Regression Methods and Performance Criteria for Small-Area Population Estimation

> P. McCullagh James V. Zidek

Technical Report No. 06

February 1985

REGRESSION METHODS AND PERFORMANCE CRITERIA

FOR SMALL-AREA POPULATION ESTIMATION

by

P. McCullagh Imperial College J. V. Zidek University of British Columbia

January 1983

ABSTRACT

In this paper a regression model is developed for the post-censal estimation of the population sizes of small-areas. The approach is nonstochastic. It is assumed that current population sizes have been determined by a function of those obtained at the last census, together with the associated values of certain symptomatic variables. As well, the current values of these variables are assumed to be in-hand. Natural properties of such a rule are shown to imply a specific, log linear form for this function. Existing models are shown to be approximations to the result.

An objective function for evaluating arbitrary estimation procedures or fitting regression models, is derived on the assumption that revenue allocation is the objective of the estimation program. By appealing to an equilibrium theory for group decision processes, it is shown that under certain conditions the appropriate criterion is given by the Kullback-Leibler discrimination function.

KEYWORDS: Ratio-correlation; difference-correlation; populations of small-areas; local-areas; log-linear models; Kullback-Leibler information; Nash solutions.

ACKNOWLEDGEMENTS

Interest in the problem considered in this paper was stimulated by a problem posed to the second author by the City of Vancouver and by discussions with Mr. Phil Mondor of its Planning Office and Professor H. Hightower, University of British Columbia.

The City of Vancouver supported related research, carried out by the second author and reported elsewhere.

The University of British Columbia and Imperial College, where the work reported herein commenced and finished, respectively, both very generously provided facilities during the period of the leaves of absence of the first and second authors, respectively.

Support for the work was provided by the Natural Sciences and Engineering and the Social Sciences and Humanities Research Councils of Canada, the latter through a Leave Fellowship to the second author.

INTRODUCTION

Government and private agencies depend on estimates of the populations of small-areas for a variety of purposes, such as revenue allocation and planning. In the United States, for example, annual, post-censal estimates are prepared under the Federal-State Co-operative Program (FSCP) for something like 39,000 municipalities, counties and so on (Kitagawa et al. 1980).

The great scale on which this activity is carried out makes the need for accurate, simply applied procedures acute. A variety of methods exist (Purcell and Kish 1979, Kitagawa et al. 1980, Zidek 1982). None are simpler and more adaptable than those which use regression models for making these estimates (Schmitt and Crosetti 1954, Morrison and Relles 1975, O'Hare 1976). These models turn easy-to-measure quantities (symptomatic variables) into estimates of population sizes which are expensive to measure. The coefficients obtained by fitting these models using two successive censuses and the measured values of all variables, one set per sub-area, (implicitly) account for migration, demographic trends and so on, so that the latter are not needed in the preparation of the estimates as they would be in, say, administrative records methods. Different symptomatic variables may be used in different sub-areas, depending on what data is readily available in each. And comparative empirical studies in the works cited above, show that these methods can be very effective, surprisingly so since the regression methodology is misappropriated in the sense that there is no conceivable experiment which generates these variable-values at random as in the conventional case. One such method, that using the ratiocorrelation model, is among the three, regression and other, which constitute the overall strategy used in the FSCP (Kitagawa et al. 1980).

In spite of the relatively high precision of regression-based methods, their simplicitly and adaptability, there have been surprisingly few models proposed and these have not been very much refined. Rosenberg (1968) points out the need to stratify local-areas by the dichotomies, first of urban versus rural and second of rapid versus slow growth. His proposal seems to have been ignored. Namboodiri (1972) argues pursuasively against the needed temporal stationarity of regression models and gives an analysis which may well show the ill-effects of multicollinearity among the symptomatic variables. A systematic residual analysis has not been published and little seems to be known about the possibly serious negative impact of influential observations in this context (see Belsey, Kuh and Welsch (1980) for a general discussion of this problem.

Existing regression models have been proposed on an ad hor basis. The major result of Section 3 is a new model which is obtained by what might be called an axionatic approach. Stochastic and approximation errors aside, it is supposed that population sizes of small-areas are assigned, hypothetically, by some unspecified function of the symptomatic variables and previous counts. The most reasonable requirements of such a rule are specified and then shown to be equivalent to a loglinear model similar but not identical to that of Morrison and Relies (1976). The model has yet to be empirically assessed.

Both the fitting of regression models and the assessment of estimation methods requires a criterion function which accumulates in some sensible way the errors made over all sub-areas. No particular choice has yet been indicated and various alternatives such as average relative absolute error are used, sometimes several in the same study. None of these criteria seem directly related to a primary objective of the program, namely, the allocation of revenue. In Section 2 is derived from first principles, a criterion function which takes account of the need to choose an allocation scheme which would be a jointly acceptable compromise to all members of the community, provided their desires obey certain weak constraints. The derivation is based on the theory of Nash (1950).

2. CRITERIA

Evaluating the performance of an estimation methodology requires answers to two fundamental questions. Against what are the answers it produces to be compared and by what criterion is the comparison to be made? This section is addressed to the second of these questions. We have no alternative to propose to the practical answer which is commonly given to the first question, namely, the corresponding census counts. The latter typically underestimate the true counts by something like 2% (Hauser 1981) and therefore seem somewhat unsatisfactory.

Kitagawa et al.(1980) describe various estimation performance criteria in terms of P., the "actual" and P., the estimated population sizes for

4.

subregion $\xi = I$, ..., n. The former would be the last available census counts corrected, possibly, as Spencer suggests (Appendix I, Kitagawa et al. 1980) for undercoverage. These criteria include average error, average relative error, number of extremely large relative errors and bias. The first two of these are obtained by putting a = I, b = 0 and a = I, b = J, respectively, in the general index $\sum_{i} |P_{ij} - \hat{P}_{ij}|^{a}/P_{ij}^{b}$.

While Kitagawa et al.(1980) refer to the need to take account of the purposes of the estimation program, no criterion has been given which does so, other than that which is derived below.

Its derivation relies on Nash's theory of bargaining as it might be applied in the present situation (Nash 1950). Let P and A denote, respectively, the population size and amount (of revenue, say) to be allocated. Let $\underline{A} = (A_1, \dots, A_n)$ be a feasible allocation of A among the n subregions and $u_{ij}(\underline{A}), j = 1, \dots, P_{ij}, i = 1, \dots, n$ represent the gain-in-value (utility) to individual j in subregion i which would result from this allocation scheme. Certain weak assumptions imply that any equilibrium solution must maximize an objective function which will now be described.

Before doing so, it should be pointed out that Nash's theory admits as potential solutions, not only the allocations, <u>A</u>, themselves but as well, all randomized mixtures of the <u>A</u>'s. Deadlocks can therefore be broken by tossing a coin, as it were. The domains of the u_{ij} 's are extended to this more general solution set by invoking the expected utility hypothesis. However, the feasible solutions, randomized and nonrandomized alike, are required to satisfy $u_{ij} \ge 0$. In this way, the Nash theory ensures that no individual can be made to suffer a net expected loss of utility as a result of the proposed allocation.

In agreement with practice, we will restrict ourselves to nonrandomized allocations and define as optimal any that are feasible and maximize the so-called Nash product given by

$$MP(\underline{A}) = \prod_{\substack{i \ j}} \prod_{\substack{i \ j}} u_{\underline{i}j}^{I/P}(\underline{A})$$
(2.1)

If the $u_{i,j}$'s are unknown, as would be the case in the situation under consideration, an approximation to (2.1) becomes necessary. If the A.'s are moderate, it is reasonable that $u_{i,j}(A)$ would equal, approximately, $A_{\underline{i}}/P_{\underline{i}}$ for all \underline{i} and \underline{j} . This approximation is supported by the assumptions that the \underline{i} -th subregion's allocation is evenly distributed among its constituents, that they receive no benefit from the allocations made to other subregions and that the utilities are linear. This approximation yields

$$MP(A) = \Pi (A_{2}/P_{2})^{P_{2}}$$
 (2.2)

where $p_{\vec{k}} = P_{\vec{k}}/P$ is the proportion of the population in subregion i for all i.

Equation (2.2) may be reduced further by letting $A_{\vec{n}} = a_{\vec{n}}A_{\vec{n}}$ $\vec{n} = 1, ..., n$. Then

$$NP(\underline{A}) = (\underline{A}/\underline{P}) \exp \left[-I(\underline{p},\underline{a})\right]$$
 (2.3)

where $I(p, a) = [p_i \log(p_i/a_i) \ge 0]$ is the Kullback-Leibler "distance" between a and p.

To maximize MP is to minimize I, i.e. to choose a = p. However, in practice p would not be known. It would seem natural then to choose for a the best available estimate, say \hat{p} of p.

The approximate *IP*-criterion given in equation (2.3) can be obtained by entirely different reasoning. Suppose a random sample of the region's current population is drawn with replacement and each individual so-obtained is classified by subregion. Let p_d be the observed sample fraction of individuals from subregion t. Then the function of ggiven in equation (2.3) is the likelihood function for these data. It would be maximized to find the optimal estimate among allowable choices of g to find the maximum likelihood estimate of the true subregion population proportions. This would be g = p unless g were constrained.

This sampling-theoretical point of view suggests a natural alternative to the criterion given in equation (2.3). If the hypothetical sample were large and the a_i 's were the "true" regional proportions, then the consistency of the (unconstrained) maximum likelihood estimator, p, would yield the approximation $\log(a_{i'}/p_{i'}) = (a_{i'}/p_{i'}-1)-(a_{i'}/p_{i'}-2)^2/2$. This, in turn would yield

$$MP(\underline{A}) = (\underline{A}/\underline{P}) \exp \left[-\frac{1}{2} \sum (a_{\underline{d}} - \underline{P}_{\underline{d}})^2 / \underline{P}_{\underline{d}}\right] \qquad (2.4)$$

Equation (2.4) suggests the minimum $\partial h'$ -squared criterion for choosing the $\{a_{ij}\}$, namely minimize

$$\chi^2 = \sum (a_{\xi} - p_{\xi})^2 / p_{\xi}$$
. (2.5)

Again, if a were unconstrained, a = p would be optimal.

This last criterion among others receives special attention in Kitagawa et al.(1980). It represents according to these authors, a compromise between $\sum (a_{\underline{c}} - p_{\underline{c}})^2$ and $\sum |a_{\underline{c}} - p_{\underline{c}}|/p_{\underline{c}}$. The first would be unduly sensitive to large individual, subregional misallocations and the second, to misallocations in small areas.

REGRESSION METHODS

Such methods yield post-censal estimates, P_{g_i} of current (time t = 2, say) population sizes, P_{g_i} , for subregions. They rely on models which involve coefficients which must be fitted, the observed current values of symptomatic variables, $S_{g_i} = (S_{g_ij}, \ldots, S_{g_ip})$, and their corresponding values at time t = 1, S_{g_ij} , when population sizes, P_{g_i} , are available for all regions, $i = 1, \ldots, n$.

To fit the coefficients, a criterion is chosen, usually least squares but possibly that given by I or χ^2 in equations (2.3) or (2.5) respectively. The times, t = 1 and 2, are taken to be successive census years so that the $\{P_{jij}\}, \{S_{jij}\}, j = 1, 2, i = 1, ..., n$ are observable. Then p_{ij} is taken to be the observed value of the "dependent" variable, $P_{2ij}/[P_{2ij}], i = 1, ..., n$ while a_{ij} is supplied by the regression model, albeit with as yet unspecified coefficients. The coefficients are then chosen to give the criterion-functional its least possible value.

Erickson (1973, 1974) proposed an Interesting variation of this scheme which might be called "sampling-regression". The time t = 2 is the present and the model is fitted to the results of a census carried out on a subsample of subregions. This approach takes account of the inevitable temporal nonstationarity of any regression model by "tuning" its coefficients to the oresent. The sampling scheme which was introduced in Section 2 to provide an interpretation of the JUP criterion could provide a (no doubt inferior) alternative to Erickson's plan. A bonus of the sampling-regression approach which is pointed out by Purcell and Kish (1979) is that, unlike the conventional approach, it carries with it a basis for inference.

It remains to choose a suitable model. As will now be shown, simple intuitive requirements lead quite easily to particular models.

For expository simplicity it will be assumed that p = 1 so that the symptomatic variables at times t = 2 and t = 2 become $S_{1\bar{t}}$ and $S_{2\bar{t}}$ respectively. The "dot" will as usual represent summation over any subscript it replaces. Thus P_g , would denote the population of the entire region at time t = 2.

Let $\hat{P}_2 = \tilde{z}(S_2, S_1, \tilde{P}_1)$ represent an unspecified post-censal regression estimation model. Here $\tilde{z} = (\tilde{x}_1, \dots, \tilde{x}_n)^T$ is the vector of subregional estimation models while $S_{j\bar{z}} = (S_{j\bar{z}}, \dots, S_{j\bar{n}})^T$, j = 1, 2, and $P_1 = (P_{1\bar{1}}, \dots, P_{1\bar{n}})^T$ are the corresponding vectors of symptomatic variables and population sizes at time t = 1." Observe that, in general, \tilde{x}_i is a function of \tilde{S}_2, \tilde{S}_1 and \tilde{P}_1 for all \tilde{z} .

Suppose, hypothetically, that subregional population sizes were to be assigned without error by means of the model, \underline{T} introduced above. Various reasonable requirements of such a rule suggest themselves and lead, as is shown below, to explicit forms for \underline{T} . These may be taken as first approximations to the actual population sizes and so used to obtain estimates. The precision of such estimates and hence the value of the models suggested by the approach outlined above, would need to be ascertained by empirical study.

Possible requirements for $\underline{2}$ are presented and discussed. Their implications are given below.

Regional population sizes may be estimated with a relatively high precision compared to those of its subregions. So they may be regarded as known. The following requirement is therefore considered as fundamental:

CTT (Controlled-to-Total). $\hat{P}_2 = P_3$.

It implies that we may regard \underline{P}_S and \underline{P}_J as vectors of proportions, each set summing to 1, rather than as vectors of population counts whenever it is expedient to do so.

Another important condition is:

PI (Permutation invariance). For all n x n permutation matrices 0

$$\underline{\tau}(4\underline{s}_{2},\ 4\underline{s}_{1},\ 4\underline{s}_{1})=\phi\ \tau(\underline{s}_{3},\ \underline{s}_{1},\ \underline{s}_{1}).$$

This condition simply ensures that the order in which the sub-regions are listed is irrelevant. It is not the same as spatial homogeneity because the symptomatic variables, S_1 and S_2 , may contain information pertaining to the geography of the sub-regions.

The following invariance requirements ensure that the model has good robustness properties. They make use of the fact that the scale on which the symptomatic variables are measured is usually arbitrary. Furthermore, if the symptomatic variables are counts computed locally (e.g. births, marriages, deaths) there may be a degree of underreporting that varies from one region to another. To some extent these effects are eliminated if the following invariance conditions are satisfied:

- RI (Regional Invariance). For all positive diagonal matrices D $\underline{T}(DS_2, DS_1, P_1) = T(S_3, S_1, P_1)$
- TI (Temporal Invariance). For all positive scalars $a_1, a_2 > 0$ $\underline{\tau}(a_2, \underline{s}_2, a_1, \underline{s}_1, \underline{s}_1) = \underline{\tau}(\underline{s}_2, \underline{s}_1, \underline{s}_1)$.

The next condition embraces a natural equivariance requirement. It derives from the recognition that the estimated growth rate in each region, i, \hat{P}_{gg}/P_{1i} , would not change even if \underline{P}_{1} 's co-ordinates were replaced by their corresponding densities, say per hectare or per square kilometre, for example. This suggests the condition $T_{i}(\underline{S}_{2}, \underline{S}_{1}, \underline{P}_{2})/(\underline{P}_{1i}) = T_{i}(\underline{S}_{2}, \underline{S}_{1}, \underline{DP}_{2})/(\underline{d}_{i}, \underline{P}_{1i})$ $i = 1, \ldots, n$, for any positive $D = \operatorname{diag}[d_{1}, \ldots, d_{n}]$. This is equivalent to

$$\begin{split} \mathbb{T}(\underline{S}_2, \underline{S}_1, \underline{DP}_1) &= D\mathbb{T}(\underline{S}_2, \underline{S}_1, \underline{P}_1), \text{ i.e. } \mathbb{T}_{\underline{t}}(\underline{S}_2, \underline{S}_1, \underline{P}_1) = P_{\underline{t}\underline{t}} \mathbb{T}_{\underline{t}}(\underline{S}_2, \underline{S}_1), \\ \vec{t} &= 1, \dots, n \quad \text{for some function } \mathbb{H}. \quad \text{However, this condition is} \\ \text{inconsistent with CTT which is considered more fundamental. And it} \\ \text{is unnecessarily strong for it is possible, as our analysis will show,} \\ \text{to derive a sufficiently small class of potential models under the} \\ \text{weaker requirement that} \quad \mathcal{Q}_{\underline{t}\underline{t}}(\underline{S}_2, \underline{S}_1, \underline{P}_1) \triangleq (\widehat{P}_{\underline{t}\underline{t}}/P_{\underline{t}\underline{t}}) + (\widehat{P}_{\underline{t}\underline{t}}/P_{\underline{t}\underline{t}}), \quad \text{the} \\ \text{relative growth rates of regions } \vec{t} \quad \text{and} \quad j, \quad \vec{t}, j = 1, \dots, n, \text{ be} \\ \text{invariant in the sense described above. This condition is} \\ \mathcal{G}_{\underline{t}\underline{t}}(\underline{S}_2, \underline{S}_1, \underline{S}_1, \underline{P}_1) = \mathcal{G}_{\underline{t}\underline{t}}(\underline{S}_2, \underline{S}_1, \underline{P}_1) \quad \vec{t}, \vec{t} = 1, \dots, n, \text{ for all positive} \\ \text{diagonal matrices, } \mathcal{D} \text{ which is equivalent to} \end{split}$$

RGRI (Relative Growth Rate Invariance). For
$$i, j = 1, ..., n$$

and all positive $D = \text{Diag}(d_1, ..., d_n)$,
 $G_{ij}(S_2, S_1, D_{ij}^p) = G_{ij}(S_2, S_1, P_1)$
where $G_{ij}(x, y, D_n) = \left[T_i(x, y, D_n)/d_i\right] \left[T_j(x, y, D_n)/d_j\right]^{-2}$
for all x, y and x in T^*s domain.

Instead of regarding, as we may because of CTT, \underline{P}_1 and \underline{P}_2 as vectors of proportions, it is more convenient to deal with the (n-1)-vector of ratios, $\underline{P}_1^{\ *}$, with elements $P_{\underline{1}\underline{0}}/P_{\underline{1}n}$, $\underline{i} = 1, \ldots, n-1$. Similar changes in \underline{T} yields $\underline{P}_2^{\ *} = \underline{T}^*(\underline{S}_2, \underline{S}_1, \underline{P}_1^{\ *})$. The advantage of $\underline{P}_2^{\ *}$ over \underline{P}_2 , the vector of proportions, is that the elements of $\underline{P}_3^{\ *}$ are unrestricted even when CTT is imposed.

Conditions RI and TI are equivalent to $\underline{\mathcal{I}}^{*}(\underline{\mathcal{S}}_{2}, \underline{\mathcal{S}}_{1}, \underline{\mathcal{P}}_{1}^{*}) = \underline{\mathcal{I}}^{*}\left[(\underline{\mathcal{S}}_{2}^{*}/\underline{\mathcal{S}}_{1}^{*}), \underline{\mathcal{I}}, \underline{\mathcal{P}}_{1}^{*}\right]$ and RGRI to $\underline{\mathcal{I}}_{\underline{\mathcal{C}}}^{*}(\underline{\mathcal{S}}_{3}, \underline{\mathcal{S}}_{1}^{*}, \underline{\mathcal{P}}_{1}^{*}) = \underline{\mathcal{P}}_{\underline{\mathcal{C}}}^{*}\underline{\mathcal{I}}_{\underline{\mathcal{C}}}^{*}(\underline{\mathcal{S}}_{3}, \underline{\mathcal{S}}_{1}, \underline{\mathcal{P}}_{1}^{*}) = \underline{\mathcal{P}}_{\underline{\mathcal{C}}}^{*}\underline{\mathcal{I}}_{\underline{\mathcal{C}}}^{*}(\underline{\mathcal{S}}_{3}, \underline{\mathcal{S}}_{1}, \underline{\mathcal{P}}_{1}^{*}) = \underline{\mathcal{P}}_{\underline{\mathcal{C}}}^{*}\underline{\mathcal{I}}_{\underline{\mathcal{C}}}^{*}(\underline{\mathcal{S}}_{3}, \underline{\mathcal{S}}_{1}, \underline{\mathcal{I}})$ $\underline{\mathcal{L}} = \underline{\mathcal{I}}, \dots, \underline{n}$. So RI, TI and RGRI combined are equivalent to

$$T_{\underline{i}}^{*}(\underline{S}_{\underline{2}}, \underline{S}_{\underline{1}}, \underline{P}_{\underline{1}}^{*}) = P_{\underline{i}}^{*} H_{\underline{i}} \left[(\underline{S}_{\underline{2}}/\underline{S}_{\underline{1}})^{*} \right]$$

(3.1)

for some function H_2 , i = 1, ..., n.

We regard CTT, PI, RI, TI and RGP.I as important: the first two would seem to be essential but it is possible to think of conditions under which the invariance requirements are not so compelling. The following conditions, although plausible, seem less important. They, or conditions like them, are needed to reduce the general model given in equation (3.1) to a more explicit, applicable form:

> TC (Temporal Coherence). Given an additional time, z = 0 $\underline{\tau}(\underline{S}_2, \underline{S}_0, \underline{P}_0) = \underline{\tau}(\underline{S}_2, \underline{S}_1, \hat{P}_1)$ where $\hat{P}_1 = \underline{\tau}(\underline{S}_1, \underline{S}_0, P_0).$

TR (Time Reversibility). For all
$$S_1$$
, S_2 , P_1
 $\underline{P}_1 = \underline{\tau}(\underline{S}_1, \underline{S}_2, \hat{P}_3)$ where
 $\hat{\underline{P}}_3 = \underline{\tau}(\underline{S}_3, \underline{S}_3, \underline{P}_3)$.

The intuitive basis for TC is clear and it nearly implies TR since the latter becomes the former when the past is reflected into the future provided $\underline{T}(\underline{S}_0, \underline{S}_1, \underline{P}_1)$ is given the value \underline{P}_0 , the subregions' known population sizes when the symptomatic variable, \underline{S}_0 , reassumes its original (time t = 0) value.

These conditions imply an easily derived explicit form for the function B_{g} of equation (3.1). The solution, f(y) = ay of Euler's functional equation f(x+y) = f(x) + f(y), f continuous, is used. The solution obtains even if the requirement $x + y \le K$ is imposed, $0 \le x, y \le K$. The result is

$$s_{i}(y) = \prod_{j=1}^{n} y_{j}^{\alpha_{ij}}$$
(3.2)

for certain constants, $-\omega < \alpha_{i,j} < =$.

If condition PI is imposed, a straightforward argument which is omitted for brevity shows that $a_{ij} = a$ or θ according as i = j or $i \neq j$ for some constant a. The model implied by all of the conditions given above is, then, in summary,

$$\hat{P}_{2\ell} = P_2, P_{1\ell}, P_{\ell} = \left[\sum_{k=1}^{n} P_{1k}, P_k\right]^{\alpha}$$

where $R_{\underline{\ell}} = S_{\underline{3}\underline{\ell}}/S_{\underline{1}\underline{\ell}}$. Of course, $P_{\underline{1}\underline{\ell}}$ and $R_{\underline{\ell}}$ can be replaced by $P_{\underline{3}\underline{\ell}}^*$ and $R_{\underline{\ell}}^*$ in equation (3.3) without changing the result. Alternagively, they can be replaced by their corresponding "shares", $P_{\underline{1}\underline{\ell}}/P_{\underline{1}}^*$ and $R_{\underline{\ell}}/R_{\underline{\ell}}^*$, respectively. In the case of more than one symptomatic variable, obvious extensions of the conditions stated above lead to a generalization of the model given in equation (3.3) :

$$\hat{P}_{2i} = P_2, \quad P_{ji} \quad \prod_{j=1}^{p} \stackrel{\alpha j}{R_{ij}} \quad \left[\begin{array}{c} n \\ \sum_{k=1}^{n} P_{jk} & \prod_{j=1}^{p} P_{kj} \end{array} \right]^{-1}$$
(3.4)

where B_{ij} , the counterpart of R_{j} , is the growth rate for symptomatic variable j, j = 1, ..., p.

An alternative form of the model derived above is obtained by taking logarithms. The result is a variant of that suggested by Morrison and Relles (1976) :

$$\log(\hat{P}_{gi}/P_{ji}) = \int_{j=1}^{p} a_j \log R_{ij}$$
 (3.5)

The basis for their choice is not given. The resulting estimate, $\hat{P}_{j_{2}}$, would not necessarily satisfy CTT and would therefore need to be "controlled" after the parameters, $\alpha_{j_{2}}$, j = 1, ..., p were fitted. The resulting estimates would, in general, differ from those obtained by applying the model given in (3.4). Neither of the models (3.4) or (3.5) has undergone a comparative empirical study.

Approximations to these models may be derived using the approximation, $\log x \approx x-2$ which is accurate if $x \approx 1$. So in stable subregions i, that is, those undergoing slow changes in population size, equation (3.5), for example, would yield the approximate model

$$\hat{P}_{2i}/P_{1i} = a_0 + \sum_{j=1}^{p} a_j R_{ij}$$
 (3.6)

This model, or that obtained by replacing all its components by their shares, is called the ratio-correlation model. It is the most commonly used of the various alternatives and is the regression component of the composite methodology used in the Federal State Co-operative Program of the United States (Kitagawa et al. 1981). Empirical studies presented in this last-cited work show that this model yields very good results in exactly those subregions where the above approximation would be accurate.

If in the model of equation (3.5) the P's are replaced by their shares and logarithms are taken, an approximate model is obtained :

$$\hat{p}_{2i} - p_{2i} = \alpha_0 + \sum_{j=2}^{p} \alpha_j (r_{2ij} - r_{2ij})$$
(3.7)

where $p_{k\bar{k}} = P_{k\bar{k}'}P_{k}$, and $r_{k\bar{k}\bar{j}} = S_{k\bar{k}\bar{j}'}S_{k,\bar{j}'}$, k = 1, 2, $\bar{k} = 1, ..., n$, and $\bar{j} = 1, ..., p$. This model would be expected to be appropriate for subregions whose population sizes comprise a relatively large fraction of that of the region. This so-called difference-correlation model was proposed by 0'Hare (1976) who provides empirical results which show that this model is marginally superior to its ratio counterpart. However, under this model, the data may well be heteorscedastic (Dr. D. Herman, personal communication) so its value remains uncertain.

CONCLUDING REMARKS

This paper presents a new regression model (equation (3.4)) for estimating the population sizes of small-areas. Its value remains to be determined by empirical study. However, the argument of Section 3 makes it a natural choice. The ratio-correlation model in current use in the U.S.A.'s Federal-State Co-operative Program would approximate the proposed model in sub-regions of stable population size for reasons given in Section 3.

In fitting and empirically evaluating a model such as derived in Section 3, a criterion or "objective function" must be specified. That obtained in Section 2 (equation (2.3)) is novel in this context. It is, approximately, the Nash criterion (Nash 1950) for determining the joint value to a group (society) of any proposed group action (revenue allocation scheme). It implies that proportional allocation on the basis of population size is optimal. Thus the value of an allocation scheme will depend on how well the best available estimates of subregional population proportions approximate the exact values.

The proposed criterion is the product of an attempt to relate the performance of an estimation method to one of the primary objectives of the estimation program, albeit under the oversimplification of linear utility functions. While this criterion may be used in evaluating any estimation procedure through the estimates of proportions, and hence approximate allocation scheme it produces, its form is, fortuitously, particularly appropriate for the regression model developed in Section 3.

Regression models used in current practice are fitted by least squares even though their performance is measured by other criteria, such as average relative absolute error. The justification for this apparent inconsistency is unknown. In any case, we would propose to use the criterion of Section 2 in both fitting and assessment.

The requirements imposed in Section 3 are the most reasonable of the various alternatives. The condition of <u>regional independence</u>, $T_{\underline{i}}(\underline{S}_{\underline{i}}, \underline{S}_{\underline{1}}, \underline{P}_{\underline{1}}) = T_{\underline{i}}(\underline{S}_{\underline{2}\underline{i}}, \underline{S}_{\underline{1}\underline{j}}, \underline{P}_{\underline{1}\underline{i}}), \underline{i} = \underline{1}, \ldots, n_s$ for example, was ruled out because it forces the model to ignore a vital piece of information, the (assumed) known regional total $P_{\underline{2}}$. If it were admitted, then PI would imply that $T_{\underline{i}} = \underline{T}$ for some \underline{T} and all \underline{i} . On top of these conditions <u>spatial coherence</u>, $\hat{j}, \hat{P}_{\underline{3}\underline{i}} = T(\hat{\underline{j}}S_{\underline{2}\underline{i}}, \hat{\underline{j}}S_{\underline{1}\underline{i}}, \hat{\underline{j}}P_{\underline{1}\underline{i}})$ for all $\underline{S}_{\underline{2}}, \underline{S}_{\underline{1}}, and \underline{P}_{\underline{1}}$, would be equivalent to $P_{\underline{2}\underline{i}} = \alpha_{\underline{2}}S_{\underline{1}\underline{i}} + \alpha_{\underline{2}}S_{\underline{2}\underline{i}} + \gamma P_{\underline{1}\underline{i}}$ for all \underline{i} . The result is not inconsistent with Ri (regional invariance) and TI (temporal invariance); however, adding this pair of conditions would lead to $\hat{P}_{\underline{2}\underline{i}} = \gamma P_{\underline{1}\underline{i}}, a model which takes no account of the symptomatic variables. Replacing the pair in question, successively by one of TC and then TR leads, respectively, to <math>\hat{P}_{\underline{3}\underline{i}} = \alpha(S_{\underline{2}\underline{i}}-S_{\underline{2}\underline{i}}) + P_{\underline{1}\underline{i}}$ and $\alpha(S_{\underline{1}\underline{i}}-\beta S_{\underline{2}\underline{i}) + \beta P_{\underline{2}\underline{i}}, \beta^2 = 1, according as TC or TR is introduced. The efficacy of such linear models is unknown.$

REFERENCES

- Belsley, D.A., Kuh, E. and Welsch, R.E. (1980). Regression Diagnostics. New York: Wiley.
- Erickson, E.P. (1973). A method for combining sample survey data and symptomatic indicators to obtain population estimates for local areas. *Demography*, 10, 137-160.
- Erickson, E.P. (1974). A regression method for estimating population changes of local areas. Jour. Amer. Statist. Assoc., 69, 867-875.
- Hauser, P.M. (1981). The census of 1980. Solentific American, v. 245, # 5, 53-61.
- Kitagawa, E.M. et al. (1980). Estimating population and income of amall areas. Washington: National Academy Press.
- Morrison, P.A. and Relles, D.A. (1976). A method of monitoring smallarea population changes in cities. Public Data Use, v.3, 10-15.
- Namboodiri, N.K. (1972). On the ratio-correlation and related methods of subnational population estimation. Demography, v.9, 443-453.
- Nash, J.F. Jr. (1950). The bargaining problem. Econometrica, 18, 155-162.
- O'Hare, N. (1976). Report on a multiple regression method for making population estimates. Demography, v.13, 369-379.
- Purcell, N.J. and Kish, L. (1979). Estimation for small domains. *Biometrico*, 35, 365-384.
- Rosenberg, H. (1968). Improving current population estimates through stratification. Land Economice, 44, 331-338.
- Schmitt, R.C. and Crosetti, A.H. (1954). Accuracy of the ratio correlation method for estimating postcensal population. Land Economics, 30, 279-281.
- Zidek, J.V. (1982). A review of methods for estimating the populations of local areas. Tech. Rep. 82-4. Instit. of Applied Maths. and Statist., U. of British Columbia.