



NATIONAL BUREAU OF STANDARDS

マン ちょう う

MRC Technical Summary Report #2836

CONDITIONING IN A MISSING DATA PROBLEM



Mathematics Research Center University of Wisconsin—Madison 610 Walnut Street Madison, Wisconsin 53705

C. J. Skinner

July 1985

(Received May 28, 1985)



DTIC FILE COPY

Sponsored by

U.S. Army Research Office P. O. Box 12211 Research Triangle Park North Carolina 27709 Approved for public release Distribution unlimited

> National Science Foundation Washington, DC 20550



UNIVERSITY OF WISCONSIN - MADISON MATHEMATICS RESEARCH CENTER

Accounter For MITS CUADE PETC TAB

Unnamoune. 3

Justificstee

Distribution/

Dist

A-1

DDC

QUALITY

Availabili, t 0 les

Avai Lingur

Special

CONDITIONING IN A MISSING DATA PROBLEM

C. J. Skinner

Technical Summary Report #2836

July 1985

ABSTRACT

 $4 \rightarrow$ Observations are recorded on variables x and y but a mechanism, which may depend on the observed x values, causes some of the y values to be missing. For three parametric examples, exact or approximate ancillary statistics are constructed. Conditioning on these ancillaries enables the missing data mechanism to be ignored under certain conditions. A correspondence is shown between these conditional procedures and the use of the observed information matrix in measuring the dispersion of the maximum likelihood estimator. Horgan = 1

AMS (MOS) Subject Classifications: 62A20, 62D05, 62F25.

Key Words: affine ancillary; ancillary statistic; conditional inference; curved exponential family; ignorability; information; missing data; survey sampling.

Work Unit Number 4 - Statistics and Probability

Sponsored by the United States Army under Contract No. DAAG29-80-C-0041. This material is based upon work supported by the National Science Foundation under Grant No. DMS-8210950, Mod. 1.

SIGNIFICANCE AND EXPLANATION

Many statistical problems can be viewed as missing data problems. Data may be missing for practical reasons, out of the control of the datacollector, or missing by design, such as in sample surveys where, for some variables, data is available for the whole population but, for other variables, data is recorded for the sample and is 'missing' for the remainder of the population.

Rubin (1976) has shown that the mechanism which causes the data to be missing (the sampling design in the survey context) can, for a wide class of situations, be ignored for Bayes or Likelihood inference but not for classical sampling distribution theory inference. This non-ignorability of the missing data mechanism can make classical inference much more complicated and even impossible if the mechanism is unknown.

In this paper we apply ideas of Barndorff-Nielsen (1980), for a class of statistical models called curved exponential families, to construct ancillary statistics for some missing data problems. We show that if classical inference is carried out conditional upon the observed values of these ancillary statistics then the missing data mechanism may be ignored, in certain situations. Some correspondence between these conditional procedures and likelihood methods is established. The approach rests heavily on examples. This is characteristic of the conditioning literature, where specific results are often more enlightening than attempts at generality.

The responsibility for the wording and views expressed in this descriptive summary lies with MRC, and not with the author of this report.

CONDITIONING IN A MISSING DATA PROBLEM

C. J. Skinner

1. INTRODUCTION

In this article we consider the possibility of conditioning, using either exact or approximate ancillary statistics, in a problem of estimation with missing data. The approach is not general but illustrative, using three examples.

Rubin (1976) showed that the mechanism which causes data to be missing may, in a large class of situations, be ignored for Bayesian or Likelihood inference but not for sampling distribution inference. We shall show how conditioning can enable this mechanism to be ignored for a wider class of situations under sampling distribution inference. Essentially we show, as in Efron and Hinkley (1978), that the conditional distribution of the maximum likelihood estimator (MLE) given an approprite ancillary corresponds, at least approximately, to the use of the inverse of the observed Fisher information matrix to measure the dispersion of the MLE. The latter procedure, being purely likelihood based, shares the ignorability properties of Likelihood inference. To provide initial motivation for the form of the conditioning procedure we compare our estimation problem with a prediction problem in survey sampling.

We assume the pairs (y_i, x_i) , i = 1, ..., N form a random sample from a bivariate distribution $p(y, x; \psi)$ belonging to a family indexed by the (generally vector) parameter ψ . We do not observe the complete data $d_C = (y_1, ..., y_N, x_1, ..., x_N)$ but only the incomplete data $d_I = (y_{i_1}, ..., y_{i_n}, s, x_1, ..., x_N)$ where $s = \{i_1, ..., i_n\}$ is a subset of size n from $U = \{1, ..., N\}$, $n \le N$. We assume that s, and hence d_I , is obtained from d_C by a selection mechanism which assigns probabilities $p(s|x_C)$, possibly dependent on $x_C = (x_1, ..., x_N)$ but not on ψ , to selecting each of the $\binom{N}{n}$ possible subsets s.

The above set-up occurs, for example, in survey sampling where a sample s is selected from a finite population U, an auxiliary variable x is known for each unit

Sponsored by the United States Army under Contract No. DAAG29-80-C-0041. This material is based upon work supported by the National Science Foundation under Grant No. DMS-8210950, Mod. 1.

i in U and is available for use in the selection of s and where y is a variable measured in the survey. The mechanism $p(s|x_C)$ is called the sample design in this context.

Consider two problems:

(A) the <u>prediction</u> of a function of (y_1, \dots, y_N) , specifically $\overline{y} = N^{-1} \sum_{i=1}^{N} y_i$, (B) the <u>estimation</u> of a parameter of the marginal distribution of y_i , specifically

 $\mu_{\mathbf{y}} = \mathbf{E}(\mathbf{y}).$

In the sample survey context, (A) is the traditional descriptive use, usually concerned …ith means and totals, whilst (B) is the analytical use, usually concerned with the estimation of underlying models such as regression models (see e.g. Hartley and Sielken, 1975). Our consideration of only \overline{y} and μ_y may be taken as special cases of this more general use.

Under the sampling distribution approach to inference, it is usual and natural in (A) to make predictive inference about \overline{y} conditional on (s, x_C) . Prediction proceeds as in the usual regression context where (x_i, y_i) , i ϵ s are known and x_i , i \neq s are new x values for which we wish to predict y. More formally, note that the conditional distribution of d_C given $d_{\overline{X}}$ depends only on the parameter ϕ which indexes the conditional distribution of y given x. Suppose we may write $\psi = (\phi, \lambda)$, where $p(y, x; \psi) = p(y|x; \phi)p(x; \lambda)$, so that the distribution of the data is

$$\mathbf{p}(\mathbf{q}^{\mathbf{i}}_{\mathbf{i}}, \mathbf{h}) = \mathbf{p}(\mathbf{x}^{\mathbf{i}}_{\mathbf{i}}, \dots, \mathbf{x}^{\mathbf{i}}_{\mathbf{i}}) = \mathbf{p}(\mathbf{x}^{\mathbf{i}}_{\mathbf{i}}, \dots, \mathbf{x}^{\mathbf{i}}_{\mathbf{i}}) = \mathbf{p}(\mathbf{x}^{\mathbf{i}}_{\mathbf{i}}, \dots, \mathbf{x}^{\mathbf{i}}_{\mathbf{i}}) = \mathbf{p}(\mathbf{x}^{\mathbf{i}}_{\mathbf{i}}, \mathbf{h}) = \mathbf{p}(\mathbf{x}^{\mathbf{$$

Then, provided φ and λ are variation independent, it may be argued (Cox and Hinkley, 1974, p. 35; Barndorff-Nielsen, 1978, p. 50) that (s,x_C) is (extended-;S-) ancillary for φ , treating λ as a nuisance parameter, and that inference about φ and hence the prediction of \overline{y} should be made conditional on (s,x_C) . (More precisely we condition on a minimal sufficient reduction of (s,x_C) , see Section 4.)

-2-

For the estimation of μ_y in (B), however, the same argument does not apply. Although (s,x_C) may be ancillary for φ , the parameter of interest μ_y is in general a function not only of φ but also of λ and will not be identifiable in the conditional distribution of d_I given (s,x_C) . Hence inference about μ_y conditional on (s,x_C) is generally inappropriate.

Now if N is large, and in the sample survey context N may be very large, the difference between \overline{y} and μ_y may be very small, even though it would appear from the discussion above that sampling distribution inference about these two quantities could proceed quite differently. This apparent 'paradox' provides one source of motivation for seeking an alternative procedure for conditional inference about μ_y .

E IN

In Bayesian or Likelihood inference about μ_y the selection mechanism only enters as a multiplicative factor free of ψ in the likelihood in (1) and hence is ignorable. In other words, making the false assumption that s is derived from U by simple random sampling rather than by the true mechanism $p(s|x_c)$, would make no difference to the inference for given d_I . In Rubin's (1976) terminology this is because we have assumed the data is 'missing at random'. The selection mechanism is not, however, ignorable for unconditional sampling distribution inference (nor for inference conditional on s but not on x_c as in Rubin, 1976).

In Section 2 we present the three examples. They are chosen to be of increasing order of complexity. In Example 1 there is an exact ancillary and the conditional distribution for inference about μ_y is exactly normal. In Example 2 there is an exact ancillary but the conditional distribution is only asymptotically normal. In Example 3 there is only an approximate ancillary.

In Section 3 we consider the prediction of \overline{y} conditional on (s, x_C) . In Section 4 we consider the estimation of μ_y on the assumption that s is a simple random sample from U. We adopt the approach of Barndorff-Nielsen (1980) to determine appropriate conditioning procedures. These are then compared with Section 3. In Section 5 we compare the conditional procedures with the use of asymptotic maximum likelihood theory using dispersion estimates based on either observed or expected Fisher information matrices. In

-3-

Section 6 we consider conditioning under a general selection mechanism $p(s|x_C)$ and discuss in what sense this mechanism is ignorable.

In asymptotic arguments we shall assume that n indexes fixed sequences $N = N_n$ and $p(s|x_c) = p_n(s|x_c)$ and that $n/N_n + f$, a constant, as $n + \infty$. Notation will be of three types. The key quantities ψ , φ , λ , μ_y , t, a are defined generally but take different forms in the different examples. The remaining statistics such as $\overline{x_s}$, \overline{x} , are defined in terms of d_I and are invariant with respect to the example. The remaining parameters, such as α , β and γ , are example-specific.

2. THE EXAMPLES

Example 1

We assume (y,x) is bivariate normal. The parameter vector is $\psi = (\psi,\lambda)$ where $\varphi = (\alpha,\beta,\sigma_{y*x}^2), \lambda = (\mu_x,\sigma_x^2)$ and $y|x \sim N(\alpha + \beta x, \sigma_{y*x}^2), x \sim N(\mu_x,\sigma_x^2)$. The parameter of interest is $\mu_y = \alpha + \beta \mu_x$. The MLE of ψ is (Anderson, 1957) $\hat{\psi} = (\overline{y}_x - b_x\overline{x}, b_y, n^{-1}SS_{y*xy}, \overline{x}, N^{-1}SS_x)$

where

$$\overline{\mathbf{y}}_{\mathbf{s}} = \mathbf{n}^{-1} \sum_{\mathbf{x}} \mathbf{y}_{\mathbf{i}}, \ \overline{\mathbf{x}}_{\mathbf{s}} = \mathbf{n}^{-1} \sum_{\mathbf{x}} \mathbf{x}_{\mathbf{i}}, \ \overline{\mathbf{x}} = \mathbf{N}^{-1} \sum_{\mathbf{x}} \mathbf{x}_{\mathbf{i}},$$
$$\mathbf{b}_{\mathbf{s}} = \sum_{\mathbf{s}} \mathbf{y}_{\mathbf{i}} (\mathbf{x}_{\mathbf{i}} - \overline{\mathbf{x}}_{\mathbf{s}}) / \mathbf{SS}_{\mathbf{x}\mathbf{s}}, \ \mathbf{SS}_{\mathbf{x}\mathbf{s}} = \sum_{\mathbf{s}} (\mathbf{x}_{\mathbf{i}} - \overline{\mathbf{x}}_{\mathbf{s}})^{2},$$
$$\mathbf{SS}_{\mathbf{y}^{*}\mathbf{x}\mathbf{s}} = \sum_{\mathbf{s}} (\mathbf{y}_{\mathbf{i}} - \hat{\alpha} - \hat{\beta} \mathbf{x}_{\mathbf{i}})^{2}, \ \mathbf{SS}_{\mathbf{x}} = \sum_{\mathbf{s}} (\mathbf{x}_{\mathbf{i}} - \overline{\mathbf{x}})^{2}.$$

Hence the MLE of μ_y is $\hat{\mu}_y = \hat{\alpha} + \hat{\beta} \cdot \overline{x} = \overline{y}_g + b_g(\overline{x} - \overline{x}_g)$ which in survey sampling is termed the regression estimator (Cochran, 1977). A minimal sufficient statistic for ψ is $t = (\overline{y}_g, \overline{x}_g, b_g, SS_{xg}, SS_{y*xg}, \overline{x}, SS_x)$. Neither $\hat{\psi}$ nor t depend on the selection mechanism because $p(s|x_c)$ is a multiplicative factor free of ψ in (1). The family of distributions for d_I indexed by ψ is a curved exponential family, labelled (7,5) by Barndorff-Nielsen (1980) since dim(t) = 7, dim(ψ) = 5.

Example 2

We assume that

 $y|x \sim N(\gamma x, \sigma^2 x)$, $x \sim Gamma(\lambda, k)$, k known

(i. $p(x_i\lambda) = \lambda^k x^{k-1} e^{-\lambda k} / \Gamma(k)$). The parameter vector is $\psi = (\varphi, \lambda)$ where $\varphi = (\gamma, \sigma^2)$. The parameter of interest is $\mu_{\gamma} = \gamma k / \lambda$. The MLE of ψ is

$$\hat{\psi} = (\overline{y}_{g}/\overline{x}_{g}, n^{-1}\sum_{g} (y_{i} - \hat{\gamma} x_{i})^{2}/x_{i}, k/\overline{x})$$

Hence the MLE of μ_y is $\mu_y = \overline{y_g x/x_g}$, the <u>ratio estimator</u> (Cochran, 1977). A minimal sufficient statistics for ψ is $t = (\overline{y_g}, \overline{x_g}, \sum_y y_i^2/x_i, \overline{x})$. The family of distributions for d_1 indexed by ψ is a (4,3) curved exponential family.

Example 3

We assume that y and x are both 0 - 1 variables with

$$Pr(y=1|x=0) = \phi_0, Pr(y=1|x=1) = \phi_1, Pr(x=1) = \lambda$$

The parameter vector is $\psi = (\varphi, \lambda)$ where $\varphi = (\varphi_0, \varphi_1)$. The parameter of interest is $\mu_y = Pr(y=1) = \lambda \varphi_1 + (1-\lambda)\varphi_0$. The MLE of ψ is $\hat{\psi} = (n_{10}/n_{.0}, n_{11}/n_{.1}, N_{.1}/N)$ where the cell counts in s and U are defined by

$$n_{\alpha\beta} = \sum_{s} y_{\underline{i}}^{\alpha} (1-y_{\underline{i}})^{1-\alpha} x_{\underline{i}}^{\beta} (1-x_{\underline{i}})^{1-\beta} \qquad \alpha, \beta = 0, 1$$
$$N_{\alpha\beta} = \sum_{1}^{N} y_{\underline{i}}^{\alpha} (1-y_{\underline{i}})^{1-\alpha} x_{\underline{i}}^{\beta} (1-x_{\underline{i}})^{1-\beta}$$

and the margins are defined by

 $n_{\boldsymbol{\alpha}\beta} = \sum_{\alpha} n_{\alpha\beta}, N_{\boldsymbol{\alpha}\beta} = \sum_{\alpha} N_{\alpha\beta} \qquad \beta = 0, 1 .$

The MLE of μ_v is thus

$$\hat{\mu}_{y} = \frac{N_{.0}}{N} \cdot \frac{n_{10}}{n_{.0}} + \frac{N_{.1}}{N} \cdot \frac{n_{11}}{n_{.1}}$$

a <u>poststratification estimator</u> (Cochran, 1977). A minimal sufficient statistic for ψ is t = (n₀₀, n₁₀, n₀₁, N_{.0}). This example also defines a (4,3) curved exponential family.

3. PREDICTION OF Y

A natural predictor of \overline{y} given d_{I} is the regression predictor:

$$\frac{1}{1} \begin{bmatrix} \sum y_{i} + \sum z(y|x = x_{i}) \end{bmatrix}$$

This predictor can be shown to be identical to $\hat{\mu}_y$ in each of our three examples. Under conditions obeyed by the examples, this predictor is the minimum variance unbiased predictor of \overline{y} conditional on (s, x_C) (Skinner, 1983). We now evaluate the conditional distribution of $\hat{\mu}_y - \overline{y}$ given (s, x_C) for each of the examples.

Example 1 (continued)

We obtain

$$\hat{\mu}_{y} = \overline{y}|_{s,x_{C}} \sim N(0, (1 - n/N + na_{1}^{2})\sigma_{y,x}^{2}/n)$$
 (2)

where $a_1 = (\bar{x}_g - \bar{x})SS_{xS}^{-\frac{1}{2}}$. Note that it is not necessary to condition on all the x_1 values and s but only on a_1 . An exact (1-a)-level conditional confidence interval for \bar{y} given a_1 is

$$\hat{\mu}_{y} \pm \left[(1 - n/N + na_{1}^{2}) SS_{y + x} / n(n-2) \right]^{\frac{1}{2}} t_{n-2} (\alpha/2)$$

where $t_{v}(\alpha)$ is the α^{th} point of Student's t-distribution with v d.f. Note that this interval does not depend on the selection mechanism. It distinguishes between 'good' samples where a_1 is small and hence \overline{y} may be more precisely predicted and 'bad' samples where a_1 is large and hence \overline{y} is more poorly predicted. Such a distinction is an essential aim of conditioning.

Example 2 (continued)

We obtain

$$\hat{\mu}_{y} = \overline{y} | s, x_{C} \sim N(0, [\overline{x}(N\overline{x} - n\overline{x}_{g})/N\overline{x}_{g}]\sigma^{2}/n) .$$
(3)

An exact $(1-\alpha)$ -level confidence level for \overline{y} conditional on $(\overline{x}, \overline{x})$ is

$$\hat{u}_{y} \pm \left\{ \left[\overline{x} (N\overline{x} - n\overline{x}_{g}) / N\overline{x}_{g} \right] \overline{\sigma}^{2} / (n-1) \right\}^{\frac{1}{2}} t_{n-1} (\alpha/2) \quad .$$

In this example a 'good' sample is one in which y is missing for small x values so that $\overline{x_s}$ is large.

Example 3 (continued)

We may write $\hat{\mu}_{y} - \overline{y} = N^{-1} (N_{.0}/n_{.0} - 1)n_{10} + N^{-1} (N_{.1}/n_{.1} - 1)n_{11} - N^{-1} \tilde{n}_{10} - N^{-1} \tilde{n}_{11}$ where $\tilde{n}_{\alpha\beta} = N_{\alpha\beta} - n_{\alpha\beta} - \alpha_{\alpha\beta} = 0, 1.$

Conditional on (s,x_{C}) the quantities $N_{.0}$, $N_{.1}$, $n_{.0}$ and $n_{.1}$ are fixed and the distribution of $\hat{\mu}_{y} - \overline{y}$ is determined by the independent Binomial distributions of n_{10} , n_{11} , \tilde{n}_{10} and \tilde{n}_{11} which possess parameters $(n_{.0},\varphi_{0})$, $(n_{.1},\varphi_{1})$, $(N_{.0} - n_{.0},\varphi_{0})$ and $(N_{.1} - n_{.1},\varphi_{1})$ respectively. An exact conditional confidence interval for \overline{y} appears intractable. Instead, suppose that $0 < \lambda$, φ_{0} , $\varphi_{1} < 1$ and that the sequence of designs is such that $n_{.0}/n$ is almost surely bounded away from 0 and 1 as $n + \infty$. Then almost surely

$$\frac{1}{2} \left(\hat{\mu}_{y} - \bar{y} \right) \left| s_{r} x_{c} \right|^{\frac{1}{2}} N(0, (1 - fa_{1})(1 - \hat{\lambda}) \varphi_{0}(1 - \varphi_{0}) / a_{1} + (1 - fa_{0}) \hat{\lambda} \varphi_{a}(1 - \varphi_{a}) / a_{n}$$
(4)

where $f = \lim(n/N)$, $a_1 = n_0 N/(nN_0)$, $a_2 = n_1 N/(nN_1)$.

A large sample conditional confidence interval may therefore be obtained substituting , φ for φ .

4. ESTIMATION OF $\mu_{\rm u}$ - MISSING VALUES SELECTED BY SIMPLE RANDOM SAMPLING

In general μ_y is not identified in the conditional distribution of d_I given (s,x_C) . For example, the transformation $\mu_y + \mu_y + \delta$, $\mu_x + \mu_x + \delta\beta^{-1}$, φ fixed leaves $p(d_I|s,x_C;\psi)$ unaffected in Example 1. Hence conditioning on (s,x_C) as in Section 3 is inappropriate. Instead we seek what Barndorff-Nielsen (1980) term a 'conditionality resolution', that is we seek an exact or approximate ancillary statistic a such that $(\hat{\psi},a)$ is a one-to-one transformation of the minimal sufficient statistic t. The conditional distribution of $\hat{\psi}$ given a is then the appropriate one for inference. For a (k,d) curved exponential family we seek a (k-d)-dimensional ancillary. Whilst $\hat{\psi}$ and t are unaffected by the missing data mechanism, the ancillarity of any statistic a

-7-

may be affected. Hence, in this section we suppose that s is obtained from U by simple random sampling of fixed size so that $\{x_i, i \approx 1, \dots, N\}$ are IID whether or not $i \in s$. We consider the general mechanism in Section 5.

Example 1 (continued)

For this (7,5) curved exponential family we seek a two-dimensional ancillary. Since the normal family is a location-scale tranformation family the following function of t:

$$a = (a_1, a_2) = [(\overline{x} - \overline{x})/SS_{xs}^{1/2}, SS_{xs}/SS_{x}]$$

is an exact ancillary, in the sense that its distribution is free of ψ . Note that the transformation t to $(\hat{\psi}, a)$ is one-one.

The distribution of ψ given a is in fact tractable although we shall only be concerned with the distribution of μ_y given a. Corresponding to (2) we have

$$\hat{\mu}_{y}|s_{r}x_{C} \sim N(\alpha + \beta \overline{x}_{r} (1 + na_{1}^{2})\sigma_{y \cdot x}^{2}/n) \qquad (5)$$

But \overline{x} is independent of a (since for example a is a function of $x_i - x_1$, i = 2,...,N) and so

$$\overline{\mathbf{x}}|\mathbf{a} \sim N(\boldsymbol{\mu}_{\mathbf{x}}, \sigma_{\mathbf{x}}^2/N) \quad . \tag{6}$$

Hence

$$\hat{\mu}_{y}|a \sim N(\mu_{y}, (1 + na_{1}^{2})\sigma_{y*x}^{2}/n + \beta^{2}\sigma_{x}^{2}/N) .$$
(7)

This provides an appropriate conditional distribution for inference about μ_y . In fact we only need to condition on a_1 and from (2) this permits us to apply the same level of conditioning in the estimation of μ_y as in the prediction of \overline{y} . In the sense of Section 1 we have therefore resolved the 'paradox' of different levels of conditioning for the estimation and prediction problems.

It only appears possible, however, to obtain a large-sample (rather than exact) $\hat{\rho}_{\nu}$ conditional interval for μ_{ν} , since no pivot based on $\hat{\mu}_{\nu}$ = μ_{ν} seems to exist.

Example 2 (continued)

For this (4,3) curved exponential family we seek a one-dimensional ancillary. Since λ is a scale parameter, an exact ancillary is given by $a = \overline{x/x_g}$. The transformation $t + (\hat{\psi}_r a)$ is one-one. It may be shown that a is independent of \overline{x} . Hence the exact

distribution of μ_y given a is obtained by integrating \overline{x} out of the joint distribution of $(\hat{\mu_y}, \overline{x})$ defined by

$$\hat{\mu}_{y} \left[\bar{x}_{r} a \sim N(\gamma \bar{x}_{r} \bar{x} a \sigma^{2} / n) \right]$$
(8)

$$\mathbf{x} = \mathbf{a} \sim \text{Gamma}(\mathbf{N}\lambda, \mathbf{N}\mathbf{k})$$
 (9)

As n + ∞

$$n^{\frac{1}{2}} (\hat{\mu}_{y} - \mu_{y}) |a^{\frac{1}{2}} N(0, ka\sigma^{2}/\lambda + k\gamma^{2} f/\lambda^{2}) .$$
 (10)

The 'paradox' of different levels of conditioning for \overline{y} and μ_y is only partially resolved here, since in (3) inference about \overline{y} is made conditional not only on a but also on \overline{x} , whereas from (8) inference about μ_y given both a and \overline{x} is not possible.

Example 3 (continued)

We again have a (4,3) curved exponential family and seek a 1-dimensional ancillary. In this example no exact ancillary, which is a function of the minimal sufficient statistic, appears to exist. Barndorff-Nielsen (1980) discusses the construction of an approximate ancillary, termed the affine ancillary, for a general curved exponential family. In this 1-dimensional case the affine ancillary is

$$\mathbf{a_{aff}} = \mathbf{k}(\hat{\psi}) [\mathbf{t} - \tau(\hat{\psi})] \tau_{\parallel}^{\mathrm{T}}(\hat{\psi})$$
(11)

where $\tau(\psi) = E(t), \tau_{\perp}(\psi)$ is a 1 × 4 vector which is orthogonal to the rows of the 3 × 4 matrix $\partial \tau(\psi)/\partial \psi$ and $k(\psi)$ is a scalar chosen such that

$$\operatorname{var}\{\mathbf{k}(\psi) [\mathbf{t} - \tau(\psi)] \tau_{j}^{\mathbf{T}}(\psi)\} = 1$$

In our example $\tau_{i}(\psi) = (1 \ 1 \ 0 \ -n/N)$, up to a scalar multiple, which implies that

$$a_{aff} = k(\lambda)(n_0/n - N_0/N)$$

where

$$k(\lambda) = [(N-n)\lambda(1-\lambda)/nN]^{-\frac{1}{2}}$$

Now let

$$T_{1} = k(\lambda)(1-\lambda)(a_{1}-1) , T_{2} = k(\lambda)\lambda(1-a_{2})$$

$$a_{1} = n_{0}N/(nN_{0}) , a_{2} = n_{1}N/(nN_{1}) .$$

Then T_1 , T_2 and a_{aff} all converge to the same (standard normal) random variable as n + ∞ . Hence in the limit as n + ∞ , conditioning on a_{aff} is equivalent to the first

-9-

order to conditioning on (a_1, a_2) . Also in the limit as $n + \infty$, N_{0}/N is independent of a_{aff} . Thus

$$n^{1/2} (N_{0} / N - \lambda) |a_{1} a_{2} \stackrel{L}{\neq} N[0, f\lambda(1-\lambda)]$$
(12)

and as in (4), almost surely

$$\frac{1}{2} \left[\hat{\mu}_{y} - (N_{\cdot 0} \phi_{0} + N_{\cdot 1} \phi_{1}) / N \right] \left| a_{1} r a_{2} r N_{\cdot 0} / N \right|^{\frac{1}{4}}$$

$$N \left[0, \phi_{0} (1 - \phi_{0}) N_{\cdot 0} / N a_{1} + \phi_{1} (1 - \phi_{1}) (1 - N_{\cdot 0} / N) / a_{1} \right] .$$

$$(13)$$

Hence, almost surely

Martine 1

$$n^{1/2} (\hat{\mu}_{y} - \mu_{y}) |_{a_{aff}} \stackrel{I_{+}}{\to} N[0, (1 - \lambda) \phi_{0}(1 - \phi_{0})/a_{1} + \lambda \phi_{1}(1 - \phi_{1})/a_{2} + (\phi_{0} - \phi_{1})^{2} \lambda (1 - \lambda) f] .$$
(14)

Note that the level of conditioning is again less than for the prediction of \overline{y} , since in (4) the conditioning is not only on a_1 and a_2 but also on $N_{.0}/N = 1-\hat{\lambda}$.

5. OBSERVED VERSUS EXPECTED FISHER INFORMATION

The asymptotic theory of maximum likelihood estimation for the regular IID case may be extended to the incomplete data structure of d_{I} (c.f. Hocking and Smith, 1968). According to this theory, two estimates of the (asymptotic) covariance matrix of $\hat{\psi}$ are $j(\hat{\psi})^{-1}$ and $i(\hat{\psi})^{-1}$ where $j(\hat{\psi})$ and $i(\hat{\psi})$ are the observed and expected Fisher information matrices respectively:

 $i(\psi) = E[j(\psi)]$, $j(\psi) = -\partial^2 \log p(d_1/\psi)/\partial\psi \partial\psi^T$.

Efron and Hinkley (1978) show that $j(\hat{\psi})^{-1}$ is often a good approximation to the conditional variance of $\hat{\psi}$ given an appropriate ancillary, a. Barndorff-Nielsen (1980) considers the extension to the multi-parameter case. We now compare $var(\hat{\mu}_y|a)$, as obtained from Section 4, with the estimates derived from $i(\hat{\psi})$ and $j(\hat{\psi})$, which are

$$v_{obs}(\hat{\mu}_{y}) = g^{*}(\hat{\psi})^{T}j(\hat{\psi})^{-1}g^{*}(\hat{\psi})$$

$$v_{exp}(\hat{\mu}_{y}) = g^{*}(\hat{\psi})^{T}j(\hat{\psi})^{-1}g^{*}(\hat{\psi})$$
where $\mu_{y} = g(\psi)$, $g^{*}(\psi) = \partial g(\psi)/\partial \psi$

-10-

Example 1 (continued)

We obtain

$$v_{obs}(\hat{\mu}_{y}) = (1 + na_{1}^{2})\hat{\sigma}_{y*x}^{2}/n + \hat{\beta}^{2}\hat{\sigma}_{x}^{2}/N$$
$$v_{exp}(\hat{\mu}_{y}) = \hat{\sigma}_{y*x}^{2}/n + \hat{\beta}^{2}\hat{\sigma}_{x}^{2}/N$$

and comparing with (7), v_{obs} is identical to $var(\hat{\mu}_y|a)$ evaluated at $\psi = \hat{\psi}$.

Example 2 (continued)

We obtain

$$v_{obs}(\hat{\mu}_y) \approx ka \hat{\sigma}^2 / \hat{\lambda} n + k \hat{\gamma}^2 / \hat{\lambda}^2 N$$

 $v_{exp}(\hat{\mu}_y) \approx k \sigma^2 / \hat{\lambda} n + k \hat{\gamma}^2 / \hat{\lambda}^2 N$

and it follows immediately from (8) and (9) that v_{obs} is identical to $var(\hat{\mu}_y|a)$ evaluated at $\psi = \hat{\psi}$.

Example 3 (continued)

We obtain

$$v_{obs}(\hat{\mu}_{y}) = (1-\hat{\lambda})\hat{\phi}_{0}(1-\hat{\phi}_{0})/na_{1} + \hat{\lambda}\hat{\phi}_{1}(1-\hat{\phi}_{1})/na_{2} + (\hat{\phi}_{1}-\hat{\phi}_{0})^{2}\hat{\lambda}(1-\hat{\lambda})/N$$
$$v_{evn}(\hat{\mu}_{y}) = (1-\hat{\lambda})\hat{\phi}_{0}(1-\hat{\phi}_{0})/n + \hat{\lambda}\hat{\phi}_{1}(1-\hat{\phi}_{1})/n + (\hat{\phi}_{1}-\hat{\phi}_{0})^{2}\hat{\lambda}(1-\hat{\lambda})/N$$

and comparing with (14), nv_{obs} is identical to $\lim var[n^{1/2}(\hat{\mu}_y - \mu_y)]a_{aff}]$ with ψ and f replaced by $\hat{\psi}$ and n/N respectively.

6. ESTIMATION OF μ_y -General MISSING DATA MECHANISM

The mechanism $p(s|x_{C})$ does not affect $\hat{\psi}$ or t but may affect the ancillarity of the statistics a discussed in Section 4. We shall first show that the exact ancillaries given for Example 1 and 2 remain exactly ancillary for a broad class of mechanisms and that the distribution $\hat{\mu}_{v}|a$ is also invariant. We then consider the more general problem.

-11-

Example 1 (continued)

<u>Condition C1</u>: the selection mechanism depends only on the configuration

 $z = [(x_3-x_1)/(x_2-x_1), (x_4-x_1)/(x_2-x_1), \dots, (x_N-x_1)/(x_2-x_1)] \text{ of } x_C.$

In other words, under C1, $p(s|x_C)$ is invariant with respect to location or scale changes in x. This condition holds for a variety of sample designs, for example in stratified random sampling when strata are determined by quantiles of x_C and in truncated sampling where the point(s) of truncation are quantiles of x_C .

Under C1, a remains ancillary. For a is a function of z and s. The distributions $p(s|z) = p(s|x_c)$ and p(z) are free of ψ . Hence the distribution of a is free of ψ . Furthermore \overline{x} is independent of z and is conditionally independent of s given z. Hence \overline{x} is independent of (s,z) and therefore of a. Thus (5) and (6) still apply and the distribution $\hat{\mu}_{y}|a$ is again given by (7).

Example 2 (continued)

<u>Condition C2</u>: the selection mechanism depends only on the ratios $w = (x_2/x_1, x_3/x_1, \dots, x_N/x_1).$

In other words, under C2, $p(s|x_C)$ is invariant with respect to scale changes of . x. This condition holds, for example, with probability proportional to size designs, where x is a size measure, which are often used in conjunction with the ratio estimator.

Under C2, a remains ancillary by a similar argument to that in Example 1. Also it may be shown that \overline{x} is independent of (s, w) and therefore of a and hence that the expressions (8), (9) and (10) remain valid.

We now turn to the general situation where the statistics a defined in Section 4 need no longer be ancillary, even approximately. We begin by attempting to construct ancillaries which may now depend on the mechanism $p(s|x_c)$. We could consider the affine ancillary but a slightly modified approach simplifies the distribution theory for $\hat{u}_y|a$. Let $\hat{\lambda}_s$ be the MLE of λ were only x_i , i $\in s$, to be observed, so for our three examples $\hat{\lambda}_s = (\bar{x}_s, n^{-1}SS_{xs})$, k/\bar{x}_s and $n_{,1}/n$ respectively. Let $T(\hat{\lambda}) = E(\hat{\lambda}_s|\hat{\lambda})$. In each of our examples $\hat{\lambda}$ is sufficient for λ_r and hence $T(\hat{\lambda})$, although critically dependent on

-12-

 $p(s|x_C)$, is free of λ . We assume that the sequence of designs is such that $T(\cdot)$ is continuous at λ and that

$$n^{\frac{1}{2}} [\hat{\lambda}_{s} - \tau(\hat{\lambda})] \stackrel{\downarrow}{=} N[0, k(\lambda)]$$

where $k(\lambda)$ is a finite positive-definite matrix. This seems a fairly weak condition although, for example, it would exclude the stratified design in Example 3 where $n_{.1}$ is set at a fixed fraction of $N_{.1}$ so that $k(\lambda) = 0$. We adopt as an approximate ancillary $a = n^{1/2} (\hat{\lambda}_{a} - r(\hat{\lambda})) k(\hat{\lambda})^{-1/2}$. (15)

Note that, when the missing values are randomly selected, a reduces to the ancillary in Section 4 for Example 3 and is asymptotically equivalent to the ancillaries for Example 1 and 2. Since $\cos(\hat{\lambda}_g - T(\hat{\lambda}), \hat{\lambda}) = 0$, $n^{\sqrt{2}}(\hat{\lambda} - \lambda)$ will be asymptotically independent of a. Now the conditional distribution of $\hat{\mu}_g$ given $(\mathbf{s}, \mathbf{x}_C)$ is unaffected by selection and in each of our examples we may write (almost surely)

$$n^{\frac{1}{2}} [\hat{\mu}_{y} - f_{0}(\hat{\lambda}, \varphi)] | \mathbf{s}, \mathbf{x}_{C} \stackrel{\text{L}}{=} \mathbb{N}[0, f_{1}(\hat{\lambda}, \hat{\lambda}_{g}, \varphi)]$$
(16)

where f_0 and f_1 are certain functions, continuous at $\hat{\lambda} = \lambda$, such that $f_0(\lambda, \varphi) = \mu_y$. Now define f_2 such that $f_1(\hat{\lambda}, \hat{\lambda}_g, \varphi) = f_2(\hat{\lambda}, a, \varphi)$. Then

$$\begin{split} & n^{\frac{1}{2}} \left[\hat{\mu}_{y} - f_{0}(\hat{\lambda}, \varphi) \right] | \mathbf{a}, \hat{\lambda} \stackrel{\text{L}}{=} \mathbb{N}[0, f_{2}(\hat{\lambda}, \mathbf{a}, \varphi)] \\ & n^{\frac{1}{2}} (\hat{\lambda} - \lambda) | \mathbf{a} \stackrel{\text{L}}{=} \mathbb{N}[0, f_{3}(\lambda)] \end{split}$$

so that $n^{1/2}(\hat{\mu}_y - \mu_y) | a \stackrel{I}{+} N[0, f_2(\lambda, a, \varphi) + f_4(\psi)]$ where $f_4(\psi) = \lim_{n \to \infty} n \operatorname{var}[f_0(\hat{\lambda}, \varphi)]$.

The estimated asymptotic variance of $n^{\frac{1}{2}}(\hat{\mu}_{y} - \mu_{y})$ given a is thus $\hat{v} = f_{2}(\hat{\lambda}, a, \hat{\varphi}) + f_{4}(\hat{\psi})$.

But from (16) $f_4(\cdot)$ and $f_2(\hat{\lambda}, a, \hat{\varphi}) = f_1(\hat{\lambda}, \hat{\lambda}_g, \hat{\varphi})$ are unaffected by the missing data mechanism and hence this mechanism is ignorable for conditional inference. In other words, if we supposed incorrectly that the missing values were randomly selected so that x_i , i \in s is distributed identically to x_i , i \notin s we would obtain the same \hat{V} .

For a given mechanism the quantity $T(\hat{\lambda})$ in (15) and hence the ancillary a may be very complicated to compute. In practice, however, this is unnecessary. We showed in

-13-

Section 4 that $n^{-1}\hat{\nabla}$ was equal to $v_{obs}(\hat{\mu}_y)$ (with n/N replaced by f) for our examples. Hence for inferential purposes we only require the straightforward computation of $\hat{\mu}_y$ and $v_{obs}(\hat{\mu}_y)$ which do not depend on $p(s|x_C)$.

Note in contrast that the estimation of the unconditional variance of $\sqrt{n}(\hat{\mu}_y \neg \mu_y)$ or the evaluation of $\nabla_{exp}(\hat{\mu}_y)$ does depend on $p(\mathbf{s}|\mathbf{x}_C)$ and can be quite intractable. This provides a further practical advantage of conditioning.

7. CONCLUSION

We have indicated how exact or approximate conditioning arguments may be applied in three examples of a missing data problem. Conditioning is attractive here for several reasons: (1) it can permit the mechanism which causes the data to be missing to be ignored, (2) it can lead to more tractable procedures, (3) it makes inference more 'datadependent' (Fisher's original motivation).

The results of this article are specific to the examples chosen although some possible generalization is suggested in Section 6 for models where t is a 1 - 1 transformation of $(\hat{\psi}, \hat{\lambda}_g)$. In this case the asymptotic maximum likelihood approach using the observed information matrix corresponds, under certain conditions, to conditioning on the ancillary defined in (15).

Acknowledgements

I am grateful to D. R. Cox for fundamental suggestions.

REFERENCES

ANDERSON, T. W. (1957). Maximum likelihood estimates for a multivariate normal

distribution when some observations are missing. J. Am. Statist. Assoc. 52, 200-3. BARNDORFF-NIELSEN, O. (1978). Information and Exponential Families. Chichester: Wiley. BARNDORFF-NIELSEN, O. (1980). Conditionality resolutions. <u>Biometrika 67, 293-310.</u> COCHRAN, W. G. (1977). <u>Sampling Techniques</u> (3rd Edition). New York: Wiley. COX, D. R. & HINKLEY, D. V. (1974). <u>Theoretical Statistics</u>. London: Chapman and Hall. EFRON, B. & HINKLEY, D. V. (1978). Assessing the accuracy of the maximum likelihood

estimator: Observed versus expected Fisher information. <u>Biometrika 65</u>, 457-82. HARTLEY, H. O. & SIELKEN, R. L. (1975). A "super-population viewpoint" for finite

population sampling. <u>Biometrics 31</u>, 411-22.

HOCKING, R. R. & SMITH, W. B. (1968). Estimation of parameters in the multivariate normal distribution with missing observations. J. Am. Statist. Assoc. 63, 159-173.
RUBIN, D. B. (1976). Inference and missing data. <u>Biometrika 63</u>, 581-92.
SKINNER, C. L. [1983). Multivariate prediction from selected samples. <u>Biometrika 70</u>, 289-92.

CJS/jvs

1. REPORT NUMBER 2. GOVY ACCESSION #2836 2. GOVY ACCESSION 4. TITLE (and Sublitie) 2. GOVY ACCESSION Conditioning in a Missing Data Problem 2. GOVY ACCESSION	
#2836 • TITLE (and Sublitio) Conditioning in a Missing Data Problem	NO. 3. RECIPIENT'S CATALOG NUMBER
Conditioning in a Missing Data Problem	
Conditioning in a Missing Data Problem	S. TYPE OF REPORT & PERIOD COVERED
Conditioning in a Missing Data Problem	Summary Report - no specific
	reporting period
	5. PERFORMING ORG. REPORT NUMBER
. AUTHOR(0)	8. CONTRACT OR GRANT NUMBER(a)
	DMS-8210950, Mod. 1
C. J. Skinner	DAAG29-80-C-0041
PERFORMING ORGANIZATION NAME AND ADDRESS	10. PROGRAM ELEMENT, PROJECT, TASK
Mathematics Research Center, University of	Work Unit Number 4 -
610 Walnut Street Wisconsin	n Statistics & Probability
Madison, Wisconsin 53705	
1. CONTROLLING OFFICE NAME AND ADDRESS	12. REPORT DATE
See Item 18 holds	July 1985
DEE TCEN TO DETOM	13. NUMBER OF PAGES
4. MONITORING AGENCY NAME & ADDRESS(II different from Controlling Offic	e) 15. SECURITY CLASS. (of this report)
	IINCLASSIFIED
	SCHEDULE
7. DISTRIBUTION STATEMENT (of the obstract onlored in Block 20, it differen	t (rom Report)
18. SUPPLEMENTARY NOTES U. S. Army Research Office	lational Science Foundation
IS. SUPPLEMENTARY NOTES U. S. Army Research Office N P. O. Box 12211 N	Mational Science Foundation Mashington, DC 20550
B. SUPPLEMENTARY NOTES U. S. Army Research Office N P. O. Box 12211 W Research Triangle Park	Mational Science Foundation Mashington, DC 20550

DD 1 JAN 73 1473 EDITION OF 1 NOV 65 IS OBSOLETE

へい

IN.

<u>.</u>

UNCLASSIFIED SECURITY CLASSIFICATION OF THIS PAGE (When Date Entered)



FILMED

9-85

DTIC