \infty$.

Proof : Multiplying (23) by $N^{\frac{1}{2}}$ and recalling (16), we see that by the Central Limit Theorem, (again with special care for P such that $\rho_P = 0$),

$$N^{-\frac{1}{2}} \underline{\hat{X}}^\tau \hat{\Sigma}^{-1} \underline{\hat{E}} \xrightarrow{\mathcal{D}} \mathcal{N}(0, B). \quad \dots \quad (30)$$

Obviously if (29) holds, then $A = B$. The proof of the last part of the theorem is straightforward. \square

Remark 3 : Theorem 2 is a satisfactory solution to problem (ii) if we can consider the predictor variables as random and can assume that the complete observations would have been independent and identically distributed, with $(\hat{\epsilon})^2$ uncorrelated with $\underline{x}_P^n \underline{x}_P^{n\tau}$. This does not seem a very heavy assumption; we already have that $\hat{\epsilon}_n$ is uncorrelated with \underline{x}_P^n .

Of course one could often be reluctant to assume that the \underline{x}^n 's are random variables at all. However, it is easy to see that (27) cannot hold under reasonable conditions in the case of Example 2, even with a different definition of the matrix B . In (27) $\underline{\hat{X}} = \hat{X}$ is now non-random, and $\underline{\hat{E}}$ is the sum of a random and a non-random component, so the left hand side of (27) itself splits into a random and a non-random part. There is no reason why the non-random part should converge at all. For instance suppose that the situation of Example 2 has arisen by starting in Example 1 and conditioning on $\underline{x}^n = \underline{x}^n$, $n = 1, 2, \dots$, as in the discussion after Example 2. Suppose the conditional distribution of $\underline{\epsilon}^n$ does not depend on \underline{x}^n .

Looking at (27), (18) and (14) we see that under the assumptions of Theorem 2, *unconditionally*, both parts of the left hand side of (27) converge in distribution to in general non-degenerate normal distributions. *Conditional* on $\underline{x}^n = x^n$, $n = 1, 2, \dots$ the random part still converges in distribution to a normal distribution with mean vector zero. The other part, now non-random, would have to converge to zero for (27) to be valid. But the probability must be zero that such x^n 's have been realized, in view of the unconditional non-degenerate limiting normal distribution.

We finally turn to Example 3, supposing it has arisen by starting from Example 1 and conditioning on the non-missing predictor variables as described previously. If A1 holds, then with probability 1 after conditioning A2 holds too, and this implies that

$$N^{-1} \hat{\underline{X}}^T \hat{\underline{\Sigma}}^{-1} \hat{\underline{X}} \rightarrow A \quad \dots (31)$$

($\hat{\underline{X}}$ is now non-random). However the validity of A3 and A4 depends on the conditional distributions of \underline{x}_M^n and \underline{e}^n given $\underline{x}_P^n = x_P^n$. Let us make the rather strong assumption (it implies for instance (29), and is itself implied by multivariate normality of $(\underline{x}^n, \underline{e}^n)$) that these are such that for *all* P and x_P^n

$$\mathcal{E}(\hat{\underline{e}}^n | \underline{x}_P^n = x_P^n) = 0 \quad \text{(c.f. (16))} \quad \dots (32)$$

$$\mathcal{E}((\hat{\underline{e}}^n)^2 | \underline{x}_P^n = x_P^n) = \sigma_P^2$$

or in words, every regression on \underline{y}^n on a group of variables from \underline{x}^n is linear and homoscedastic. Looking at (20), this now implies that

$$\mathcal{E}(\underline{\hat{\beta}} | \underline{\hat{X}} = \hat{\underline{X}}) = \underline{\beta}$$

$$\mathcal{E}((\underline{\hat{\beta}} - \underline{\beta})(\underline{\hat{\beta}} - \underline{\beta})^T | \underline{\hat{X}} = \hat{\underline{X}}) = (\hat{\underline{X}}^T \hat{\underline{\Sigma}}^{-1} \hat{\underline{X}})^{-1}. \quad \dots (33)$$

By (31) and (33), $\underline{\hat{\beta}}$ is consistent; but more importantly, (33) gives new motivation for using $(\hat{\underline{X}}^T \hat{\underline{\Sigma}}^{-1} \hat{\underline{X}})^{-1}$, with α_{MP} and σ_P^2 replaced by estimates, as an approximate covariance matrix for $\underline{\hat{\beta}}$. Does such a simple argument also give asymptotic normality of $N^{1/2}(\underline{\hat{\beta}} - \underline{\beta})$?

By (25), what is needed is a central limit theorem for $N^{-1/2} \hat{\underline{X}}^T \hat{\underline{\Sigma}}^{-1} \hat{\underline{E}}$ which by (32) can be written as a sum of N independent zero mean random vectors. Moreover, by (31) the covariance matrix of this sum converges to A . So we need only add a Lindebergh-type condition ensuring that each term in the sum is asymptotically negligible in order to guarantee asymptotic normality of $\underline{\hat{\beta}}$. This condition is going to involve the conditional distributions of

\hat{e}^n given $x_P^n = x_P^n$, which could depend on P^n and x_P^n in a very complicated way. For simplicity we might assume them only to depend on P^n ; we have already assumed this for the conditional expectations and variances. However this is rather close to assuming multivariate normality of x^n as the following special case, $K = 2$, shows. The new assumption is equivalent to assuming that x_P^n and $e^n + \beta_M^T(x_M^n - \alpha_{MP}x_P^n)$ are independent for each P . Taking $P = \{1, 2\}$ and $M = \phi$, e^n and x^n are independent; taking $P = \{1\}$ and then $P = \{2\}$ we find that x_1^n is independent of $\beta_2(x_2^n - \alpha_{21}x_1^n)$ and x_2^n of $\beta_1(x_1^n - \alpha_{12}x_2^n)$. By the theorem of Skitovich (see Kagan, Linnik and Rao (1973) Theorem 3.1.1) it now follows that if all the coefficients involved here are non-zero then x^n is bivariate normally distributed.

Of course, if (x^n, e^n) is multivariate normally distributed, then conditional on $\hat{X} = \hat{X}$, $\hat{\beta}$ has the $\mathcal{N}(\beta, (\hat{X}^T \hat{\Sigma}^{-1} \hat{X})^{-1})$ distribution and the required results can be obtained very easily.

CONCLUSION

Though under reasonable conditions the estimators considered are consistent—it is not even necessary to assume the covariates are random—fairly strong conditions are needed to justify the use of $(\hat{X}^T \hat{\Sigma}^{-1} \hat{X})^{-1}$ as an approximate covariance matrix for $\hat{\beta}$: namely randomness of the covariates, independence between the N observations, and uncorrelatedness of $(\hat{e}^n)^2$ and $x_P^n x_P^{nT}$. It is worth pointing out that small sample simulation results in the literature are nearly always based on a multivariate normal distribution for (x^n, e^n) ; see for instance Little (1979).

REFERENCES

- AFIFI, A. A. and ELASHOFF R. M. (1966): Missing observations in multivariate statistics—I. Review of the literature. *J. Amer. Statist. Ass.*, **61**, 595-604.
- (1967): Missing observations in multivariate statistics—II. Point estimation in simple linear regression. *J. Amer. Statist. Ass.*, **62**, 10-29.
- (1969a): Missing observations in multivariate statistics—III. Large sample analysis of simple linear regression. *J. Amer. Statist. Ass.*, **64**, 337-358.
- (1969b): Missing observations in multivariate statistics—IV. A note on simple linear regression. *J. Amer. Statist. Ass.*, **64**, 359-365.
- BEALE, E. M. L. and LITTLE R. J. (1975): Missing values in multivariate analysis. *J. R. Statist. Soc.*, (B) **37**, 129-145.
- DAGENAIS, M. G. (1973): The use of incomplete observations in multiple regression analysis: A generalized least squares approach. *J. Econometrics*, **1**, 317-328.
- DEMPTER, A. P., N. M. LAIRD and D. B. RUBIN (1977): Maximum likelihood from incomplete data via the EM algorithm. *J. Roy. Statist. Soc.* (B), **39**, 1-38.

- GILL, R. D. (1977): Consistency of maximum likelihood estimators of the factor analysis model, when the observations are not multivariate normally distributed. *Recent Developments in Statistics*, J. R. Barra et al. (eds.), North-Holland, Amsterdam.
- HARTLEY, H. O. and R. R. HOCKING (1971): The analysis of incomplete data. *Biometrics*, **27**, 783-823.
- KAGAN, A. M., Y. V. LINNIK and C. R. RAO (1973): *Characterization Problems in Mathematical Statistics*, Wiley, New York.
- LITTLE, R. J. A. (1979): Maximum likelihood inference for multiple regression with missing values: a simulation study. *J. Roy Statist. Soc.*, (B) **41**, 76-87.

Paper received : March, 1982.

Revised : December, 1983.