

*Krystyna Pruska**

TESTS FOR RATIO OF TWO MEANS IN CASE OF SMALL AREAS

Abstract

The relations between characteristics for subpopulations and for the whole population are very important in small area investigations.

In the paper there are proposed testing procedures for verification of hypothesis which says that there is no difference between the ratio of small area mean and population mean for analysed variable and auxiliary variable. The properties of one considered procedure are investigated with the use of simulation methods.

Key words: small area, synthetic estimator.

I. INTRODUCTION

In small area statistics estimation of unknown parameters for subpopulation is considered as general problem. Different estimators are constructed and applied and their properties are investigated. The synthetic estimators are considered, too. They can be used when some assumptions are true. In this paper we consider possibility of verification whether these assumptions are fulfilled.

II. SYNTHETIC ESTIMATORS

In statistical literature different definitions of synthetic estimators are given (see: Dol, 1991; Bracha, 1996; Särndal et al., 1997; Kordos, 1999, Domański and Pruska, 2001). Generally, their construction is possible when some relations between parameters for subpopulation and population are constant.

* Professor, Chair of Statistical Methods, University of Łódź.

In the paper there are considered finite population only.

We will introduce the following notations:

Y – investigated variable,

X – auxiliary variable,

TY_p – total value of variable Y for population,

TX_p – total value of variable X for population,

$TY_{.g}$ – total value of variable Y for stratum g of population,

$TX_{.g}$ – total value of variable X for stratum g of population,

TY_{hg} – total value of variable Y for stratum g and small area h of population,

TX_{hg} – total value of variable X for stratum g and small area h of population,

TY_h – total value of variable Y for small area h of population,

where $g = 1, \dots, G$; $h = 1, \dots, H$ and G is number of strata in the population, H is number of small areas in the population.

If we assume (see Dol, 1991):

$$\frac{TY_{.g}}{TX_{.g}} = \frac{TY_{hg}}{TX_{hg}} \quad \text{for } g = 1, \dots, G \quad \text{and } h = 1, \dots, H \quad (1)$$

then we can consider the synthetic estimator of total value of variable Y for small area h of the following form:

$$\hat{T}Y_h = \sum_{g=1}^G \hat{\beta}_g TX_{hg} \quad (2)$$

where $\hat{\beta}_g$ is estimator of value $\frac{TY_{.g}}{TX_{.g}}$. There are different forms of statistic $\hat{\beta}_g$ (see: Dol, 1991).

Estimator (2) can be used when assumptions (1) are fulfilled. In empirical investigations we ought to verify whether the conditions are true. We may apply the estimator of TX_{hg} instead TX_{hg} in formula (2).

III. FORMULATING OF HYPOTHESIS

We can consider verification of possibility of the use of synthetic estimator as the verification of a suitable statistical hypothesis. Synthetic estimator (2) is constructed for the population which is divided into strata. Assumptions

(1) can be verified on the basis of independent samples drawn from each strata. All hypotheses which have the form:

$$H_{0g} : \frac{TY_{.g}}{TX_{.g}} = \frac{TY_{hg}}{TX_{hg}} \quad (3)$$

for $g = 1, \dots, G$ and for the fixed h from the set $\{1, \dots, H\}$ against an alternate hypothesis:

$$H_{1g} : \sim H_{0g} \quad (4)$$

can be verified analogously as hypothesis:

$$H_0 : \frac{TY_P}{TX_P} = \frac{TY_{MO}}{TX_{MO}} \quad (5)$$

against hypothesis:

$$H_1 : \sim H_0, \quad (6)$$

where TY_{MO} and TX_{MO} are total values for variables Y and X for fixed small area.

In this paper we will consider the following equivalent form of hypothesis H_0 :

$$H_0 : \frac{\mu_{YP}}{\mu_{XP}} = \frac{\mu_{YMO}}{\mu_{XMO}} \quad (7)$$

where μ_{YP} , μ_{XP} , μ_{YMO} , μ_{XMO} are means for variables Y and X for population and for small area, respectively.

IV. TEST PROCEDURES

Test statistic for the verification of hypothesis (7) can be random variable:

$$Z = \frac{\bar{Y}_P}{\bar{X}_P} - \frac{\bar{Y}_{MO}}{\bar{X}_{MO}} \quad (8)$$

or its function where \bar{Y}_P , \bar{X}_P , \bar{Y}_{MO} , \bar{X}_{MO} are sample means for variables Y and X for population and for small area, respectively.

Determining the distribution of statistic Z is difficult when we do not know the distribution of variables Y and X .

Now the test procedure for verification of hypothesis (7) will be proposed. In this procedure the conditional distributions (for fixed values of statistics \bar{X}_{MO} , \bar{X}_{P-MO} or \bar{Y}_{MO} , \bar{X}_{P-MO} respectively) of the following statistics are used:

$$U_Y = \frac{\bar{Y}_{MO} - \frac{\bar{X}_{MO}}{\bar{X}_{P-MO}} \cdot \bar{Y}_{P-MO}}{\sqrt{\frac{S_{YMO}^2}{NMO} + \frac{\bar{X}_{MO}^2}{\bar{X}_{P-MO}^2} \frac{S_{YP-MO}^2}{NPMO}}} \quad (9)$$

and

$$U_X = \frac{\bar{X}_{MO} - \frac{\bar{Y}_{MO}}{\bar{Y}_{P-MO}} \cdot \bar{X}_{P-MO}}{\sqrt{\frac{S_{XMO}^2}{NMO} + \frac{\bar{Y}_{MO}^2}{\bar{Y}_{P-MO}^2} \frac{S_{XP-MO}^2}{NPMO}}} \quad (10)$$

where

\bar{Y}_{P-MO} , \bar{X}_{P-MO} are sample means for set which is difference between population and small area and for variables Y and X , respectively;

S_{XMO}^2 , S_{XP-MO}^2 are sample variances for variable X for small area and for set which is difference between population and small area, respectively;

S_{YMO}^2 , S_{YP-MO}^2 are sample variances for variable Y for small area and for set which is difference between population and small area, respectively;

NMO is the number of these elements of sample from population, which belong to small area;

$NPMO$ is the number of these elements of sample from population, which do not belong to small area.

The test algorithm is the following:

1. We draw independently NP -element sample from the whole population. The elements of the sample belonging to small area are the sample for the small area and the elements which do not belong to small area are the sample for set which is difference between population and small area, respectively.

2. We determine the value of the following statistics: \bar{X}_{MO} , \bar{X}_{P-MO} , S_{XMO}^2 , S_{XP-MO}^2 , \bar{Y}_{MO} , \bar{Y}_{P-MO} , S_{YMO}^2 , S_{YP-MO}^2 , U_X , U_Y .

3. We verify whether $|u_X| \geq 1,96$ or $|u_Y| \geq 1,96$ where u_X and u_Y are values of variables U_X , U_Y , respectively. If one inequality is not true then we reject hypothesis (7).

We can notice that:

$$P(|U_x| \geq 1.96 \text{ or } |U_y| \geq 1.96) = P(|U_x| \geq 1.96) + P(|U_y| \geq 1.96) - P(|U_x| \geq 1.96 \text{ and } |U_y| \geq 1.96) \leq P(|U_x| \geq 1.96) + P(|U_y| \geq 1.96) = 0.05 + 0.05 = 0.1,$$

when we consider the probability for suitable conditional distribution of statistics U_X and U_Y which are asymptotic normally.

If we know total values of auxiliary variable for population and small area (i.e. we know TX_P and TX_{MO}) then we can use the following random variable as test statistic:

$$U = \frac{\bar{Y}_{MO} - \frac{TX_{MO}}{TX_{P-MO}} \cdot \bar{Y}_{P-MO}}{\sqrt{\frac{S_{YMO}^2}{NMO} + \frac{TX_{MO}^2}{TX_{P-MO}^2} \frac{S_{Y_{P-MO}}^2}{NPMO}}}. \quad (11)$$

Asymptotic distribution of statistic (11) is normal $N(0; 1)$. In this case rejection region of hypothesis (7) is determined in the classic way. We reject hypothesis when the value of statistic (11) calculated on the basis of sample belongs to the rejection region.

V. MONTE CARLO ANALYSIS OF PROPOSED TEST PROCEDURE PROPERTIES

Monte Carlo analysis deals with the first presented test procedure, it means the case when the total values of auxiliary variable for population and small area are unknown.

The aim of the conducted experiments was determining the number of cases in which the hypothesis (7) was rejected in 1000 repetitions for fixed distribution of population.

The experiments were conducted in the following way:

1. Creating the population consisting of 50 000 values of variable (Y, X) which are generated from fixed distribution or two fixed distributions.

2. Determining the small area which consists of 5000 elements.

3. Drawing n -element sample from population ($n = 2000, 2500$).

4. Verification of hypothesis (7) by means of test whose statistics are variables (9) and (10).

5. Conducting 1000 repetitions of stages 3. and 4.

6. Determining the number of cases in which the hypothesis (7) was rejected in 1000 repetitions.

The results of Monte Carlo experiments are presented in Tables 1. and 2.

Table 1. Number of cases of rejection of hypothesis (7) among 1000 experiments for the same distribution of (Y, X) in population and small area

No.	Distribution of variable (Y, X)	Size of sample from population	Minimal size of sample from small area	Average size of sample from small area	Maximal size of sample from small area	Number of rejection of H_0
1	$N(m, \Sigma)$, $m^T = [100; 20]$ $\Sigma = \begin{bmatrix} 100 & 36 \\ 36 & 16 \end{bmatrix}$	2000	156	200	251	123
		2500	207	250	299	113
2	$N(m, \Sigma)$, $m^T = [10; 20]$ $\Sigma = \begin{bmatrix} 1 & 0,8 \\ 0,8 & 1 \end{bmatrix}$	2000	156	200	251	136
		2500	207	250	299	140
3	$X \sim N(10; 1)$ $Y = [X]$	2000	149	200	247	0
		2500	196	250	305	0
4	$X \sim N(60; 12)$ $Y = [X]$	2000	149	200	247	0
		2500	196	250	305	0
5	$X \sim P_5$ $Y = X + Z$ $Z \sim N\left(\frac{5}{2}; \frac{\sqrt{5}}{10}\right)$	2000	161	200	252	0
		2500	188	250	312	1
6	$X \sim P_{10}$ $Y = X + Z$ $Z \sim N\left(5; \frac{\sqrt{10}}{10}\right)$	2000	158	200	253	0
		2500	209	250	293	0

Source: Author's calculations.

Table 2. Number of cases of rejection of hypothesis (7) among 1000 experiments for the different distributions of (Y, X) in population and small area

No.	Distribution of variable (Y, X) in small area	Distribution of variable (Y, X) out of small area	Size of sample from population	Minimal size of sample from small area	Average of sample from small area	Maximal size of sample from small area	Number of cases of H_0
1	$N(m, \Sigma),$ $m^T = [120; 25]$ $\Sigma = \begin{bmatrix} 100 & 36 \\ 36 & 16 \end{bmatrix}$	$N(m, \Sigma),$ $m^T = [100; 20]$ $\Sigma = \begin{bmatrix} 100 & 36 \\ 36 & 16 \end{bmatrix}$	2000	156	200	251	120
			2500	207	250	299	113
2	$N(m, \Sigma),$ $m^T = [12; 21]$ $\Sigma = \begin{bmatrix} 1 & 0,8 \\ 0,8 & 1 \end{bmatrix}$	$N(m, \Sigma),$ $m^T = [10; 20]$ $\Sigma = \begin{bmatrix} 1 & 0,8 \\ 0,8 & 1 \end{bmatrix}$	2000	156	200	251	1000
			2500	207	250	299	1000
3	$X = U + 2$ $Y = V + 1$ $U \sim N(10; 1)$ $V = [U]$	$X \sim N(10; 1)$ $Y = [X]$	2000	149	200	247	1000
			2500	196	250	305	1000
4	$X = U + 2$ $Y = V + 1$ $U \sim N(60; 12)$ $V = [U]$	$X \sim N(60; 12)$ $Y = [X]$	2000	149	200	247	0
			2500	196	250	305	0
5	$X = U + 2$ $Y = V + 1$ $U \sim P_5$ $V = U + Z$ $Z \sim N\left(\frac{5}{2}; \frac{\sqrt{5}}{10}\right)$	$X \sim P_5$ $Y = X + Z$ $Z \sim N\left(\frac{5}{2}; \frac{\sqrt{5}}{10}\right)$	2000	161	200	252	1000
			2500	188	250	312	1000
6	$X = U + 2$ $Y = V + 1$ $U \sim P_{10}$ $V = U + Z$ $Z \sim N\left(\frac{5}{2}; \frac{\sqrt{5}}{10}\right)$	$X \sim P_{10}$ $Y = X + Z$ $Z \sim N\left(\frac{5}{2}; \frac{\sqrt{5}}{10}\right)$	2000	158	200	253	1000
			2500	209	250	293	1000

Source: Author's calculations.

We can notice that the considered test procedure does not reject hypothesis (7) or rejects in a few cases among 1000 repetitions when hypothesis is true. If hypothesis (7) is not true and distributions of investigated variable in population and small area differ significantly then hypothesis (7) is rejected in all repetitions for a given case. If the distributions do not differ significantly then hypothesis (7) is rejected in some repetitions or in no repetitions. Such results are not typical for significance tests. The proposed test procedure can be modified and its properties ought to be investigated.

VI. FINAL REMARKS

Different methods for the verification of hypothesis about relations between total values for subpopulation and population can be constructed. Some propositions are presented in the paper, another one (using bootstrap method) is given in the Pruska's paper (2002). The problems, considered in the paper, are important in investigations of small area. It seems necessary to continue the conducted analyses.

REFERENCES

- Bracha Cz. (1996), *Teoretyczne podstawy metody reprezentacyjnej*, PWN, Warszawa.
- Dol W. (1991), *Small Area Estimation. A Synthesis between Sampling Theory and Econometrics*, Wolters Noordhoff, Groningen.
- Domański Cz., Pruska K. (2001), *Metody statystyki małych obszarów*, Wyd. Uniwersytetu Łódzkiego, Łódź.
- Kordos J. (1999), Problemy estymacji danych dla małych obszarów, *Wiadomości Statystyczne*, **1**, 85–101.
- Pruska K. (2002), Statystyczna weryfikacja możliwości zastosowania estymatorów syntetycznych w badaniach małych obszarów, referat wygłoszony w Łagowie na konferencji „Statystyka regionalna w jednoczącej się Europie” (2–5.09.2002).
- Särndal C., Swensson B., Wretman J. (1997), *Model Assisted Survey Sampling*, Springer-Verlag, New York.

Krystyna Pruska

**TESTY DLA STOSUNKU DWÓCH ŚREDNICH
W PRZYPADKU MAŁYCH OBSZARÓW**

Streszczenie

Relacje pomiędzy charakterystykami podpopulacji i całej populacji są bardzo ważne w badaniach małych obszarów.

W pracy tej zaproponowane są procedury testowe, służące do weryfikacji hipotezy o równości stosunku średniej dla małego obszaru do średniej z populacji dla analizowanej zmiennej i zmiennej pomocniczej. Własności jednej z zaproponowanych procedur badane są za pomocą metod symulacyjnych.