# SYNTHETIC ESTIMATION WHEN ONLY PARTIAL SYMPTOMATIC INFORMATION IS AVAILABLE

Thomas Gerig, North Carolina State University
Teresa Lopez Aivarez, Mexico City

## 1. INTRODUCTION

Sample surveys are frequently designed on large domains of study. But administrative decisions are being increasingly based on data on subdivisions of the original domains of study. These subdivisions are called small domains or subdomains. It has been found convenient to combine several sources of information in an effort to obtain efficient small domain estimates. The commonly available sources of data are the administrative registers or population census and the survey which is planned on the large domain. The most popular and simplest approach for obtaining small domain estimators of proportions is the so called method of synthetic estimation which is described next. The available data from census records are the proportions of individuals belonging to subgroup j within small domain i. Where the subgroups $j=1,\ldots J$ are formed according to one or more symptomatic categorical variables (e.g. sex, race). Then it can be assumed that the units from the population are grouped according to three classifications. The first classification is that corresponding to the small domains, the second one to the subgroups and the third one to the variable of interest. Let the classifications be, respectively, A with I categories, $A_1$, $A_2\ldots A_I$; B with J categories: $B_1,B_2,\ldots B_J$ and C with K categories $C_1,C_2,\ldots C_K$. We may also think of A, B and C as categorical random variables, and write, say, $A_i$ to denote the event $\{A = A_i\}$. The standard assumption of synthetic estimation is that $P(C_k|A_i \cap B_j) = P(C_k|B_j)$, (independent of i). Let $\pi_{ij} = P(B_j|A_i)$ for $i=1,\ldots,I$, $j=1,\ldots,J$ which are known from census records and $\alpha_{jk} = P(C_k|B_j)$ for $j=1,\ldots,J$ and $k=1,\ldots,K$ which can be estimated from survey data. The parameter to be estimated is: $P(C_k|A_i)$ = Probability that a unit belongs to the category k of the variable of interest, given that it comes from domain i; which will be denoted by $\phi_{k,i}$. This can be expressed in terms of the $\alpha_{jk}$'s and the $\pi_{ij}$'s and under the standard assumption in the following form

$$P(C_k|A_i) = \sum_{j=1}^{j} \alpha_{jk} \pi_{ij} \qquad (1.1)$$

The synthetic estimator of $P(C_k|A_i)$ is:

$$\hat{\phi}_{k,i} = \hat{P}(C_k|A_i) = \sum_{j=1}^{J} \hat{\alpha}_{jk} \pi_{ij} \qquad (1.2)$$

where $\hat{\alpha}_{jk}$ are estimates of $\alpha_{jk}$.

## 2. PARTIAL SYMPTOMATIC INFORMATION

It commonly happens that census data are not available for all combinations of the symptomatic variables for each small domain. When this occurs the standard formula (1.2) can not be applied. However, certain linear combinations of the proportions $\pi_{ij}$ are known (for each small domain i). Thus the following restrictions are imposed on the parameters. Let A be a known matrix $(L \times J)$ of full row rank and $\underline{b}_i$ be a known vector $(L \times 1)$, $i=1,\ldots,I$, then the parameters must satisfy

$$A \underline{\pi}_i = \underline{b}_i \quad \text{for} \quad i=1,\ldots,I \qquad (2.1)$$

where $\underline{\pi}'_i = (\pi_{i1},\ldots,\pi_{iJ})$, the first row of A consists of units, and the first element of $\underline{b}_i$ is unity. The purpose of this work is to estimate (1.1) subject to (2.1), the restriction that $\sum_k \alpha_{jk} = 1$ and the restrictions that $0 < \alpha_{jk} < 1$ and $0 < \pi_{ij} < 1$. Two methods of estimation are developed. The first one is the maximum likelihood procedure which estimates jointly the $\alpha_{jk}$'s and the $\pi_{ij}$'s. The second one consists of estimating the $\pi_{ij}$'s subject to some linear constraints by the Iterative Proportional Fitting (IPF) procedure developed by Deming and Stephan (1940), and estimate the $\alpha_{jk}$'s by the ML method. Both estimators are obtained in the case that a stratified random sampling on the large domain is assumed, where the strata are the small domains. Let $n_{ijk}$ be the number of units falling into the (i,j,k) cell and denote $\underline{n}'_i = (n_{i11},\ldots,n_{ijk})$ for $i=1,\ldots,I$, $n_{i++} = \sum_j \sum_k n_{ijk}$ and $n = n_{+++} \sum_i n_{i++}$. The $n_{i++}$ is fixed for $i = 1,\ldots,I$ and $n_{i++}$ is considered of the form $n_{i++} = \mu_i n$ where $\mu_i$ are constants such that $\Sigma\mu_i = 1$ and $\mu_i > 0$. Under the standard assumption for each i, $\underline{n}'_i$ has a multinomial distribution with parameter $\underline{P}'_i = (\pi_{i1} \alpha_{11}, \ldots,\pi_{iJ} \alpha_{JK})$ and the multinomial distributions are independent.

## 3. MAXIMUM LIKELIHOOD ESTIMATOR OF $\phi_{k,i}$

Let $\underline{\theta}' = (\underline{\alpha}'_1, \underline{\alpha}'_2,\ldots,\underline{\alpha}'_k, \underline{\pi}'_1,\ldots,\underline{\pi}'_I)$ where

$$\underline{\alpha}'_k = (\alpha_{1k}, \alpha_{2k},\ldots,\alpha_{jk}) \quad k = 1,\ldots,K \quad \text{and}$$

$$\underline{\pi}'_i = (\pi_{i1},\ldots,\pi_{iJ}) \qquad i = 1,\ldots,I.$$

The likelihood function of $\theta$ is

$$L(\underline{\theta}) = \text{Constant} \prod_{j=1}^{J} \prod_{k=1}^{K} \alpha_{jk}^{n_{+jk}} \times \prod_{i=1}^{I} \prod_{j=1}^{J} \pi_{ij}^{n_{ij+}} \qquad (3.1)$$

Then, the following result holds true.

Proposition 3.1 Let $n_{ij+}$ and $n_{+jk}$ be positive for every i,j and k. The maximum of the logarithm of the likelihood function (3.1) subject to $A \underline{\pi}_i = \underline{b}_i$ for $i=1,\ldots,I$, $\sum_k \alpha_{jk} = 1$ for $j=1,\ldots,J$ and the restriction that all the $\alpha$'s and $\pi$'s are strictly between 0 and 1, exists and it is unique.

Proof: See Lopez Alvarez (1982).

The estimator of $\phi_{k,i}$ is $\hat{\phi}_{k,i} = \sum_{j=1}^{J} \hat{\alpha}_{jk} \hat{\pi}_{ij}$.

Asymptotic properties of $\hat{\phi}_{k,i}$ will be given next.

**Notation 3.1**

$\underline{\theta}_R = (\alpha_1', \ldots, \alpha_J', \underline{\pi}_1^{(2)'}, \ldots, \underline{\pi}_I^{(2)'})$

$\underline{\alpha}_j' = (\alpha_{j1}, \alpha_{j2}, \ldots, \alpha_{jK-1})$

$\underline{\pi}_i^{(2)'} = (\pi_{i1}^{(2)}, \ldots, \pi_{iJ-L}^{(2)})$

$\underline{\pi}_i^{(1)'} = (\pi_{i1}^{(1)}, \ldots, \pi_{iL}^{(1)})$

$D_{\ell i}$ — diag $\underline{\pi}_i^{(\ell)}$ for $\ell = 1,2$ and $i = 1, \ldots, I$

$D_i$ = diag $[\underline{\pi}_i]$

$R = P \begin{bmatrix} U \\ -I \end{bmatrix}$

$T_i = \mu_i [U' \; D_{1i}^{-1} \; U + D_{2i}^{-1}]$

$V_i = D_i - D_i \; A' \; (A \; D_i \; A')^{-1} \; A \; D_i$

$S_j^{-1} = \dfrac{1}{\overline{\pi}_{\cdot j}} [- \underline{\alpha}_j \; \underline{\alpha}_j' + \text{diag}\{\alpha_{j1}, \; 1=1, \ldots, K-1\}]$

$\overline{\pi}_{\cdot j} = \sum_{i=1}^{I} \pi_{ij} \; \mu_i$

$\hat{\underline{\theta}}_1' = (\hat{\underline{\alpha}}_1', \ldots, \hat{\underline{\alpha}}_J', \hat{\underline{\pi}}_1', \ldots, \hat{\underline{\pi}}_J')$ the M.L.E. of $\underline{\theta}$

$\underline{\theta}_R' = (\hat{\underline{\alpha}}_1', \ldots, \hat{\underline{\alpha}}_J', \hat{\underline{\pi}}_1^{(2)'}, \ldots, \hat{\underline{\pi}}_I^{(2)'})$ the ML.E. of $\underline{\theta}_R$

**Lemma 3.2**

Using notation 3.1 then $\mu_i \; R \; T_i^{-1} \; R' = V_i$

Proof. See Lopez Alvarez (1982).

**Proposition 3.3**

i) $\hat{\underline{\theta}} \xrightarrow{P} \underline{\theta}$ and

ii) $\hat{\underline{\theta}}$ is $AN(\underline{\theta}, \frac{1}{n} \Sigma(\underline{\theta}))$ where

$$\Sigma(\theta) = \begin{bmatrix} \begin{matrix} S_1^{-1} & -S_1^{-1} \; 1 \\ & 1'S_1^{-1} \; 1 \end{matrix} & & 0 & \\ & \ddots & & 0 \\ 0 & & \begin{matrix} S_j^{-1} & - \; S_j^{-1} \; 1 \\ & 1'S_j^{-1} \; 1 \end{matrix} & \\ & & & \begin{matrix} \frac{1}{\mu_1} \; V_1 & 0 \\ 0 & \ddots \frac{1}{\mu_I} \; V_I \end{matrix} \end{bmatrix}$$ (3.4)

Proof. See Lopez Alvarez (1982).

**Proposition 3.4**

$\hat{\phi}_{k,i} = \hat{\phi}_{k,i}(\hat{\underline{\theta}}) = \sum_j \hat{\alpha}_{jk} \hat{\pi}_{ij}$ is $AN(\phi_{k,i}(\underline{\theta}), \frac{1}{n} \sigma_{k,i}^2)$

where

$$\sigma_{k,i}^2 = \sigma_{k,i}^2(\underline{\theta}) = \sum_{j=1}^{J} \frac{\pi_{ij}^2 \; \alpha_{jk}(1-\alpha_{jk})}{\overline{\pi}_j} + \underline{\alpha}_k' \; \frac{1}{\mu_i} \; V_i \; \underline{\alpha}_k$$ (3.5)

for $k = 1, \ldots, K$, $i = 1, \ldots, I$.

Proof. See Lopez Alvarez (1982).

Considering the asymptotic distribution of $(\hat{\phi}_{11}, \ldots, \hat{\phi}_{K-1,J})$ the following asymptotic covariances

can be obtained.

i) $k = k'$ and $i \neq i'$ (same category, different domain)

$$\text{Cov}(\hat{\phi}_{k,i}, \hat{\phi}_{k,i'}) = \sum_{j=1}^{J} \pi_{ij} \; \pi_{i'j} \; \frac{\alpha_{jk}(1-\alpha_{jk})}{\overline{\pi}_{\cdot j}}$$

ii) $k \neq k'$ and $i = i'$ (different category, same domain)

$$\text{Cov}(\hat{\phi}_{k,i}, \hat{\phi}_{k',i}) = \sum_{j=1}^{J} \pi_{ij}^2 \; \frac{(-\alpha_{jk}\alpha_{jk'})}{\overline{\pi}_{\cdot j}} + \underline{\alpha}_k' \; V_i \; \underline{\alpha}_k$$

iii) $k \neq k'$ and $i \neq i'$ (different category, different domain)

$$\text{Cov}(\hat{\phi}_{k,i}, \hat{\phi}_{k',i'}) = \sum_{j=1}^{J} \pi_{ij} \; \pi_{i'j} \; \frac{(-\alpha_{jk} \; \alpha_{jk'})}{\overline{\pi}_{\cdot j}}.$$

Computational aspects of the estimator are discussed next. The maximum likelihood estimator of $\underline{\pi}_i$ can be obtained, when all $n_{ij+}$'s and $n_{+jk}$'s are positive, solving the likelihood equations developed from (3.2). This system of equations doesn't have a solution in closed form, so it is necessary to use an iterative method. The Newton-Rhapson method can be used to compute $\hat{\underline{\pi}}_i$. In the case when some $n_{ij+}$'s or $n_{+jk}$'s are zero. finding the maximum likelihood estimator turns out to be a problem of nonlinear programming, since it is necessary to maximize the nonlinear function (3.2) subject to the linear constraints (3.3). To solve this problem the gradient projection given in Bazaraa and Shetty (1979) was used. This method projects the gradient, which is the direction of steepest ascent, in such a way that the objective function is improved and at the same time feasibility is maintaned. Both algorithms were programmed in Fortran IV, and both require a preliminary estimate of $\pi$.

**4. IPF ESTIMATOR**

The $\hat{\alpha}_{jk}$'s are estimated by $\hat{\alpha}_{jk} = \dfrac{n_{+jk}}{n_{+j+}}$.

To estimate $\pi_{ij}$ for $j = 1, \ldots, J$, $i$ fixed, subject to $A\underline{\pi}_i = \underline{b}_i$, one of the major algorithms for the analysis of cross classified frequency counts is used, which is known as iterative proportional fitting (IPF). Denote the estimator by $\tilde{\pi}_{ij}$. Freeman and Kock (1976) mention $\tilde{\pi}_{ij}$ satisfies (or can be obtained by solving) the following equations

$$\begin{aligned} A\tilde{\underline{\pi}}_i &= \underline{b}_i \\ \overline{A} \ln \{\tilde{\underline{\pi}}_i\} &= \overline{A} \ln \frac{\underline{r}_i}{} \end{aligned}$$ (4.1)

where $\overline{A}$ is an orthocomplement of $A$, $(\overline{A} \; A' = 0)$, and

$$\underline{r}_i' = \frac{n_{i1+}}{n_{i++}}, \ldots, \frac{n_{iJ+}}{n_{i++}}$$

The estimator of $\phi_{k,i}$, which will be called IPF estimator and denoted by $\tilde{\phi}_{k,i}$ is:

$$\tilde{\phi}_{k,i} = \sum_{j=1}^{J} \hat{\alpha}_{jk} \; \tilde{\pi}_{ij}$$

**Notation 4.1**

$\hat{\underline{\theta}}' = (\hat{\underline{\alpha}}_1', \hat{\underline{\alpha}}_2', \ldots, \hat{\underline{\alpha}}_J', \tilde{\underline{\pi}}_1', \tilde{\underline{\pi}}_1', \ldots, \tilde{\underline{\pi}}_I')$, where

$\hat{\underline{\alpha}}_j = (\hat{\alpha}_{j1}, \hat{\alpha}_{j2}, \ldots, \hat{\alpha}_{jK})$ for $j = 1, \ldots, J$

$\tilde{\underline{\pi}}_i' = (\tilde{\pi}_{i_1}, \tilde{\pi}_{i_2}, \ldots, \tilde{\pi}_{iJ})$

$\underline{e}'_k$ elementary vector of dimension K

$\underline{e}'_j$ elementary vector of dimension J

$\tilde{D}_i = \text{diag } [\tilde{\underline{\pi}}_i]$

$D_{\underline{P}_1} = \text{diag } [\underline{P}_1] ; \quad \underline{P}_1 = (P_{111}, \ldots, P_{1JK}).$

The asymptotic distribution of $\hat{\theta}'$ is obtained next, and from this the asymptotic distribution of $\tilde{\phi}_{k,i}$ can be easily obtained.

## Proposition 4.1

$\tilde{\theta}$ is AN$(\theta, \frac{1}{n} \Sigma(\underline{\theta}))$, where $\Sigma(\underline{\theta})$ is given by (3.4).

**Proof.** See Lopez Alvarez (1982).

It has been shown that $\hat{\theta}$ and $\tilde{\theta}$ have the same asymptotic distribution, therefore $\tilde{\phi}_{k,i} = \phi_{k,i}(\tilde{\underline{\theta}})$ has the same asymptotic distribution of $\hat{\phi}_{k,i} = \phi_{k,i}(\hat{\underline{\theta}})$ given in proposition 3.4. A consistent estimator of the asymptotic variance of $\tilde{\phi}_{k,i}$ is $\sigma^2_{k,j}(\tilde{\underline{\theta}})$, where $\sigma^2_{k,i}(\underline{\theta})$ is given by (3.5). Notice that $\tilde{\phi}_{k,i}$ is a BAN estimator, since $\hat{\phi}_{k,i}$ is BAN. Computational aspects of $\tilde{\pi}_{k,i}$ are discussed next.

The vector $\tilde{\pi}_i$ may be computed using the subroutine CTLLF of the International Mathematical and Statistical Libraries, Inc. (1981). This subroutine adjusts frequency tables to some given set of marginal constraints using the IPF method. The IPF method doesn't provide a means of computing an estimate for a cell that is empty in the original table, that is if $n_{ij+} = 0$ then the estimate of $\pi_{ij}$ is zero. This leads to the fact that if certain relative frequencies are zero, the IPF method might not converge to $\tilde{\pi}_i$ which satisfies the equations (4.1). In this case there is no estimate of $\phi_{k,i}$ by the IPF method. Another problem with having null cells is that the matrix $A\tilde{D}_iA'$ is singular. To compute an estimate of the variance of the estimator the Moore Penrose inverse $(A\tilde{D}_iA')^+$ was substituted in the expression of $\sigma^2_{k,i}(\underline{\theta})$.

## 5. MONTE CARLO STUDY

In order to compare the small sample properties of the estimates, evaluate the usefulness in small samples of the asymptotic variance formulas of the estimates and explore the effects of departure from the basic assumption of synthetic estimation, a Monte Carlo study was performed. Data which simulated 1,000 replications of an experiment were constructed using known values of $\pi$'s and $\alpha$'s, and hence of $\phi_{k,i}$. Estimates of $\phi_{k,i}$ were calculated by both methods in each experiment. In this Monte Carlo study it is considered that the population is divided according to two domains, with the domain of interest being number 1, that is i=1,2. Also the population is divided in eight subgroups j=1,...,8. The categorical variable of interest has two categories; k=1,2. The simulation sample size is 3,000. Two multinomial distributions were generated for each replication, one with parameters $\underline{p}'_1$ and $n_{1++}$, and the other with parameters $\underline{p}'_2$ and $3000 - n_{1++}$. The vectors $\underline{p}_i$ are of the form

$$\underline{p}'_i = (\pi_{i1}\alpha_{11}, \pi_{i2}\alpha_{21}, \ldots, \pi_{i8}\alpha_{81}, \pi_{i1}\alpha_{12}, \pi_{i2}\alpha_{22}, \ldots,$$
$$\pi_{i8}\alpha_{82}) \text{ for } i = 1,2.$$

The multinomial distributions were generated using the subroutine GGMTN of the IMSL (1981) package.

Three different sample sizes were considered, which represent 1%, 4% and 10% of the simulation sample size, so according to Purcell and Kish (1979), domain 1 can be considered a small domain for these cases. Other simulation parameters are: the number of restrictions of $\pi_i$; the size of $P(C_1|A_1)$ and the departure from the basic assumption of synthetic estimation. The departure from the standard assumption is obtained by considering:

$$\alpha_{ijk} = \alpha_{jk} + \delta_{ik} \quad \text{where } 0 < |\delta_{ik}| < 1 \text{ and } \sum_{k=1}^{K} \delta_{ik} = 0.$$

Some results of the Monte Carlo study are shown in tables 5.1, 5.2 and 5.3. On the basis of the results given in table 5.1 we conclude that, the IPF estimator is unbiased when there is no departure from the basic assumption of synthetic estimation and when the sample size of the small domain represents 1%, 4% or 10%. The ML estimator sometimes shows bias, even when there is no departure from the basic assumption. When there is moderate departure the ML estimator is biased. For large departure both estimators are biased, notice that the bias is greater in the cases that $\phi$ is small. Also it can be seen that the ML estimator is more efficient than the IPF estimator, especially when the sample size of the small domain represents 1% of the simulation sample size. A close look at tables 5.2 and 5.3 shows that the average estimated asymptotic standard deviation (a.s.d.) usually overestimates the standard deviation, regardless of the method of estimation. As it is expected the asymptotic standard deviation always underestimates the standard deviation of the estimators. Finally, notice that the s.d of the ML estimator is always smaller than the s.d of the IPF estimator.

## 6. CONCLUSIONS

The estimator which is obtained using the IPF method is easier to compute than the ML estimator due to the simplicity of the IPF algorithm. It was found that the asymptotic distributions of the ML estimator and the IPF estimator are the same, that is, both are best asymptotically normal estimators. The results of the Monte Carlo study indicate that both estimators are biased when there is a large departure from the basic assumption of synthetic estimation. The ML estimator is sometimes biased even when there is no departure from the assumption and the sample size of the small domain represents 4% of the simulation sample size. A possible reason why the ML estimator is biased is that the algorithm used to solve the nonlinear programming problem may not perform well in certain cases when some cells have null counts. The Monte Carlo study also showed that the ML estimator is more efficient than the IPF estimator.

## BIBLIOGRAPHY

BAZARAA, M.S. and Shetty, C.H. (1979). *Nonlinear Programming Theory and Algorithms.* John Wiley and Sons.

BRADLEY, R.P. and Gart, J.J. (1962). *The asymptotic properties of ML estimators when sampling from associated populations.* Biometrika, 49, 205-214.

DEMING, W.E. and Stephan, F.F. (1940). *On a least squares adjustment of a sampled frequency table when the expected marginal totals are known.*

Ann. Math. Statist. 11, 427–444.

FREEMAN, D.H. and Koch, G.G. (1976). *An asymptotic covariance structure for testing hypothesis on ranked contingency tables from complex sample surveys*. 1976 Proceedings of the Social Statistical Section, ASA, 330–333.

INTERNATIONAL MATHEMATICAL AND STATISTICAL LIBRARIES, INC. (1981) *IMSL Library, Ninth Edition*, Houston: International Mathematical and Statistical Section, ASA, 330–333.

KHATRI, C.G. (1966). *A note on manova model applied to problems in growth curve.* Annals of the Institute of Statistical Mathematics 18(1); 75–86.

LOPEZ ALVAREZ, T. (1982). *Synthetic estimation when only partial symptomatic information is available*. Ph.D. Thesis. North Carolina State University at Raleigh.

MAKELAINEN, T., Schmidt, K. and Styan, G.P.H. (1981). *On the existence and uniqueness of the maximum likelihood estimate of a vector valued parameter on fixed size samples*. Annals of Statistics, 76, 758–767.

PURCELL, N.J. and Kish, L. (1979). *Estimation for small domains*. Biometrics, 35, 365–384.

TABLE 5.1

| $P(C_1|A_1)$ | Number of Restrictions | Sample size of small domain | Percent of bias ML | IPF | eff(IPF,ML) |
|---|---|---|---|---|---|
| | | No Departure | | | |
| .10 | 3 | 1% | 1 | 0* | .71 |
| .10 | 3 | 4% | 0 | 0 | .93 |
| .10 | 4 | 1% | 2 | 0 | .81 |
| .10 | 4 | 4% | 1 | 0 | .98 |
| .10 | 3 | 10% | 0 | 0 | .99 |
| .47 | 3 | 1% | 0 | 0 | .80 |
| .47 | 3 | 4% | 0 | 0 | .99 |
| .47 | 3 | 10% | 0 | 0 | 1.00 |
| | | Moderate Departure | | | |
| .47 | 3 | 1% | 1 | 0 | .87 |
| .47 | 3 | 4% | 1 | 0 | .94 |
| | | Large Departure | | | |
| .12 | 3 | 1% | 8 | 10 | .71 |
| .12 | 3 | 4% | 8 | 9 | .89 |
| .12 | 4 | 4% | 8 | 9 | .83 |
| .46 | 3 | 1% | 2 | 1 | .95 |
| .46 | 3 | 4% | 2 | 1 | .90 |

*0 bias indicates a non significant t-test.

TABLE 5.2

Estimated asymptotic s.d, asymptotic s.d and empirical s.d of ML and IPF estimators when there is no departure from the assumption of synthetic estimation.

| Estimator | $P(C_1|A_1)$ | Number of restrictions | Sample size of small domain | s.d. $\times 10^3$ | Average estimated a.s.d$\times 10^3$ | a.s.d. $\times 10^3$ |
|---|---|---|---|---|---|---|
| ML | .10 | 3 | 1% | 7.185 | 8.550 | 6.387 |
| IPF | .10 | 3 | 1% | 8.718 | 8.706 | 6.387 |
| ML | .10 | 3 | 4% | 6.710 | 6.890 | 6.327 |
| IPF | .10 | 3 | 4% | 6.952 | 6.970 | 6.327 |
| ML | .10 | 4 | 1% | 7.115 | 8.446 | 6.387 |
| IPF | .10 | 4 | 1% | 8.614 | 8.495 | 6.387 |
| ML | .10 | 4 | 4% | 6.849 | 6.884 | 6.327 |
| IPF | .10 | 4 | 4% | 6.819 | 6.919 | 6.327 |
| ML | .47 | 3 | 1% | 12.524 | 14.748 | 11.42 |
| IPF | .47 | 3 | 1% | 14.674 | 14.939 | 11.421 |
| ML | .47 | 3 | 4% | 12.452 | 12.458 | 11.240 |
| IPF | .47 | 3 | 4% | 12.476 | 12.413 | 11.240 |

TABLE 5.3

Estimated asymptotic s.d and s.d of the ML and IPF estimators
when there is departure from the assumption of synthetic esti-
mation.

| Estimator | $P(C_1 \mid A_1)$ | Number of restrictions | Sample size of small domain | s.d. $\times 10^3$ | Average estimated a.s.d. $\times 10^3$ |
|---|---|---|---|---|---|
| | | | Moderate Departure | | |
| ML | .47 | 3 | 1% | 13.042 | 14.789 |
| IPF | .47 | 3 | 1% | 15.023 | 14.989 |
| ML | .47 | 3 | 4% | 12.125 | 12.166 |
| IPF | .47 | 3 | 4% | 12.269 | 12.214 |
| | | | Large Departure | | |
| ML | .12 | 3 | 1% | 7.497 | 8.468 |
| IPF | .12 | 3 | 1% | 9.141 | 8.630 |
| ML | .12 | 3 | 4% | 6.682 | 6.913 |
| IPF | .12 | 3 | 4% | 6.846 | 6.990 |
| ML | .12 | 4 | 4% | 6.848 | 6.899 |
| IPF | .12 | 4 | 4% | 7.047 | 6.938 |
| ML | .46 | 3 | 1% | 12.625 | 14.738 |
| IPF | .46 | 3 | 1% | 14.982 | 14.931 |
| ML | .46 | 3 | 4% | 12.293 | 12.154 |
| IPF | .46 | 3 | 4% | 12.383 | 12.198 |