

CZEN PIN i DZAN DZO-I (Wrocław)

O ESTYMACJI LICZNOŚCI POPULACJI ZA POMOCĄ METODY
ŁOWIENIA I ZNAKOWANIA

Wstęp. W naukach przyrodniczych często spotykamy się z problemem estymacji liczności populacji. Przykładem może być estymacja liczności ryb określonego gatunku i wielkości w jeziorze, estymacja ilości widzów telewizyjnych w danym kraju itp. Do tego celu może służyć metoda „łowienia i znakowania”. Jest ona obecnie powszechnie znana i używana w różnych dziedzinach nauk przyrodniczych. Metoda ta polega na losowaniu pewnej ilości elementów z całej populacji, znakowaniu elementów wylosowanych i zwracaniu do populacji z powrotem, ponownym losowaniu i zwracaniu itd. Z otrzymanych wyników buduje się estymator liczności całej populacji. Na ogół liczność populacji biologicznej jest zmienna w czasie. Powodem mogą być śmierć i urodzenie się nowych elementów populacji, emigracja i imigracja itd. Aby więc badać takie zagadnienia, trzeba konstruować różne modele statystyczne. Modele takie skonstruowano do różnych celów, o czym można się dowiedzieć np. z pracy Chapmana [1].

W niniejszej pracy zajmiemy się populacją o stałej liczności. W tym przypadku wygodnie jest sformułować zagadnienie korzystając z modelu urny. Wyobraźmy sobie urnę, w której znajduje się nieznaną ilość kul. Oznaczmy tę ilość przez N . Chcemy ją oszacować.

H. Steinhaus zaproponował następujący sposób postępowania: Najpierw losujemy bez zwracania m_1 kul i zaopatrujemy każdą z nich w znaczek. Jest to pierwsza próbka. Po oznakowaniu wrzucamy te kule razem do urny. Następnie losujemy bez zwracania m_2 kul i znaczymy każdą z nich, bez względu na to, czy kula ta już była, czy nie była znakowana. Jest to druga próbka. Wrzucamy kule z powrotem do urny i w ten sam sposób postępujemy dalej aż do n -tej próbki. W ten sposób otrzymamy ciąg liczb V_1, V_2, \dots, V_n , gdzie V_i oznacza ilość kul, na których jest i znaczków. Oczywiście $\sum_{i=1}^n iV_i = \sum_{i=1}^n m_i = m$. Liczbę m nazywamy *licznością próbki*. Chcemy na podstawie V_1, \dots, V_n zbudować estymator parametru N .

Oznaczmy przez V łączną ilość różnych kul wylosowanych w trakcie pobierania n podpróbek. Mamy więc $V = \sum_{i=1}^n V_i$. Powstaje pytanie: Jaki jest rozkład zmiennej losowej V ?

1. Rozkład zmiennej V . Możemy założyć, że wszystkie kule w urnie są ponumerowane. Oznaczmy ich ilość przez N . Możemy zatem mówić o i -tej kuli ($i \leq N$). Oznaczmy przez R_i ilość znaczków na i -tej kuli.

R_i są zmiennymi losowymi i $\sum_{i=1}^n R_i = m$.

Prawdopodobieństwo tego, że pewna kula zostanie wylosowana w i -tym losowaniu j -tej podpróbki, jest równe $1/(N-i+1)$. Ale wszystkie podpróbki są wylosowane niezależnie. A więc

$$(1) \quad P[V = v; R_{k_i} = r_{k_i} (i = 1, 2, \dots, v), R_k = 0 \text{ dla } k \notin \{k_1, \dots, k_v\}] = \\ = C(r_{k_1}, r_{k_2}, \dots, r_{k_v}) \prod_{i=1}^v \left(\frac{1}{N(N-1) \dots (N-m_i+1)} \right),$$

gdzie r_{k_i} ($i = 1, \dots, v$), są liczbami całkowitymi spełniającymi warunki $r_{k_i} > 0$, $\sum_{i=1}^v r_{k_i} = m$, a $C(r_1, r_2, \dots, r_l)$ jest zdefiniowane jako liczba wszystkich możliwych rozmieszczeń l różnych kul w n podpróbkach tak, żeby i -ta kula powtarzała się r_i razy i żeby żadna kula nie powtarzała się w jednej i tej samej podpróbce. Oznaczmy

$$(2) \quad G(v) = \sum_{\substack{r_1 \geq 1, \dots, r_v \geq 1 \\ r_1 + \dots + r_v = m}} C(r_1, r_2, \dots, r_v).$$

Wtedy

$$(3) \quad P[V = v | N] = \\ = \binom{N}{v} \sum_{\substack{r_1 \geq 1, \dots, r_v \geq 1 \\ r_1 + \dots + r_v = m}} P[V = v; R_1 = r_1, \dots, R_v = r_v, R_k = 0 \text{ dla } k > v] = \\ = \binom{N}{v} G(v) \prod_{i=1}^v \left(\frac{1}{N(N-1) \dots (N-m_i+1)} \right).$$

W przypadku $m_1 = m_2 = \dots = m_n = 1$, funkcja $G(v)$ daje się łatwo wyrachować. Mamy wtedy

$$G(v) = \sum_{\substack{r_1 \geq 1, \dots, r_v \geq 1 \\ r_1 + \dots + r_v = n}} \frac{n!}{r_1! r_2! \dots r_v!}.$$

Porównując współczynniki rozwinięcia funkcji

$$(e^x - 1)^v = \sum_{i=0}^v (-1)^i \binom{v}{i} e^{(v-i)x}$$

w szereg potęgowy, otrzymujemy

$$G(v) = \sum_{i=0}^v (-1)^i \binom{v}{i} (v-i)^n.$$

A więc ostatecznie

$$(4) \quad P[V = v | N] = \binom{N}{v} \left(\frac{1}{N}\right)^n \sum_{i=0}^v (-1)^i \binom{v}{i} (v-i)^n.$$

W przypadku, gdy równość $m_1 = m_2 = \dots = m_n = 1$ nie zachodzi, znalezienie dokładnego rozkładu zmiennej losowej V jest trudne. Połóżmy $L = m - V$. W następnej sekcji znajdziemy asymptotyczny rozkład zmiennej L , gdy $m \rightarrow \infty$ i $m/N \rightarrow 0$.

2. Graniczny rozkład zmiennej losowej L .

Z wzoru (3) wynika

$$\begin{aligned} P[L = l | N] &= P[V = m - l | N] = \\ &= \binom{N}{m-l} G(m-l) \prod_{i=1}^n \left(\frac{1}{N(N-1) \dots (N-m_i+1)} \right). \end{aligned}$$

Udowodnimy najpierw

LEMAT 1.

(a) Jeśli $m / \sum_{i \neq j}^n m_i m_j \rightarrow 0$, to

$$O(\underbrace{2, 2, \dots, 2}_{l \text{ razy}}) = \left(\sum_{i=1}^{n-1} \sum_{j=i+1}^n m_i m_j \right)^l [1 + O(m / \sum_{i \neq j}^n m_i m_j)];$$

(b), jeśli $\sum_{j=1}^k r_j = l + k$ i $r_j \geq 2$ ($j = 1, 2, \dots, k$), to

$$O(r_1, r_2, \dots, r_k) \leq \left(\sum_{i=1}^{n-1} \sum_{j=i+1}^n m_i m_j \right)^{l-1} m.$$

Dowód. (a) Z definicji $O(r_1 = 2, \dots, r_l = 2)$ jest ilością wszystkich możliwych rozmieszczeń l różnych liczb w n podpróbce tak, żeby każda liczba powtarzała się dwa razy i żeby żadna liczba nie powtarzała się w jednej i tej samej podpróbce. Jest rzeczą jasną, że liczbę 1 możemy

umieścić w $\sum_{i=1}^{n-1} \sum_{j=i+1}^n m_i m_j$ różnych sposobów. Po umieszczeniu liczby 1, możemy umieścić liczbę 2 w

$$\begin{aligned} \sum_{i=1}^{n-1} \sum_{j=i+1}^n (m_i - \varepsilon_i^1)(m_j - \varepsilon_j^1) &= \sum_{i=1}^{n-1} \sum_{j=i+1}^n m_i m_j - \sum_{i \neq j} \varepsilon_i^1 m_j + \sum_{i=1}^{n-1} \sum_{j=i+1}^n \varepsilon_i^1 \varepsilon_j^1 = \\ &= \sum_{i=1}^{n-1} \sum_{j=i+1}^n m_i m_j \left(1 + O\left(m / \sum_{j \neq i}^n m_i m_j\right)\right) \end{aligned}$$

różnych sposobów, gdzie

$$\varepsilon_i^1 = \begin{cases} 1, & \text{jeżeli w } i\text{-tej podpróbce znajduje się liczba 1,} \\ 0, & \text{jeżeli tak nie jest.} \end{cases}$$

Podobnie, jeśli mamy już rozmieszczone liczby $1, 2, \dots, k-1$, to liczbę k możemy umieścić w

$$(5) \quad \sum_{i=1}^{n-1} \sum_{j=i+1}^n (m_i - \varepsilon_i^1 - \varepsilon_i^2 - \dots - \varepsilon_i^{k-1})(m_j - \varepsilon_j^1 - \varepsilon_j^2 - \dots - \varepsilon_j^{k-1}) = \\ = \sum_{i=1}^{n-1} \sum_{j=i+1}^n m_i m_j \left(1 + O\left(m / \sum_{i \neq j}^n m_i m_j\right)\right)$$

różnych sposobów, gdzie

$$\varepsilon_i^t = \begin{cases} 1, & \text{jeżeli w } i\text{-tej podpróbce znajduje się liczba } t, \\ 0, & \text{jeżeli tak nie jest} \end{cases}$$

($t = 1, 2, \dots, k-1$; $i = 1, 2, \dots, n$). Z wzoru (5) wynika pierwsza część lematu.

(b) Przejdźmy do dowodu drugiej części lematu. Rozpatrzmy najpierw przypadek $l = k+1$. Wówczas z założenia, że $r_j \geq 2$ ($j = 1, 2, \dots, k$) wynika, że $\max_{1 \leq i \leq k} r_i = 3$. Bez naruszenia ogólności możemy założyć, że $r_1 = 3$, a więc $r_2 = r_3 = \dots = r_k = 2$. W tym przypadku możemy liczbę 1 umieścić w $\sum_{i=1}^{n-2} \sum_{j=i+1}^{n-1} \sum_{s=j+1}^n m_i m_j m_s$ różnych sposobów, a żadnej z pozostałych $k-1$ liczb nie możemy umieścić w więcej niż $\sum_{i=1}^{n-1} \sum_{j=i+1}^n m_i m_j$ sposobów. Mamy więc

$$(6) \quad C(3, \underbrace{2, \dots, 2}_{l-2 \text{ razy}}) \leq \left(\sum_{i=1}^{n-2} \sum_{j=i+1}^{n-1} \sum_{s=j+1}^n m_i m_j m_s\right) \left(\sum_{i=1}^{n-1} \sum_{j=i+1}^n m_i m_j\right)^{l-2} \leq \\ \leq m \left(\sum_{i=1}^{n-1} \sum_{j=i+1}^n m_i m_j\right)^{l-1}.$$

Dla dowolnego układu (r_1, r_2, \dots, r_k) spełniającego warunki $\sum_{i=1}^k r_i = l+k$, $r_j \geq 2$ ($j = 1, 2, \dots, k$), mamy

$$(7) \quad C(r_1, r_2, \dots, r_k) \leq C(\underbrace{3, 2, \dots, 2}_{l-2 \text{ razy}}).$$

Z wzorów (6) i (7) wynika druga część lematu.

LEMAT 2. Jeśli l jest ustalone i $m / \sum_{i \neq j}^n m_i m_j \rightarrow 0$, to

$$(8) \quad G(m-l) = \frac{(m-l)!}{l!} \left(\sum_{i=1}^{n-1} \sum_{j=i+1}^n m_i m_j \right)^l \left(1 + O \left(m / \sum_{i \neq j}^n m_i m_j \right) \right).$$

Dowód. Udowodnimy najpierw, że

$$(9) \quad G(v) = \begin{cases} \sum_{i=1}^{m-v} \frac{v!}{i!} \sum_{\substack{r_1 \geq 2, \dots, r_i \geq 2 \\ r_1 + \dots + r_i = m-v+i}} C(r_1, r_2, \dots, r_i) & \text{dla } v < m, \\ v! & \text{dla } v = m. \end{cases}$$

Z definicji liczba $C(r_1, \dots, r_l)$ oznacza ilość wszystkich możliwych rozmieszczeń l różnych liczb w n podpróbkach tak, żeby liczba i powtarzała się r_i razy i żeby żadna liczba nie powtarzała się w jednej i tej samej podpróbce, natomiast

$$G(v) = \sum_{\substack{r_1 \geq 1, \dots, r_v \geq 1 \\ r_1 + \dots + r_v = m}} C(r_1, r_2, \dots, r_v).$$

Jeżeli więc $v = m$, to $G(v) = C(r_1 = 1, \dots, r_m = 1) = m!$. Aby udowodnić równość (9) dla $v \leq m$, zauważmy, że

$$C(r_1, \dots, r_{v-i}, \underbrace{1, \dots, 1}_{i \text{ razy}}) = i! C(r_1, \dots, r_{v-i}),$$

i że

$$C(r_1, \dots, r_i, \dots, r_j, \dots, r_v) = C(r_1, \dots, r_j, \dots, r_i, \dots, r_v) \quad (i, j = 1, 2, \dots, v).$$

Załóżmy teraz, że wśród liczb r_1, r_2, \dots, r_v jest $v-i$ jedynek. Pozostałych i liczb, oraz $v-i$ jedynek można rozmieścić na miejscach $1, 2, \dots, v$

na $\binom{v}{i}$ sposobów. Stąd więc wynika, że

$$\begin{aligned}
 (10) \quad u(i) &\stackrel{\text{def}}{=} \sum C(r_1, r_2, \dots, r_v) = \\
 &= \binom{v}{i} \sum_{\substack{s_1 \geq 2, \dots, s_i \geq 2 \\ s_1 + s_2 + \dots + s_i = m - v + i}} C(s_1, s_2, \dots, s_i, \underbrace{1, \dots, 1}_{v-i \text{ razy}}) = \\
 &= \binom{v}{i} \sum_{\substack{s_1 \geq 2, \dots, s_i \geq 2 \\ s_1 + \dots + s_i = m - v + i}} (v-i)! C(s_1, s_2, \dots, s_i) = \frac{v!}{i!} \sum_{\substack{s_1 \geq 2, \dots, s_i \geq 2 \\ s_1 + \dots + s_i = m - v + i}} C(s_1, s_2, \dots, s_i),
 \end{aligned}$$

gdzie pierwsza suma rozciąga się na wszystkie te układy (r_1, \dots, r_v) ($r_j \geq 1$, $\sum_{j=1}^n r_j = m$), które zawierają dokładnie $v-i$ jedynek. Korzystając z definicji funkcji $G(v)$ i z wzoru (10) otrzymujemy ostatecznie wzór (9). Podstawiając we wzorze (9) $v = m-l$ i korzystając z lematu 1 otrzymamy tezę lematu 2.

LEMAT 3. *Jeśli $m \rightarrow \infty$ i*

$$\begin{aligned}
 \text{(i)} \quad \frac{m^2}{\sum_{i \neq j} m_i m_j} &\leq k, & \text{(ii)} \quad \frac{\sum_{i=1}^{n-1} \sum_{j=i+1}^n m_i m_j}{N} &= \lambda,
 \end{aligned}$$

to

$$\frac{N!}{(N-m+l)!} \prod_{i=1}^n \left(\frac{1}{N(N-1)\dots(N-m_i+1)} \right) = \frac{1}{N^l} e^{-\lambda} \left[1 + O\left(\frac{m}{N}\right) \right].$$

Dowód. Z nierówności

$$\sum_{j=1}^{m_i} \log(N-j) \leq \int_{N-m_i}^N \log x dx \leq \sum_{j=0}^{m_i-1} \log(N-j)$$

i wzoru

$$\int_{N-m_i}^N \log x dx = \log \frac{N^N}{(N-m_i)^{N-m_i}} - m_i$$

wynika, że

$$\log(N(N-1)\dots(N-m_i+1)) = \log \frac{N^N}{(N-m_i)^{N-m_i}} - m_i + O\left(\frac{m_i}{N}\right).$$

Mamy zatem

$$\begin{aligned}
 (11) \quad \sum_{i=1}^n \log(N(N-1)\dots(N-m_i+1)) &= \\
 &= \log \frac{N^{nN}}{(N-m_1)^{N-m_1} \dots (N-m_n)^{N-m_n}} - m + O\left(\frac{m}{N}\right)
 \end{aligned}$$

oraz

$$(12) \quad \log \frac{N!}{(N-m+l)!} = \log \frac{N^N}{(N-m+l)^{N-m+l}} - m+l + O\left(\frac{m}{N}\right).$$

A wiec na podstawie wzorow (11) i (12)

$$\begin{aligned} (13) \quad & \log \left[\frac{N!}{(N-m+l)!} \prod_{i=1}^n \frac{1}{N(N-1)\dots(N-m_i+1)} \right] = \\ & = -\log N^l + \sum_{i=1}^n (N-m_i) \log \left(1 - \frac{m_i}{N}\right) - (N-m+l) \log \left(1 - \frac{m-l}{N}\right) + \\ & \quad + l + O\left(\frac{m}{N}\right) = \\ & = -l \log N + \sum_{i=1}^n (N-m_i) \left(-\frac{m_i}{N} - \frac{m_i^2}{2N^2}\right) - \\ & \quad - (N-m+l) \left(-\frac{m-l}{N} - \frac{(m-l)^2}{2N^2}\right) + l + O\left(\frac{m}{N}\right) = \\ & = -l \log N + \frac{1}{2N} \sum_{i=1}^n m_i^2 - \frac{1}{2N} \left(\sum_{i=1}^n m_i\right)^2 + O\left(\frac{m}{N}\right) = \\ & = -l \log N - \frac{1}{N} \sum_{i=1}^{n-1} \sum_{j=i+1}^n m_i m_j + O\left(\frac{m}{N}\right), \end{aligned}$$

c. b. d. o.

Twierdzenie 1. *Jeżeli $m \rightarrow \infty$ w ten sposob, że*

$$(i) \quad \frac{m^2}{\sum_{i \neq j}^n m_i m_j} \leq k = \text{const}, \quad (ii) \quad \frac{\sum_{i=1}^{n-1} \sum_{j=i+1}^n m_i m_j}{N} = \lambda,$$

to

$$(14) \quad P[L = l | N] = \frac{\lambda^l}{l!} e^{-\lambda} \left[1 + O\left(\frac{m}{N}\right)\right].$$

Dowód. Twierdzenie jest prostą konsekwencją lematow 2 i 3.

3. Estymator \bar{N} liczności N . W niniejszym ustępie będziemy cały czas zakładali, że warunki twierdzenia 1 są spełnione. Łatwo wtedy wykazać, że

$$(15) \quad E\left(\frac{1}{L+1}\right)^k \rightarrow \sum_{l=0}^{\infty} \frac{1}{(l+1)^k} \cdot \frac{\lambda^l}{l!} e^{-\lambda}, \quad k > 0.$$

Dla $k = 1$ mamy

$$(16) \quad E\left(\frac{1}{L+1}\right) \approx \sum_{l=0}^{\infty} \frac{\lambda}{(l+1)!} e^{-\lambda} = \frac{1}{\lambda} (1 - e^{-\lambda}) = \\ = \frac{N}{\sum_{i=1}^{n-1} \sum_{j=i+1}^n m_i m_j} \left[1 - e^{-\left(\sum_{i=1}^{n-1} \sum_{j=i+1}^n m_i m_j\right)/N} \right].$$

Zdefiniujmy estymator \bar{N} :

$$(17) \quad \bar{N} \stackrel{\text{def}}{=} \frac{\sum_{i=1}^{n-1} \sum_{j=i+1}^n m_i m_j}{L+1}.$$

D. G. Chapman ([2]) znalazł, dla tego samego postępowania, metodą największej wiarygodności estymator

$$\bar{N}_c = \frac{\sum_{i=1}^{n-1} \sum_{j=i+1}^n m_i m_j}{L}.$$

Nie znalazł on jednak rozkładu zmiennej losowej L i nie mógł zbadać głębiej własności tego estymatora.

Zajmijmy się bliżej estymatorem \bar{N} . W przybliżeniu

$$(18) \quad E(\bar{N}) \approx N(1 - e^{-\lambda}) = N \left[1 - e^{-\left(\sum_{i=1}^{n-1} \sum_{j=i+1}^n m_i m_j\right)/N} \right].$$

Ponadto z równości

$$\frac{1}{(l+1)^2} = \sum_{r=2}^R \frac{(r-2)!}{(l+1) \dots (l+r)} + \frac{(R-1)!}{(l+1)^2 (l+2) \dots (l+R)}$$

oraz

$$E\left[\frac{1}{(l+1) \dots (l+r)}\right] = \sum_{l=0}^{\infty} \frac{\lambda^l}{(l+r)!} e^{-\lambda} = \\ = \frac{1}{\lambda^r} \left[1 - e^{-\lambda} \left(1 + \lambda + \dots + \frac{\lambda^r}{(r-1)!} \right) \right],$$

otrzymujemy wzór

$$(19) \quad E(\bar{N} - N)^2 = N^2 \left(\frac{1}{\lambda} + \frac{2}{\lambda^2} + \frac{6}{\lambda^3} + O\left(\frac{1}{\lambda^4}\right) \right) =$$

$$= N^2 \left[\frac{N}{\sum_{i=1}^{n-1} \sum_{j=i+1}^n m_i m_j} + \frac{2N^2}{\left(\sum_{i=1}^{n-1} \sum_{j=i+1}^n m_i m_j\right)^2} + \frac{6N^3}{\left(\sum_{i=1}^{n-1} \sum_{j=i+1}^n m_i m_j\right)^3} + \right.$$

$$\left. + O\left(N^4 / \left(\sum_{i=1}^{n-1} \sum_{j=i+1}^n m_i m_j\right)^4\right) \right].$$

Z wzoru (18) widzimy, że estymator \bar{N} nie jest nieobciążony, ale dla dużych λ obciążenie tego estymatora jest bardzo małe. Z wzoru (19) wynika, że jeśli chcemy, żeby średni błąd kwadratowy tego estymatora był mały, musimy powiększyć licznosc próbki tak, żeby λ było duże, np. rzędu kilkuset.

Łatwo również obliczyć, że dla $\left(\sum_{i=1}^{n-1} \sum_{j=i+1}^n m_i m_j\right)/N \geq \log N$ zachodzi nierówność $|E(\bar{N}) - N| \leq 1$.

Dla przypadku $m_1 = m_2 = m_3 = \dots = m_n = k$ wzór (17) sprowadza się do

$$(20) \quad \bar{N}_{n,k} = \frac{k^2 n(n-1)}{2(L+1)}$$

i

$$(21) \quad E(\bar{N}_{n,k} - N)^2 = N^2 \left[\frac{2N}{k^2 n(n-1)} + \frac{8N^2}{k^4 n^2(n-1)^2} + O\left(\frac{N^3}{k^6 n^3(n-1)^3}\right) \right].$$

Z wzoru (21) widzimy, że dla ustalonej licznosci próbki kn wyrażenie $E(\bar{N}_{n,k} - N)^2$ jest tym mniejsze, im mniejsza jest liczba k . Optymalny przypadek jest dla $k = 1$. Mamy wtedy

$$\bar{N}_{n,1} = \frac{n(n-1)}{2(L+1)}$$

oraz

$$(22) \quad E(\bar{N}_{n,1} - N)^2 = N^2 \left[\frac{2N}{n(n-1)} + \frac{8N^2}{n^2(n-1)^2} + O\left(\frac{N^3}{n^3(n-1)^3}\right) \right].$$

W najprostszym przypadku, gdy $n = 2$, $\bar{N} = m_1 m_2 / (L+1)$. Przypuśćmy, że $m_1 = m_2 = k$. Mamy wtedy

$$\bar{N}_{2,k} = \frac{k^2}{L+1}$$

oraz

$$(23) \quad E(\bar{N}_{2,k} - N)^2 = N^2 \left[\frac{N}{k^2} + \frac{2N^2}{k^4} + O\left(\frac{N^3}{k^6}\right) \right].$$

Dla tej samej liczności próbki $2k$

$$(24) \quad E(\bar{N}_{2k,1} - N)^2 = N^2 \left[\frac{N}{k(2k-1)} + \frac{2N^2}{k^2(2k-1)^2} + O\left(\frac{N^3}{k^3(2k-1)^3}\right) \right]$$

otrzymujemy zatem

$$(25) \quad \frac{E(\bar{N}_{2,k} - N)^2}{E(\bar{N}_{2k,1} - N)^2} \approx 2.$$

Z wzoru (25) widzimy, że na tym samym poziomie dokładności (tzn. tej samej dyspersji) w postępowaniu dwupróbkowym trzeba będzie wylosować 2 razy więcej kul niż w postępowaniu pojedynczym.

4. Przedziałowa estymacja liczności N . W praktyce na ogół nie chodzi nam o estymację punktową liczności populacji. Wolelibyśmy orzekać, że liczność danej populacji jest mniej więcej tyle a tyle, niż mówić, że liczność danej populacji jest np. sto, sto jeden, sto pięć, a nie inna. Innymi słowy, chcemy znaleźć taki przedział o końcach będących zmiennymi losowymi, o którym możemy z jakimś współczynnikiem ufności powiedzieć, że zawiera liczbę N .

Z twierdzenia 1 i z centralnego twierdzenia granicznego wynika, że dla $\lambda \rightarrow \infty$

$$(26) \quad P\left[\frac{(L-\lambda)^2}{\lambda} \leq a^2 \mid N\right] = P\left[\left|\frac{L-\lambda}{\sqrt{\lambda}}\right| \leq a \mid N\right] \rightarrow \Phi(a) - \Phi(-a),$$

gdzie

$$\Phi(a) = \frac{1}{\sqrt{2\pi}} \int_{-\infty}^a e^{-t^2/2} dt.$$

Z wzoru (26) wnioskujemy, że jeśli

$$(27) \quad \begin{aligned} A &= \frac{1}{2} [2L + a^2 + a\sqrt{4L + a^2}], \\ B &= \frac{1}{2} [2L + a^2 - a\sqrt{4L + a^2}], \end{aligned}$$

to

$$P[B \leq \lambda \leq A \mid N] \approx \Phi(a) - \Phi(-a).$$

A więc

$$(28) \quad P\left[\frac{\sum_{i=1}^{n-1} \sum_{j=i+1}^n m_i m_j}{A} \leq N \leq \frac{\sum_{i=1}^{n-1} \sum_{j=i+1}^n m_i m_j}{B} \mid N\right] \approx \Phi(a) - \Phi(-a).$$

Dla $a = 2$ i dla poziomu ufności 0,95 mamy

$$(29) \quad P \left[\frac{\sum_{i=1}^{n-1} \sum_{j=i+1}^n m_i m_j}{(\sqrt{L+1}+1)^2} \leq N \leq \frac{\sum_{i=1}^{n-1} \sum_{j=i+1}^n m_i m_j}{(\sqrt{L+1}-1)^2} \right] = 0,95.$$

Wtedy długość przedziału

$$\left(\frac{\sum_{i=1}^{n-1} \sum_{j=i+1}^n m_i m_j}{(\sqrt{L+1}+1)^2}, \frac{\sum_{i=1}^{n-1} \sum_{j=i+1}^n m_i m_j}{(\sqrt{L+1}-1)^2} \right)$$

jest równa

$$(30) \quad \frac{\sum_{i=1}^{n-1} \sum_{j=i+1}^n m_i m_j}{L^2} 4\sqrt{L+1} \approx$$

$$\approx \frac{4 \sum_{i=1}^{n-1} \sum_{j=i+1}^n m_i m_j}{L^{3/2}} = \frac{4 \sum_{i=1}^n m_i}{L} \left(\frac{\sum_{i=1}^{n-1} \sum_{j=i+1}^n m_i m_j}{\sum_{i=1}^n m_i \sqrt{L}} \right).$$

W przypadku $m_1 = m_2 = \dots = m_n = k$

$$P \left[\frac{k^2 n(n-1)}{2(\sqrt{L+1}+1)^2} \leq N \leq \frac{k^2 n(n-1)}{2(\sqrt{L+1}-1)^2} \right]$$

i długość przedziału $(k^2(n-1)n/2(\sqrt{L+1}+1)^2, k^2(n-1)n/2(\sqrt{L+1}-1)^2)$ jest równa

$$\frac{2k^2 n(n-1)}{L^{3/2}} = \frac{nk}{L} \cdot \frac{2k(n-1)}{\sqrt{L}}.$$

Przykład. Jeżeli dla $k = 1$, $n = 2000$ otrzymaliśmy $L = 200$, to z warunku

$$P[8676 \leq N \leq 11515] = 0,95$$

długość przedziału jest $11515 - 8676 = 2839$.

Dziękujemy dr Stanisławowi Trybule za pomoc w przygotowaniu ostatecznej wersji niniejszej pracy.

Prace cytowane

[1] D. G. Chapman, *The estimation of biological populations*, Ann. Math. Statist. 125 (1954), str. 1-15.

[2] — *Inverse, multiple and sequential sample censuses*, Biometrics 8 (1952), str. 288-306.

Praca wpłynęła 12. 4. 1960

ЧЕНЬ ПИНЬ и ДЗАН ДЗО-И (Вроцлав)

ОБ ОЦЕНКЕ ЧИСЛЕННОСТИ СОВОКУПНОСТИ МЕТОДОМ
ЛОВЛИ И ОТМЕТОК

РЕЗЮМЕ

В статье рассматривается следующая схема выбора: Выбирается n подвыборок численностью соответственно m_1, m_2, \dots, m_n . В каждой из них выбор производится без возвращений. Перед выбором $(k+1)$ -ой подвыборки отмечаются и возвращаются в совокупность все элементы предыдущей подвыборки.

Пусть $m = \sum_{i=1}^n m_i$ и пусть V обозначает число разных элементов в выборке (т.е. во всех n подвыборках вместе). В статье доказывается, что случайная величина $L = m - V$ имеет в пределе распределение Пуассона. Пользуясь этой теоремой мы исследуем свойства оценки

$$\bar{N} = \frac{\sum_{i=1}^{n-1} \sum_{j=i+1}^n m_i m_j}{L+1}$$

численности N совокупности. Дается также способ оценки этой численности методом доверительных пределов.

CZEN PIN and DZAN DZO-I (Wrocław)

ON ESTIMATING THE SIZE OF POPULATION BY CAPTURE-MARK
METHOD

SUMMARY

The following sampling scheme is considered in this paper. We draw n subsamples of sizes m_1, \dots, m_n , each one without replacement. Before drawing the $(k+1)$ -st subsample we put back to the population all the elements of k -th subsample.

Let $m = \sum_{i=1}^n m_i$, and let V denote the number of distinct elements in the whole sample (i.e. in n subsamples altogether). It is proved here that the

random variable $L = m - V$ has the limiting Poisson distribution. Using this theorem the properties of the estimate

$$\bar{N} = \frac{\sum_{i=1}^{n-1} \sum_{j=i+1}^n m_i m_j}{L+1}$$

of the population size N are investigated. The paper also gives the method of construction of confidence intervals for the population size.
