

S. TRYBUŁA (Wrocław)

OCENA FREKWENCJI W POPULACJI ELEMENTÓW NALEŻĄCYCH DO KLAS NIEREPREZENTOWANYCH W PRÓBCE<sup>(1)</sup>

Założmy, że populacja generalna jest podzielona na skończoną lub przeliczalną ilość klas. Założmy dalej, że istnieje przepis pozwalający rozstrzygnąć, czy dwa elementy dowolnie wybrane z populacji należą do różnych klas, czy nie. Z populacji pobieramy, zgodnie ze schematem Bernoulliego, próbkę o liczebności  $n$ . Na ogół nie wszystkie klasy będą reprezentowane w próbce. Jak ocenić frekwencję w populacji elementów należących do klas w próbce niereprezentowanych? Podamy klasę statystyk, które będą nieobciążonymi estymatorami tej frekwencji.

Oznaczmy przez  $v_{k,n}^{(i)}$  zmienną losową zdefiniowaną w następujący sposób:  $v_{k,n}^{(i)} = 1$  wtedy i tylko wtedy, gdy w danej próbce o liczebności  $n$   $i$ -ta klasa jest reprezentowana  $k$  razy. W pozostałych przypadkach  $v_{k,n}^{(i)} = 0$ . Jeżeli przez  $p_i$  oznaczymy nieznaną frekwencję  $i$ -tej klasy w populacji, to wartość oczekiwana  $E(v_{k,n}^{(i)})$  będzie równa

$$(1) \quad E(v_{k,n}^{(i)}) = P(v_{k,n}^{(i)} = 1) = \binom{n}{k} p_i^k (1 - p_i)^{n-k}.$$

Oznaczmy przez  $v_{k,n}$  sumę  $\sum_i v_{k,n}^{(i)}$ ; wtedy

$$(2) \quad E(v_{k,n}) = \sum_i E(v_{k,n}^{(i)}) = \binom{n}{k} \sum_i p_i^k (1 - p_i)^{n-k}.$$

Zdefiniujmy zmienną losową  $\mu_{k,n}^{(i)}$ ;  $\mu_{k,n}^{(i)} = p_i$ , gdy w danej próbce o liczebności  $n$   $i$ -ta klasa jest reprezentowana  $k$  razy, i  $\mu_{k,n}^{(i)} = 0$  w pozostałych przypadkach. Zmienna  $\mu_{k,n} = \sum_i \mu_{k,n}^{(i)}$  jest frekwencją elementów

---

<sup>(1)</sup> Frekwencja elementów w populacji należących do klas niereprezentowanych w próbce jest zmienną losową. Dlatego też termin *estymator* oznaczający statystykę pozwalającą ocenić wartość tej frekwencji jest niepoprawny. Anglicy nazywają taką statystykę *predict function*. Może należałoby wymyślić nowy termin, np. *predyktor*. Termin *prognoza*, spotykany w polskich publikacjach, uważam w danym przypadku za niezręczny.

w populacji należących do klas  $k$  razy reprezentowanych w próbce. Jej wartość oczekiwana jest równa

$$(3) \quad E(\mu_{k,n}) = \sum_i E(\mu_{k,n}^{(i)}) = \sum_i p_i P(\mu_{k,n}^{(i)} = p_i) = \binom{n}{k} \sum_i p_i^{k+1} (1-p_i)^{n-k}.$$

Z wzoru (2) wynika

$$(4) \quad \frac{\binom{m}{s}}{\binom{N+m+1}{r+s+1}} E(v_{r+s+1, N+m+1}) = \sum_i p_i^{r+1} (1-p_i)^{N-r} \left[ \binom{m}{s} p_i^s (1-p_i)^{m-s} \right],$$

gdzie  $m \geq s$ . Połóżmy

$$(5) \quad \kappa_{r,N}^{(m)} = \sum_{s=0}^m \frac{\binom{N}{r} \binom{m}{s}}{\binom{N+m+1}{r+s+1}} v_{r+s+1, N+m+1}.$$

Z wzoru (4) otrzymamy

$$\begin{aligned} E(\kappa_{r,N}^{(m)}) &= \binom{N}{r} \sum_{s=0}^m \sum_i p_i^{r+1} (1-p_i)^{N-r} \left[ \binom{m}{s} p_i^s (1-p_i)^{m-s} \right] = \\ &= \binom{N}{r} \sum_i p_i^{r+1} (1-p_i)^{N-r} \left[ \sum_{s=0}^m \binom{m}{s} p_i^s (1-p_i)^{m-s} \right] = \\ &= \binom{N}{r} \sum_i p_i^{r+1} (1-p_i)^{N-r}, \end{aligned}$$

a więc ostatecznie

$$(6) \quad E(\kappa_{r,N}^{(m)}) = E(\mu_{r,N}) \quad (m = 0, 1, 2, \dots).$$

Zmienna losowa  $\kappa_{r,N}^{(m)}$  jest kombinacją liniową zmiennych  $v_{r+s+1, N+m+1}$ ; to jest ilości klas mających w próbce o licznosci  $N+m+1$ ,  $r+s+1$  reprezentantów. Zmienna  $\mu_{r,N}$  jest łączną frekwencją elementów należących do klas mających  $r$  reprezentantów w próbce o licznosci  $N$ . Gdy  $m \ll N$ , zmienna losowa może służyć do oceny frekwencji elementów, należących do klas  $r$  razy reprezentowanych w danej próbce o licznosci  $N$ . W szczególności statystyka

$$\kappa_{0,N}^{(m)} = \sum_{s=0}^m \frac{\binom{m}{s}}{\binom{N+m+1}{s+1}} v_{s+1, N+m+1} \quad (m = 0, 1, 2, \dots)$$

może służyć do oceny frekwencji elementów niereprezentowanych w próbie. Przypuśćmy np., że z populacji słów zawartych w polskich dziennikach z 1957 roku pobrano, zgodnie ze schematem Bernoulliego, próbkę o liczności 10000. Aby ocenić frekwencję słów, których w próbie nie spotkano, losuje się np. 5 dodatkowych słów z populacji ( $m = 4$ ) i oblicza się wartość statystyki  $\kappa_{0,10000}^{(4)}$ . Wzór (6) powiada nam, że statystyka  $\kappa_{0,10000}^{(4)}$  jest nieobciążonym estymatorem wartości oczekiwanej tej frekwencji.

Dla  $r = 0$  i  $m = 0, 1$  mamy

$$\kappa_{0,N}^{(0)} = \frac{1}{N+1} \nu_{1,N+1},$$

$$\kappa_{0,N}^{(1)} = \frac{1}{N+2} \nu_{1,N+2} + \frac{2}{(N+1)(N+2)} \nu_{2,N+2}.$$

Ogólnie biorąc, przy małych  $r$  i  $N \gg m$  współczynniki przy  $\nu_{r+s+1, N+m+1}$  we wzorze (5) szybko maleją, gdy  $s \rightarrow m$ .

Można obliczyć wariancję zmiennej  $\kappa_{r,N}^{(m)}$ . Wystarczy skorzystać z wzorów

$$E(\nu_{k,n}^{(i)})^2 = \binom{n}{k} p_i^k (1-p_i)^{n-k},$$

$$E(\nu_{k,n}^{(i)} \nu_{l,n}^{(i)}) = 0, \quad \text{gdy } k \neq l,$$

$$E(\nu_{k,n}^{(i)} \nu_{l,n}^{(j)}) = \frac{n!}{k!l!(n-k-l)!} p_i^k p_j^l (1-p_i-p_j)^{n-k-l} \quad (i \neq j)$$

i rozbić  $E[\kappa_{r,N}^{(m)} - E(\kappa_{r,N}^{(m)})]^2$  na kombinację liniową wartości oczekiwanych wyrażeń  $\nu_{k,n}^{(i)} \nu_{l,n}^{(j)}$ . Otrzymany rezultat jest jednak bardzo skomplikowany i nie będę go tutaj podawał.

Statystyki  $\kappa_{r,N}^{(m)}$  mają tę zaletę, że przy ich użyciu zbędna jest informacja o ilości klas w populacji. Często, jak np. w podanym wyżej przykładzie, informacji takiej brak.

Zmienne losowe  $\nu_{k,n}$  były badane przez różnych autorów, np. Gooda [1]. Posłużyły one, między innymi, do oceny frekwencji  $p_i$  oraz do oceny ilości klas w próbie o liczności  $\alpha N$  ( $\alpha > 1$ ), gdy dana jest próbka o liczności  $N$ .

#### Praca cytowana

[1] I. J. Good, *The population frequencies of species and the estimation of population parameters*, Biometrika 40 (1953).

INSTYTUT MATEMATYCZNY POLSKIEJ AKADEMII NAUK

Praca wpłynęła 11. 4. 1958

С. Т Р Ы Б У Л А (Вроцлав)

**ОЦЕНКА ЧАСТОТЫ В СОВОКУПНОСТИ ЭЛЕМЕНТОВ,  
ПРИНАДЛЕЖАЩИХ КЛАССАМ, НЕ ИМЕЮЩИМ  
ПРЕДСТАВИТЕЛЕЙ В ВЫБОРКЕ**

## РЕЗЮМЕ

Если из совокупности, разделенной на конечное или счетное число классов, выберем согласно схеме Бернулли выборку, состоящую из  $N$  элементов, то ожидаемое значение частоты элементов, принадлежащих классам имеющим по  $r$  представителей в выборке, будет равно

$$m_{r,N} = \binom{N}{r} \sum_i p_i^{r+1} (1-p_i)^{N-r} \quad (r = 0, 1, 2, \dots),$$

где  $p_i$  есть (неизвестная) частота в совокупности элементов, принадлежащих  $i$ -му классу. Обозначим через  $v_{r,N}$  число классов с  $r$  представителями в выборке, состоящей из  $N$  элементов. Доказывается, что

$$x_{r,N}^{(m)} = \binom{N}{r} \sum_{\delta=0}^m \frac{\binom{m}{\delta}}{\binom{N+m+1}{r+\delta+1}} v_{r+\delta+1, N+m+1} \quad (m = 0, 1, 2, \dots)$$

являются несмещенными оценками параметра  $m_{r,N}$ . Приводятся возможные приложения этого результата, в частности к оценке частоты элементов, принадлежащих к классам, которые не имеют представителей в выборке. Заметим, что для составления статистик  $x_{r,N}^{(m)}$  не нужны сведения о числе классов в совокупности.

S. T R Y B U Ł A (Wrocław)

**THE ESTIMATION OF FREQUENCY IN A POPULATION OF ELEMENTS  
BELONGING TO CLASSES NOT REPRESENTED IN THE SAMPLE**

## SUMMARY

If we take a sample of size  $N$  according to the Bernoulli scheme from a population divided into a finite or denumerable number of classes, then the expected value of the frequency, in the population, of elements belonging to classes represented  $r$  times in the sample will be

$$m_{r,N} = \binom{N}{r} \sum_i p_i^{r+1} (1-p_i)^{N-r} \quad (r = 0, 1, 2, \dots),$$

where  $p_i$  is the (unknown) frequency of the elements belonging to the  $i$ -th class in the population. Denote by  $v_{r,N}$  the number of classes represented  $r$  times in a given sample of size  $N$ . It can be shown that

$$\kappa_{r,N}^{(m)} = \binom{N}{r} \sum_{\delta=0}^m \frac{\binom{m}{\delta}}{\binom{N+m+1}{r+\delta+1}} v_{r+\delta+1, N+m+1} \quad (m = 0, 1, 2, \dots)$$

are unbiased estimators of the parameter  $m_{r,N}$ .

The possible applications of this result are given, in particular to the estimation of frequency of elements belonging to classes which are not represented in the sample.

It will be observed that in order to form the statistics  $\kappa_{r,N}^{(m)}$  no information concerning the number of classes in the population is needed.

---