

On influence of clustering population on estimation accuracy of population totals vector

Janusz L. Wywiał, Grzegorz Sitek

Department of Statistics, Econometrics and Mathematics
University of Economics in Katowice, Poland,
janusz.wywial@ue.katowice.pl, grzegors12@wp.pl

ECDA, Wrocław, 27-29 September 2017

Contents:

1. Simple cluster sample vector estimator.
2. Relative efficiency.
3. Clustering algorithms.
4. Accuracy analysis.
5. Conclusion.
6. References.

Simple cluster sample vector estimator

Basic notation

- U is population of size N ,
- m is number of variables observed in U ,
- $\mathbf{y}_k = [y_{k,1} \dots y_{k,m}]$ where $k \in U$,
- $\mathcal{D} = \{U_1, \dots, U_h, \dots, U_G\}$ is partition of U into clusters U_h ,
 $h = 1, \dots, G$, $\bar{N} = N/G$,

$$\bar{\mathbf{y}} = [\bar{y}_1 \dots \bar{y}_m] = \sum_{k \in U} \mathbf{y}_k / N, \quad \mathbf{y}_U = N\bar{\mathbf{y}} = \sum_{k \in U} \mathbf{y}_k = [y_{U,1} \dots y_{U,m}],$$

$$y_{U,i} = \sum_{k \in U} y_{k,i}, \quad \mathbf{C} = [c_{i,j}], \quad c_{i,j} = \sum_{k \in U} (y_{k,i} - \bar{y}_i)(y_{k,j} - \bar{y}_j) / (N-1),$$

$$\mathbf{R} = \mathbf{D}^{-1/2} \mathbf{C} \mathbf{D}^{-1/2} = [r_{i,j}], \quad \mathbf{D} = [v_i], \quad r_{i,j} = \frac{c_{i,j}}{\sqrt{v_i v_j}}, \quad v_i = c_{i,i}.$$

Simple cluster sample vector estimator

Basic notation

$$\bar{\mathbf{y}}_{U_h} = \sum_{k \in U_h} \mathbf{y}_k / N_h, \quad \bar{\mathbf{y}}_{U_h} = [\bar{y}_{U_h,1} \dots \bar{y}_{U_h,m}], \quad \bar{y}_{U_h,i} = \sum_{k \in U_h} y_{k,i} / N_h,$$

$$\mathbf{y}_{U_h} = N_h \bar{\mathbf{y}}_{U_h} = \sum_{k \in U_h} \mathbf{y}_k, \quad \mathbf{y}_{U_h} = [y_{U_h,1} \dots y_{U_h,m}], \quad y_{U_h,i} = \sum_{k \in U_h} y_{k,i},$$

$$\mathbf{C}_{U_h} = [c_{U_h,i,j}], \quad c_{U_h,i,j} = \sum_{k \in U_h} (y_{k,i} - \bar{y}_{U_h,i})(y_{k,j} - \bar{y}_{U_h,j}) / (N_h - 1),$$

$$\bar{\mathbf{y}}_D = \sum_{h=1}^G \mathbf{y}_{U_h} / G = \mathbf{y}_U / G = [\bar{y}_{D,1} \dots \bar{y}_{D,m}], \quad \bar{y}_{D,i} = \sum_{h=1}^G y_{U_h,i} / G = y_{U,i} / G,$$

$$\mathbf{y}_D = G \bar{\mathbf{y}}_D = \sum_{h=1}^G \mathbf{y}_{U_h} = \mathbf{y}_U, \quad \mathbf{C}_D = [c_{D,i,j}],$$

$$c_{D,i,j} = \sum_{h=1}^G (y_{U_h,i} - \bar{y}_{D,i})(y_{U_h,j} - \bar{y}_{D,j}) / (G - 1), \quad i, j = 1, \dots, m.$$

Simple cluster sample vector estimator

Properties of the estimator

$$\tilde{\mathbf{y}}_S = \frac{G}{g} \sum_{h \in S} \sum_{k \in U_h} \mathbf{y}_k = \frac{G}{g} \sum_{h \in S} \mathbf{y}_{U_h}, \quad \mathbf{V}(\tilde{\mathbf{y}}_S) = \frac{G(G-g)}{g} \mathbf{C}_D, \quad (1)$$

$$\mathbf{V}(\tilde{\mathbf{y}}_S) = \frac{G(G-g)}{g} \bar{N} \mathbf{C} \left(\mathbf{I}_m + \frac{N-G}{G-1} \Delta \right) + \frac{G(G-g)}{g} \mathbf{A} \quad (2)$$

$$\text{homogeneity matrix:} \quad \Delta = \mathbf{I}_m - \mathbf{C}^{-1} \mathbf{C}_*, \quad (3)$$

$$\mathbf{A} = [a_{i,j}]; \quad a_{i,j} = \frac{1}{G-1} \sum_{h=1}^G (N_h - \bar{N}) N_h \bar{y}_{U_h,i} \bar{y}_{U_h,j}, \quad (4)$$

$$\mathbf{C}_* = [c_{*i,j}], \quad c_{*i,j} = \frac{1}{N-G} \sum_{h=1}^G \sum_{k \in U_h} (y_{k,i} - \bar{y}_{U_h,i})(y_{k,j} - \bar{y}_{U_h,j}), \quad (5)$$

The eigenvalues of Δ and the diagonal elements of Δ take values from $\left[-\frac{G-1}{N-G}; 1\right]$.

Simple cluster sample vector estimator

Relative efficiency

The relative efficiency coefficient (Rao and Scott (1981):

$$deff(\tilde{\mathbf{y}}_S) = \lambda \left(\mathbf{V}(\mathbf{y}_S)^{-1} \mathbf{V}(\mathbf{t}_S) \right) \propto \lambda \left(\mathbf{C}^{-1} \mathbf{C}_D \right). \quad (6)$$

where $\lambda(\dots)$ is maximal eigenvalue and ordinary estimator

$$\mathbf{y}_S = \frac{N}{n} \sum_{k \in S} \mathbf{y}_k, \quad \mathbf{V}(\mathbf{y}_S) = \frac{N(N-n)}{n} \mathbf{C} \quad (7)$$

Estimator $\tilde{\mathbf{y}}_S$ is not worse than \mathbf{y}_S if and only if $\mathbf{V}(\tilde{\mathbf{y}}_S) - \mathbf{V}(\mathbf{y}_S)$ is non-positive definite and all eigenvalues of $\mathbf{V}(\mathbf{y}_S)^{-1} \mathbf{V}(\mathbf{t}_S)$ take values from $[0; 1]$.

$$deff(\tilde{\mathbf{y}}_S) = 1 + \lambda \left(\frac{N-G}{G-1} \Delta + \frac{1}{N} \mathbf{C}^{-1} \mathbf{A} \right). \quad (8)$$

If $N_h = \text{const}$ for all $h = 1, \dots, G$:

$$0 \leq deff(\tilde{\mathbf{y}}_S) = 1 + \frac{N-G}{G-1} \lambda(\Delta) \leq \frac{N-1}{G-1}. \quad (9)$$

Clustering algorithms

Systematic algorithm \mathcal{D}_1

- Let $\mathbf{y}_k > \mathbf{0}$ for all $k = 1, \dots, N$,
- evaluation of squared distances $d_k = \mathbf{y}_k \mathbf{y}_k^T$ of \mathbf{y}_k from the zero vector $\mathbf{0}$ for all $k \in U$,
- let us assume that $d_k \leq d_{k+1}$ for $k = 1, \dots, N - 1$,
- h -th cluster is identified by such $k \in U_h$ that $k = (i - 1)G + h$, for $i = 1, \dots, M$ and $h = 1, \dots, G$,
- this leads to: $d_{U_h} \leq d_{U_{h+1}}$ for $h = 1, \dots, G - 1$ where $d_{U_h} = \sum_{k \in U_h} d_k$.

Clustering algorithms

Systematic algorithm \mathcal{D}_2

- Let $d_k = (\mathbf{y}_k - \bar{\mathbf{y}})(\mathbf{y}_k - \bar{\mathbf{y}})^T$ be the squared distance of \mathbf{y}_k from vector $\bar{\mathbf{y}}$ for all $k \in U$,
- let $d_k \leq d_{k+1}$ for $k = 1, \dots, N - 1$.
- when M is even and $N = MG$, then

$$U_h = \left\{ (h-1)\frac{M}{2} + i; N - (h-1)\frac{M}{2} - i + 1 \right\}$$

for $h = 1, \dots, G$ and $i = 1, \dots, M/2$.

- Particularly, if $M = 2$ and $N = MG$,

$$U_h = \{h; N - h + 1\}$$

for $h = 1, \dots, G$.

Clustering algorithms

Permutation algorithm \mathcal{D}_3

- Let $\mathcal{D}^{(0)} = \{U_1^{(0)}, \dots, U_G^{(0)}\}$ be any start partition of population into clusters of the same sizes,
- in the t -h ($t=0,1,\dots$) iteration partition $\mathcal{D}^{(t)} = \{U_1^{(t)}, \dots, U_G^{(t)}\}$ is generated through permutation population elements of U ,
- for assumed $t = T$, \mathcal{D}_3 is treated as optimal when

$$\mathcal{D}_3 = \arg(\min_{\{t=1,\dots,T\}}(\lambda(\Delta(\mathcal{D}^{(t)}))))). \quad (10)$$

Clustering algorithms

Algorithm \mathcal{D}_4

- Let $\mathcal{D}^{(0)} = \{U_1^{(0)}, \dots, U_G^{(0)}\}$ be any start partition the population into clusters of not necessary of the same size,
- let $f : U \rightarrow \mathcal{D}^{(t)}$, $f_t(k) = h$, if and only if $k \in U_h^{(t)}$.
- in iteration $t + 1$ we randomly choose number k_* from $1, \dots, N$,
- element k_* is moved from the cluster $h_{\#} = f_t(k_*)$ to cluster h_* . h_* is randomly drawn from $\{h : h = 1, \dots, G; h \neq h_{\#}\}$. This leads to new partition $\mathcal{D}^{(t+1)}$,
- let $\lambda_{t+1} = \lambda(\mathbf{C}_{\mathcal{D}^{(t+1)}})$. If $\lambda_{t+1} < \lambda_t$, then $\mathcal{D}^{(t+1)}$ is the current partition and we start the iteration $t + 2$ of the algorithm,
- if $\lambda_{t+1} \geq \lambda_t$, then we start the stage $t + 2$ of the algorithm from partition $\mathcal{D}^{(t)}$;
- the algorithm is stopped when number of the iteration reaches assumed level T ;
- this algorithm minimizes $deff(\tilde{\mathbf{y}}_S)$.

Clustering algorithms

Algorithm \mathcal{D}_5

- $\mathcal{D}^{(t)} = \{U_1^{(t)}, \dots, U_G^{(t)}\}$ is resulted of t -th iteration where $t = (l - 1)N + k$, $k = 1, \dots, N$, $l = 1, 2, \dots$;
- let $\lambda_t = \lambda(\mathbf{C}_{\mathcal{D}^{(t)}})$ and let $f : U \rightarrow \mathcal{D}^{(t)}$, $f_t(l) = h \Leftrightarrow l \in U_h^{(t)}$;
- in stage $t + 1$ element $k \in U_h^{(t)}$, where $h = f_t(k)$, is moved to clusters $U_z^{(t)}$, $z \neq h$, $z = 1, \dots, G$ and calculated the following

$$(k, \underline{z}) = \arg(\min_{\{z=1, \dots, G, z \neq f_t(k)\}} (\lambda(\mathbf{C}_{\mathcal{D}^{(t)}}(k, z)))) \quad (11)$$

- $\lambda(\mathbf{C}_{\mathcal{D}^{(t)}}(k, z))$ is evaluated for the partition $\mathcal{D}^{(t)}$ in which clusters $U_z^{(t)}$ and $U_h^{(t)}$ are replaced by $\{U_z^{(t)} \cup \{k\}\}$ and $\{U_h^{(t)} - \{k\}\}$, respectively, and $h = f_t(k)$;

Clustering algorithms

Algorithm \mathcal{D}_5 , continuation

- If $\lambda(\mathbf{C}_{\mathcal{D}^{(t)}}(\underline{z})) < \lambda_t$, then $\lambda_{t+1} = \lambda(\mathbf{C}_{\mathcal{D}^{(t+1)}})$ and $\mathcal{D}^{(t+1)}$ is equal to $\mathcal{D}^{(t)}$ where clusters $U_{\underline{z}}^{(t)}$ and $U_h^{(t)}$ are replaced by $U_{\underline{z}}^{(t+1)} = \{U_{\underline{z}}^{(t)} \cup \{k\}\}$ and $U_h^{(t+1)} = \{U_h^{(t)} - \{k\}\}$, respectively;
- if $\lambda(\mathbf{C}_{\mathcal{D}^{(t)}}(\underline{z})) \geq \lambda_t$, then $\mathcal{D}^{(t+1)} = \mathcal{D}^{(t)}$ and $\lambda_{t+1} = \lambda_t$;
- the iteration process is stopped when $\lambda_{t+N} = \lambda_t$ or the number of the iterations attains the preassigned level T .

Clustering algorithms

Algorithm \mathcal{D}_6

- in iteration $t + 1$ element $k \in U_h^{(t)}$, where $h = f_t(k)$, is moved to clusters $U_z^{(t)}$, $z \neq h$, $z = 1, \dots, G$ and calculated:

$$(\underline{k}, \underline{z}) = \arg \left(\min_{\{k \in U\}} \min_{\{z \neq f_t(k), z=1, \dots, G\}} (\lambda(\mathbf{C}_{\mathcal{D}^{(t)}}(k, z))) \right) \quad (12)$$

- $\lambda(\mathbf{C}_{\mathcal{D}^{(t)}}(k, z))$ is evaluated for $\mathcal{D}^{(t)}$ in which clusters $U_z^{(t)}$ and $U_h^{(t)}$ are replaced by $\{U_z^{(t)} \cup \{k\}\}$ and $\{U_h^{(t)} - \{k\}\}$, respectively, and $h = f_t(k)$;
- if $\lambda(\mathbf{C}_{\mathcal{D}^{(t)}}(\underline{k}, \underline{z})) < \lambda_t$, then $\lambda(\mathbf{C}_{\mathcal{D}^{(t+1)}}) = \lambda(\mathbf{C}_{\mathcal{D}^{(t)}}(\underline{k}, \underline{z}))$ and $\mathcal{D}^{(t+1)}$ is equal to $\mathcal{D}^{(t)}$ where $U_z^{(t)}$ and $U_h^{(t)}$ are replaced by $U_{\underline{z}}^{(t+1)} = \{U_{\underline{z}}^{(t)} \cup \{k\}\}$ and $U_h^{(t+1)} = \{U_h^{(t)} - \{k\}\}$, respectively;
- the process is stopped when $\lambda(\mathbf{C}_{\mathcal{D}^{(t)}}(\underline{k}, \underline{z})) \geq \lambda_t$.

Accuracy analysis

Data on Swedish municipalities from Särndal C.E., et al.
Variables y_1 and y_2 are the real estate values and number of municipal employees, respectively.

Table 1. Relative efficiencies.

n	(M,g)	\mathcal{D}_1	\mathcal{D}_2	\mathcal{D}_3	\mathcal{D}_4	\mathcal{D}_5	\mathcal{D}_6
1	2	3	5	6	7	8	9
16	(2,8)	0.99	1.11	0.82	0.56	0.64	0.77
16	(4,4)	1.10	2.15	0.75	0.18	0.44	0.58
16	(8,2)	1.17	4.25	0.62	0.05	0.18	0.44
28	(2,14)	0.99	1.11	0.82	0.56	0.64	0.77
28	(4,7)	1.10	2.15	0.75	0.18	0.44	0.58
28	(14,2)	1.31	7.35	0.50	0.02	0.04	0.20
48	(2,24)	0.99	1.11	0.82	0.56	0.64	0.77
48	(4,12)	1.10	2.15	0.75	0.18	0.44	0.58
48	(8,6)	1.17	4.25	0.62	0.05	0.18	0.44

Source: Own calculations.

Conclusions

- Only under \mathcal{D}_1 and \mathcal{D}_2 the accuracy of \mathbf{y}_S is not less than the accuracy $\tilde{\mathbf{y}}_S$ for all (M, g) .
- \mathcal{D}_4 leads to the most efficient estimation based on $\tilde{\mathbf{y}}_S$.
- \mathcal{D}_3 leads to the most efficient estimation based on $\tilde{\mathbf{y}}_S$, when we assume that the population is split into clusters of the same sizes.
- For \mathcal{D}_1 and \mathcal{D}_2 the efficiency of $\tilde{\mathbf{y}}_S$ decreases, when number of clusters g decreases under fixed n .
- For \mathcal{D}_3 - \mathcal{D}_6 the efficiency of $\tilde{\mathbf{y}}_S$ increases, when number of clusters g decreases under fixed n .
- For instance, under \mathcal{D}_4 , when $(M, g) = (2, 14)$ and $(M, g) = (14, 2)$, $deff(\tilde{\mathbf{y}}_S) = 0.56$ and $deff(\tilde{\mathbf{y}}_S) = 0.02$, respectively.

Reference

- Borovkov A.A. (1981). *Mathematical Statistics. Estimation. Testing Hypotheses.* (in Russian) Nauka Moskva.
- Rao, C.R. (1973). *Linear Statistical Inference and Its Applications.* John Wiley & Sons, Inc. New York.
- Rao, J.N.K., Scott A.J. (1981). The analysis of categorical data from complex sample surveys: Chi-squared tests of goodness of fit and independence in two-way tables. *Journal of the American Statistical Association* vol. 76, no. 374, pp. 221-230.
- Särndal C.E., Swensson B., Wretman J. (1992). *Model Assisted Survey Sampling* Springer-Verlag, New York.
- Wywił J.L. (2003). *Some Contributions to Multivariate Methods in Survey Sampling.* University of Economics in Katowice. [http : // www . ue . katowice . pl / fileadmin / user _upload / wydawnictwo / Darmowe _ E – Booki](http://www.ue.katowice.pl/fileadmin/user_upload/wydawnictwo/Darmowe_E – Booki)

Thank you very much