

A. PLUCIŃSKA (Warszawa)

WARIANCJA ŚREDNIEJ W PEWNYM SCHEMACIE LOSOWANIA

W badaniach metodą reprezentacyjną mamy często do czynienia z następującą sytuacją: Populację generalną stanowi s arkuszy, a i -ty arkusz ($i = 1, 2, \dots, s$) zawiera u_i pozycji. Celem badania reprezentacyjnego jest oszacowanie wyrażenia

$$(1) \quad \bar{u} = \frac{1}{s} \sum_{i=1}^s u_i,$$

gdzie s jest znane, ale u_i są nieznane.

Za przykład może służyć populacja generalna złożona z arkuszy spisowych spisu powszechnego ludności, przy czym każdy arkusz zawiera kilka osób, mianowicie osoby prowadzące wspólne gospodarstwo domowe. Metodą reprezentacyjną chcemy ustalić średnią ilość osób w gospodarstwie.

W rozważanej sytuacji dogodnie jest nieraz stosować następujący schemat losowania: Traktujemy arkusz jako jednostkę losowania i losujemy dla każdego arkusza, czy ma wejść do próby czy nie, przy czym losowanie jest niezależne, a każdy arkusz ma jednakowe prawdopodobieństwo p ($0 < p < 1$) dostania się do próbki. W rozważanym schemacie mamy s jednostek losowania.

Oznaczamy przez X_i ($i = 1, 2, \dots, s$) liczbę elementów indywidualnych wylosowanych w i -tym losowaniu. Mamy

$$(2) \quad P(X_i = u_i) = p, \quad P(X_i = 0) = 1 - p = q \quad (i = 1, 2, \dots, s).$$

Oznaczamy przez X liczbę indywidualnych elementów w próbie. Mamy

$$(3) \quad X = \sum_{i=1}^s X_i.$$

Niech N oznacza liczbę jednostek losowania, które weszły do próbki. N jest zmienną losową o rozkładzie dwumianowym

$$(4) \quad P(N = n) = \binom{s}{n} p^n q^{s-n}.$$

Za estymator nieznannej wielkości \bar{u} określonej wzorem (1) przyjmujemy

$$(5) \quad Y = X/N.$$

Celem niniejszej noty jest znalezienie wariancji $D^2(Y)$ zmiennej losowej Y i porównanie jej z wariancją średniej arytmetycznej z próbki o stałej liczbie ps elementów wybranych z populacji zgodnie ze schematem Bernoulliego. Ta ostatnia wariancja równa jest, jak wiadomo, $D^2(u)/ps$, gdzie

$$D^2(u) = \frac{1}{s} \sum_{i=1}^s (u_i - \bar{u})^2.$$

W dalszych rozważaniach będziemy jednak pomijali ewentualność, że $N = 0$, gdyż prawdopodobieństwo tego zdarzenia wynoszące q^s jest przy tych wartościach s , z którymi mamy do czynienia w praktyce, bardzo małe. Ściśle mówiąc znajdziemy wariancję $D^2(Y|N > 0)$, tj. wariancję zmiennej losowej Y pod warunkiem, że $N > 0$. Udowodnimy

TWIERDZENIE. Niech zmienne losowe X i N mają rozkłady określone wzorami (2)-(4) i niech Y będzie dane wzorem (5). Wówczas zachodzą relacje

$$(a) \quad E(Y|N > 0) = \bar{u},$$

$$(b) \quad D^2(Y|N > 0) = \frac{D^2(u)}{1-q^s} \sum_{n=1}^s \frac{q^{s-n} - q^s}{n},$$

przy czym dla $s \rightarrow \infty$ zachodzi równość asymptotyczna

$$(c) \quad D^2(Y|N > 0) \cong D^2(u)/ps.$$

Dowód. Mamy

$$(6) \quad E(Y|N > 0) = \frac{\sum_{n=1}^s E(Y|N = n)P(N = n)}{1-q^s} = \frac{\bar{u}(1-q^s)}{1-q^s} = \bar{u}.$$

Teza (a) została więc udowodniona. Widzimy, że zmienna losowa ($Y|N > 0$) jest nieobciążonym estymatorem parametru \bar{u} .

Obliczmy teraz drugi moment. Mamy

$$(7) \quad E(Y^2|N > 0) = \frac{\sum_{n=1}^s E(Y^2|N = n)P(N = n)}{1-q^s}.$$

Dla każdego n

$$E(Y^2|N = n) = D^2(Y|N = n) + [E(Y|N = n)]^2,$$

a więc korzystając z wzorów (5) i (6) otrzymujemy

$$E(Y^2|N = n) = D^2(u)/n + \bar{u}^2.$$

Stąd i z wzoru (7) wynika

$$E(Y^2|N > 0) = \frac{\sum_{n=1}^s \left[\frac{D^2(u)}{n} + \bar{u}^2 \right] P(N = n)}{1 - q^s} = D^2(u) \frac{\sum_{n=1}^s \frac{1}{n} P(N = n)}{1 - q^s} + \bar{u}^2.$$

Ostatecznie więc

$$D^2(Y|N > 0) = \frac{D^2(u) \sum_{n=1}^s \frac{1}{n} P(N = n)}{1 - q^s}.$$

Oznaczmy

$$(8) \quad \varphi(p) = \sum_{n=1}^s \frac{1}{n} P(N = n) = \sum_{n=1}^s \frac{1}{n} \binom{s}{n} p^n (1-p)^{s-n}.$$

W celu obliczenia tej sumy różniczkujemy ją najpierw względem p :

$$\begin{aligned} \varphi'(p) &= \sum_{n=1}^s \binom{s}{n} p^{n-1} (1-p)^{s-n} - \sum_{n=1}^s \frac{s-n}{n} \binom{s}{n} p^n (1-p)^{s-n-1} = \\ &= \frac{1}{p(1-p)} [1 - (1-p)^s] - \frac{s}{1-p} \varphi(p). \end{aligned}$$

Otrzymaliśmy równanie różniczkowe liniowe rzędu pierwszego. Rozwiązaniem równania jednorodnego

$$\varphi'(p) = -\frac{s}{1-p} \varphi(p)$$

jest

$$\varphi(p) = c(1-p)^s.$$

Po uzmiennieniu stałej i dokonaniu drobnych przekształceń otrzymujemy

$$c'(p) = \frac{1}{p(1-p)^{s+1}} [1 - (1-p)^s].$$

Aby obliczyć $\int \frac{dp}{p(1-p)^{s+1}}$ korzystamy z wzoru

$$\frac{1}{p(1-p)^{s+1}} = \frac{1}{p} + \frac{1}{(1-p)^{s+1}} + \dots + \frac{1}{1-p},$$

którego słuszność łatwo wykazać indukcyjnie, i znajdujemy

$$c = \sum_{n=1}^s \frac{(1-p)^{n-1-s}}{s+1-n} + A,$$

gdzie A jest stałą całkowania.

Wstawiając wartość c otrzymujemy

$$(9) \quad \varphi(p) = \sum_{n=1}^s \frac{(1-p)^{n-1}}{s+1-n} + A(1-p)^s.$$

Stałą A wyznaczamy przyjmując w wyrażeniu (8) $p = 0$. Stąd otrzymujemy $\varphi(0) = 0$, a następnie porównujemy z (9).

Ostatecznie

$$A = \sum_{n=1}^s \frac{1}{s+1-n},$$

$$\varphi(p) = \sum_{n=1}^s \frac{(1-p)^{n-1} - (1-p)^s}{s+1-n} = \sum_{n=1}^s \frac{q^{s-n} - q^s}{n}.$$

Tym samym wykazany został wzór (b).

Przechodzimy do dowodu wzoru (c).

Oszacujemy sumę

$$\sum_{n=1}^s \frac{q^{s-n}}{n} = \sum_{n=0}^{s-1} \frac{q^n}{s-n}.$$

W tym celu skorzystamy z rozwinięcia

$$\frac{1}{s-n} = \frac{1}{s} + \frac{n}{s^2} + \frac{n^2}{s^3} + \dots,$$

które daje

$$\sum_{n=0}^{s-1} \frac{q^n}{s-n} = \frac{1}{s} \sum_{n=0}^{s-1} q^n + \frac{1}{s^2} \sum_{n=0}^{s-1} nq^n + \dots = \frac{1}{s} \sum_{n=0}^{s-1} q^n + O\left(\frac{1}{s}\right) \cong \frac{1}{sp}.$$

Zauważmy dalej, że

$$\lim_{s \rightarrow \infty} q^s \sum_{n=1}^s \frac{1}{s+1-n} = O\left(\frac{1}{s}\right),$$

mamy bowiem

$$\lim_{s \rightarrow \infty} sq^s \sum_{n=1}^s \frac{1}{s+1-n} < \lim_{s \rightarrow \infty} s^2 q^s = 0.$$

Jeżeli ponadto uwzględnimy, że $\lim_{s \rightarrow \infty} (1-q^s) = 1$, to będziemy mieli

$$D^2(Y|N > 0) \cong D^2(u)/ps.$$

Wzór (c) został więc udowodniony.

Widzimy więc, że dla dużych wartości s jednakowa jest efektywność omawianego w niniejszej pracy estymatora ($Y|N > 0$) określonego przez (5) i średniej arytmetycznej z próbki o stałej liczbie ps elementów wybranych z populacji. Asymptotyczną wariancję średniej dla rozważanego typu schematu losowania znalazł Fiz [1].

Dla niezbyt dużych wartości s średnia arytmetyczna z próbki o stałej liczbie ps elementów jest estymatorem efektywniejszym. Łatwo to sprawdzić przez zbadanie różnicy wariancji obu estymatorów, która równa jest

$$\begin{aligned} R &= \frac{D^2(u)}{ps} - \frac{D^2(u)}{1-q^s} \sum_{n=1}^s \frac{q^{s-n} - q^s}{n} = \\ &= \frac{D^2(u)}{sp(1-q^s)} \left[1 - q^s - ps \left(\frac{1}{s} \sum_{n=1}^s q^{s-n} + \sum_{n=1}^s \frac{s-n}{sn} q^{s-n} - \sum_{n=1}^s \frac{q^s}{n} \right) \right] = \\ &= C \sum_{n=1}^s \frac{-(s-n)q^{1-n} + sq^s}{n} < C \sum_{n=1}^s [-(s-n)q^{s-n} + sq^s] = \\ &= C \sum_{k=0}^{s-1} (-kq^k + sq^s) = C \sum_{k=0}^{s-1} q^k (-k + sq^{s-k}) < C \sum_{k=0}^{s-1} (-k + sq^{s-k}) = \\ &= C \left(-\frac{s(s-1)}{2} + sq \frac{1-q^s}{1-q} \right), \end{aligned}$$

gdzie

$$C = \frac{D^2(u)}{s^2 p(1-q^s)}.$$

Zbadamy dla jakich wartości q różnica $R < 0$, czyli kiedy są spełnione nierówności

$$\begin{aligned}sq(1-q^s)/(1-q) &< s(s-1)/2, \\q - q^{s+1} &< (s-1)/2 - q(s-1)/2, \\q - 2q^{s+1}/(s+1) &< (s-1)/(s+1).\end{aligned}$$

Zatem $R < 0$, przynajmniej gdy $q < (s-1)/(s+1)$.

Nierówność $q < (s-1)/(s+1)$ obejmuje wszystkie praktycznie ważne przypadki. Za pomocą podobnych oszacowań można wykazać, że estymator (5) jest efektywniejszy, gdy

$$q - 2q^{s+1}/(s+1) > (s-1)/(s+1).$$

Nierówność tę spełnia np. wartość $q = (s-1)/s$.

Praca cytowana

[1] M. Fisz, *Efektywność dwóch schematów losowania*, Studia i prace statystyczne 1 (1951).

Praca wpłynęła 12. 2. 57

А. ПЛУЦИНСКАЯ (Варшава)

ДИСПЕРСИЯ СРЕДНЕГО В НЕКОТОРОЙ СХЕМЕ ВЫБОРОВ

РЕЗЮМЕ

Генеральная совокупность состоит из s единиц выбора, каждая i -ая единица ($i = 1, 2, \dots, s$) содержит u_i индивидуальных элементов. Целью репрезентативного исследования является оценка выражения (1)

$$\bar{u} = \frac{1}{s} \sum_{i=1}^s u_i.$$

В качестве оценки этой величины принимаем Y , определенное формулой (5). N есть случайная величина равная числу единиц выбора, которые вошли в состав выборки и которая по определению больше нуля. Доказывает, что

- 1) Y является несмещенной оценкой,
- 2) дисперсия $D^2(Y)$ определена формулой (b),
- 3) асимптотическая эффективность оценки Y и среднего арифметического выборка с постоянным числом rs элементов одинаковы, но для небольших значений s среднее арифметическое выборки с постоянным числом rs элементов является более эффективной оценкой, чем оценка Y , по крайней мере для $q < (s-1)/(s+1)$.

A. PLUCIŃSKA (Warszawa)

VARIANCE OF THE MEAN IN A CERTAIN SAMPLING SCHEME

SUMMARY

The general population consists of s sampling units, every i -th unit ($i = 1, 2, \dots, s$) contains u_i individual elements. The object of the survey is the estimation of expression (1).

$$\bar{u} = \frac{1}{s} \sum_{i=1}^s u_i.$$

We take Y defined by (5) as the estimator of that quantity. N is a random variable denoting the number of sampling units which have entered into the sample, that number being greater than zero by hypothesis. We prove that

- 1) Y is an unbiased estimator,
 - 2) the variance $D^2(Y)$ is defined by (b),
 - 3) the asymptotic effectiveness of the estimator Y and that of the arithmetic mean of a sample with a constant number ps of elements are identical, but for not very large s the arithmetic mean of a sample with the constant number ps of elements is a more effective estimator than the estimator Y , at least for $q < (s-1)/(s+1)$.
-