

## STUDIAREA PROBABILITĂȚII $\pi$ ÎN METODA RĂSPUNSURILOR RANDOMIZATE

*Andrei POȘTARU, Nicolae PRODAN, Olga BENDERSCHI*

*Universitatea de Stat din Moldova*

Este examinată probabilitatea  $\pi$  de a obține răspunsuri sincere în sondaje. Sunt expuse câteva metode de estimare a lui  $\pi$ : metoda intervalelor de încredere, metoda verosimilității maxime, verificarea ipotezelor statistice simple etc.

**Cuvinte-cheie:** sondaje, probabilitate, statistică, estimație, ipoteze statistice, interval de încredere, răspunsuri randomizate.

### **RANDOMIZED RESPONSE METHOD: ESTIMATIONS OF THE PROBABILITY $\pi$**

The probability  $\pi$  to get honest answers in surveys is considered. There are presented several methods for  $\pi$  estimating: confidence interval method, maximum likelihood estimation, statistical hypotheses testing etc.

**Keywords:** surveys, probability, statistics, estimation, statistical hypotheses, confidence interval, randomized response.

### **Introducere**

Rezultatele cercetărilor sociologice, de regulă, conțin erori, uneori foarte grave, care pot reduce la zero veridicitatea lor. Oamenii nu sunt sinceri în răspunsurile lor la întrebări intime, la fel și la întrebări ce țin de consumul de droguri, violența în familie, eschivarea de la achitarea corectă a impozitelor etc. Astfel, problema sporirii sincerității răspunsurilor este actuală [1].

### **1. Modelul de sondaj al lui Stanley Warner**

Stanley Warner [3] a propus un model de sondaj care sporește substanțial veridicitatea răspunsurilor. Una dintre variantele posibile constă în următoarele. Presupunem că se efectuează un sondaj pentru a determina amploarea în societate a unui fenomen  $C$  (negativ). Să notăm prin  $\pi$  probabilitatea ca un individ luat la întâmplare să fie afectat de fenomenul  $C$ . Pentru evaluarea acestei probabilități se poate lua un eșantion de  $N$  indivizi și dacă  $k$  dintre ei sunt afectați de fenomenul  $C$ , atunci putem spune  $\pi \approx \frac{k}{N}$ . Spre regret, numărul  $k$  este necunoscut. Prin întrebări directe este imposibil să aflăm valoarea adevărată a lui  $k$ . Se propune următorul procedeu.

Fiecărui dintre  $N$  respondenți  $i$  se înmânează o fișă cu două întrebări:

1. Este adevărat că sunteți afectat de fenomenul  $C$ ?
2. Este adevărat că nu sunteți afectat de fenomenul  $C$ ?

Fiecare respondent va răspunde doar cu „da” sau „nu” și doar la una dintre aceste întrebări, întrebare „aleasă” la întâmplare. Cum se poate organiza aceasta?

De exemplu, odată cu fișa fiecărui respondent  $i$  se înmânează și un zar pe care urmează să-l arunce în sus, o singură dată. Acei care obțin, de exemplu, cel mult două puncte, vor răspunde la prima întrebare, iar cei care obțin trei puncte și mai multe vor răspunde la întrebarea a doua. Totodată, nimeni nu trebuie să cunoască rezultatul aruncării și, în consecință, nici la care din întrebare se dă răspuns. Astfel, fiecare respondent răspunde doar la una din întrebări: la prima întrebare cu probabilitatea  $p = \frac{1}{3}$ , iar la a doua cu probabilitatea

$1 - p = \frac{2}{3}$ . Bineînțeles, putem alege și alte probabilități, în loc de zar utilizând alte mecanisme (tabele de numere aleatoare etc.). Cazurile  $p = 1$  și  $p = 0$  conduc la cazul de sondaj clasic.

În condițiile în care este pus respondentul acum, el ar trebui (cel puțin din punct de vedere psihologic) să se simtă în afara oricărui „pericol” venit din cauza unui răspuns sincer.

Într-adevăr, ambele răspunsuri posibile, și „da” și „nu”, pot avea o interpretare „pozitivă”: răspunsul „da” poate fi atribuit celei de a doua întrebări, iar „nu” primei întrebări.

Astfel, putem spera că în urma unui astfel de sondaj se vor obține rezultate mult mai veridice.

În [2] pentru probabilitatea  $\pi$ , de care suntem interesați, noi am stabilit formula

$$\pi = \frac{\frac{n}{N} - 1 + p}{2p - 1}; \left( p \neq \frac{1}{2} \right). \quad (1)$$

Aici  $N$  este numărul respondenților,  $n$  este numărul răspunsurilor de „da” date de ei, iar  $p$  reprezintă probabilitatea cu care un respondent răspunde la prima întrebare. Tot acolo, în [2], este examinată probabilitatea  $\pi$  ca funcție de  $n$ ,  $n$  fiind interpretat ca o variabilă aleatoare.

În prezenta lucrare este studiată în continuare probabilitatea  $\pi$ , sunt expuse câteva metode de estimare, este testată o ipoteză statistică privind valoarea lui  $\pi$ .

Numărul  $n$  de „da” poate fi privit ca o variabilă aleatoare. Se constată cu ușurință că  $n$  are repartiție binomială  $n \sim B(N, \psi)$ , unde  $\psi = 1 - p + (2p - 1)\pi$ , și, deci,  $Mn = N(1 - p + (2p - 1)\pi)$ ;  $M$  este simbolul valorii medii.

Astfel, pentru probabilitatea  $\pi$  obținem formula (exactă)

$$\pi = \frac{\frac{Mn}{N} - 1 + p}{2p - 1}, \quad (2)$$

formulă ce diferă de (1) prin aceea că  $n$  este înlocuit cu valoarea lui medie.

În continuare vom examina rolul lui  $\psi$  la stabilirea valorii probabilității  $\pi$ . Amintim că  $\psi$  reprezintă probabilitatea răspunsului „da” al unui respondent și este o funcție de  $\pi$ .

Mai întâi menționăm că: a) dacă  $p < \frac{1}{2}$ , atunci  $\psi = 1 - p + (2p - 1)\pi \in [p, 1 - p]$ ;

b) dacă  $p > \frac{1}{2}$ , atunci  $\psi = 1 - p + (2p - 1)\pi \in [1 - p, p]$ .

## 2. Unele metode de evaluare a probabilității $\pi$

Vom expune unele metode de evaluare a lui  $\pi$ , atunci când se cunoaște numărul răspunsurilor de „da” obținute.

**Metoda 1.** Presupunem că se efectuează un sondaj în care se obțin  $n^*$  de „da”. Dacă  $p < \frac{1}{2}$  și  $n^* \in [Np, N(1 - p)]$ , sau dacă  $p > \frac{1}{2}$  și  $n^* \in [N(1 - p), Np]$ , atunci (în ambele cazuri) drept estimatie a

funcției parametrice  $\psi(\pi) = 1 - p + (2p - 1)\pi$  putem lua statistica  $\hat{\psi} = \frac{n^*}{N}$ . Prin urmare, în acest caz

pentru probabilitatea  $\pi$  obținem estimatia  $\hat{\pi} = \frac{\hat{\psi} - 1 + p}{2p - 1}$ .

**Metoda 2.** Pentru parametrul  $\psi$ , ca probabilitate, putem construi un interval de încredere.

Se știe că frecvența relativă  $\frac{n^*}{N}$  este o estimatie eficientă a probabilității  $\psi$ . Alegând un prag de semnificație  $\alpha$ , construim intervalul de încredere

$$\left( \frac{n^*}{N} - u_{1-\frac{\alpha}{2}} \sqrt{\frac{\frac{n^*}{N} \left(1 - \frac{n^*}{N}\right)}{N}}, \frac{n^*}{N} + u_{1-\frac{\alpha}{2}} \sqrt{\frac{\frac{n^*}{N} \left(1 - \frac{n^*}{N}\right)}{N}} \right)$$

(aici  $u_{1-\frac{\alpha}{2}}$  este „ $1 - \frac{\alpha}{2}$  - cuantila” repartiției normale  $N(0,1)$ ).

Acest interval acoperă valoarea parametrului  $\psi$  cu probabilitatea  $1 - \alpha$ . Mai important pentru noi este să construim un interval de încredere pentru probabilitatea  $\pi$ . Alegând un prag de semnificație  $\alpha$  și ținând cont de legătura dintre  $\pi$  și  $\psi$ , putem scrie:

$$\frac{n^*}{N} - u_{1-\frac{\alpha}{2}} \sqrt{\frac{\frac{n^*}{N} \left(1 - \frac{n^*}{N}\right)}{N}} < 1 - p + (2p-1)\pi < \frac{n^*}{N} + u_{1-\frac{\alpha}{2}} \sqrt{\frac{\frac{n^*}{N} \left(1 - \frac{n^*}{N}\right)}{N}},$$

adică

$$\frac{n^*}{N} + p - 1 - u_{1-\frac{\alpha}{2}} \sqrt{\frac{\frac{n^*}{N} \left(1 - \frac{n^*}{N}\right)}{N}} < (2p-1)\pi < \frac{n^*}{N} + p - 1 + u_{1-\frac{\alpha}{2}} \sqrt{\frac{\frac{n^*}{N} \left(1 - \frac{n^*}{N}\right)}{N}}.$$

În funcție de valoarea lui  $p$ , conchidem:

a) dacă  $p < \frac{1}{2}$ , atunci

$$\frac{1}{2p-1} \left[ \frac{n^*}{N} + p - 1 + u_{1-\frac{\alpha}{2}} \sqrt{\frac{\frac{n^*}{N} \left(1 - \frac{n^*}{N}\right)}{N}} \right] < \pi < \frac{1}{2p-1} \left[ \frac{n^*}{N} + p - 1 - u_{1-\frac{\alpha}{2}} \sqrt{\frac{\frac{n^*}{N} \left(1 - \frac{n^*}{N}\right)}{N}} \right];$$

b) dacă  $p > \frac{1}{2}$ , atunci

$$\frac{1}{2p-1} \left[ \frac{n^*}{N} + p - 1 - u_{1-\frac{\alpha}{2}} \sqrt{\frac{\frac{n^*}{N} \left(1 - \frac{n^*}{N}\right)}{N}} \right] < \pi < \frac{1}{2p-1} \left[ \frac{n^*}{N} + p - 1 + u_{1-\frac{\alpha}{2}} \sqrt{\frac{\frac{n^*}{N} \left(1 - \frac{n^*}{N}\right)}{N}} \right].$$

**Metoda 3.** Pentru  $\pi$  putem construi o estimatie de verosimilitate maximă.

Amintim că  $\psi$  este funcție de probabilitatea  $\pi$ :  $\psi = 1 - p + (2p-1)\pi$ . Evident,  $P(n = n^*) = C_N^{n^*} \psi^{n^*} (1-\psi)^{N-n^*}$ .

Conform teoremei Moivre-Laplace,

$$P(n = n^*) \approx \frac{1}{\sqrt{N\psi(1-\psi)}} \cdot \varphi\left(\frac{n^* - N\psi}{\sqrt{N\psi(1-\psi)}}\right),$$

astfel probabilitatea  $P(n = n^*)$  este o funcție  $g(\pi)$ .

Ca estimatie a lui  $\pi$  putem lua acea valoare  $\pi^*$  a lui  $\pi$ , pentru care probabilitatea  $P(n = n^*)$  este maximă. Deoarece funcția  $g(\pi)$ , ca funcție de  $\pi$ , este destul de complicată și greu de cercetat la maximum, se va proceda în felul următor: în intervalul  $(0,1)$  cu un pas oarecare  $h$  (de exemplu, cu  $h = 0,01$ ) examinăm pe rând valorile lui  $\pi$ , calculând pentru fiecare din ele valoarea funcției  $g(\pi)$ . Valoarea  $\pi^*$  va fi acea valoare a lui  $\pi$ , pentru care  $g(\pi)$  ia valoarea maximă.

**Metoda 4.** Amintim că  $\psi$  reprezintă probabilitatea ca un respondent să răspundă cu „da”;  $\psi = 1 - p + (2p-1)\pi$  și deci  $\pi = \frac{\psi - (1-p)}{2p-1}$ . Observăm că legătura dintre  $\psi$  și  $\pi$  este bijectivă.

În continuare, pentru a determina valoarea (aproximativă) a lui  $\psi$  și, prin urmare, și a lui  $\pi$ , problema respectivă este formulată ca o problemă de verificare a unei ipoteze statistice privind parametrul repartiției binomiale (am menționat deja că  $\psi$  este parametrul unei repartiții binomiale).

Într-adevăr, problema privind determinarea probabilității  $\psi$  poate fi privită ca o problemă de verificare a ipotezelor, anume:

Să se verifice ipoteza nulă

$$H_0: \psi = \psi_0$$

cu alternativa:

$$H_1: \psi > \psi_0, \psi_0 \in (0,1).$$

Pentru realizarea practică a acestei idei vom formula următoarea serie de probleme:

Să se verifice ipoteza nulă  $H_0: \psi = kh$  cu alternativa  $H_1: \psi > kh$  ( $k = 1, 2, \dots, \left[\frac{1}{h}\right]$ ) sau, pe scurt:

$$H_0: \psi = kh,$$

$$H_1: \psi > kh.$$

Aici  $h \in (0,1)$  este un număr mic (de exemplu,  $h = 0,01$ ). La rezolvarea acestei serii de probleme ne vom conduce de următorul algoritm:

Se ia  $k = 1$  și, aplicând un test statistic corespunzător (care urmează să fie descris), se rezolvă problema:

$$H_0: \psi = h,$$

$$H_1: \psi > h.$$

Dacă ipoteza  $H_0$  este acceptată, atunci problema privind valoarea probabilității  $\psi$  este rezolvată:  $\psi = h$ . Dacă ipoteza  $H_0$  este respinsă și, prin urmare, este acceptată alternativa  $H_1$  ( $\psi > h$ ), atunci se formulează problema:

$$H_0: \psi = 2h,$$

$$H_1: \psi > 2h.$$

Se aplică testul menționat și se rezolvă această problemă. Dacă ipoteza nulă  $H_0$  este acceptată, atunci problema despre valoarea probabilității  $\psi$  este rezolvată:  $\psi = 2h$ . În caz contrar este acceptată alternativa  $H_1$ , se ia  $k = 3$  și se testează ipoteza  $H_0: \psi = 3h$  cu alternativa  $H_1: \psi > 3h$  etc.

Sunt posibile două variante de încheiere a acestui algoritm:

a) pentru un  $k_0$  oarecare ipoteza  $H_0$  este acceptată și valoarea probabilității  $\psi$  este stabilită:  $\psi = k_0h$ ;

b) toate ipotezele  $H_0: \psi = kh, k = 1, 2, \dots, \left[\frac{1}{h}\right]$  sunt respinse și, prin urmare, seria noastră de probleme se

încheie cu acceptarea ipotezei  $H_1: \psi > \left[\frac{1}{h}\right] \cdot h$ .

În acest caz, problema privind valoarea probabilității  $\pi$  încă nu este rezolvată, doar se poate constata că  $\psi > \left[\frac{1}{h}\right] \cdot h$ , adică  $|1 - \psi| < h$ . Deci, putem considera că  $\psi \approx 1$ , dacă  $h$  este suficient de mic. Dar putem micșora „pasul”, luând, de exemplu, un nou „pas”  $h^* = 0,01h$  și trecând la rezolvarea unei serii noi de probleme:

$$H_0: \psi = kh^*,$$

$$H_1: \psi > kh^*, k = 1, 2, \dots, \left[\frac{1}{h^*}\right].$$

Și de data aceasta vom aplica algoritmul expus mai sus. Din nou putem avea cazurile a) și b), menționate deja. În cazul b) vom lua  $\psi \approx 1$  sau vom continua testarea cu un pas mai mic, rezolvând o serie nouă de probleme.

În continuare vom descrie testul statistic care, în algoritmul expus mai sus, se aplică la verificarea ipotezelor de tipul:

$$H_0: \psi = \psi_0$$

cu alternativa

$$H_1: \psi > \psi_0, \psi_0 \in (0,1).$$

Fie că  $n^*$  dintre cei  $N$  respondenți au răspuns „da”. Conform teoremei integrale Moivre-Laplace, dacă ipoteza  $H_0$  este adevărată, atunci pentru valori mari ale lui  $N$  ( $N > 50$ ) repartiția variabilei aleatoare

$$Z = \frac{n^* - N\psi_0}{\sqrt{N\psi_0(1-\psi_0)}}, \text{ sau, echivalent, a lui } Z = \frac{\frac{n^*}{N} - \psi_0}{\sqrt{\frac{\psi_0(1-\psi_0)}{N}}}$$

este aproximativ normală cu parametrii 0 și 1.

Această variabilă aleatoare este luată în calitate de statistică a testului; regiunea critică a testului statistic pentru pragul de semnificație  $\alpha$  este mulțimea valorilor lui  $Z$ , determinate de inegalitatea  $z > u_{1-\alpha}$ , unde

$u_{1-\alpha}$  este „ $1-\alpha$  – cuantila” funcției de repartiție  $\Phi(x) = \frac{1}{\sqrt{2\pi}} \int_{-\infty}^x e^{-\frac{1}{2}u^2} du$ ”:  $\Phi(u_{1-\alpha}) = 1 - \alpha$ .

**Metoda 5.** La efectuarea a  $N$  experimente Bernoulli cu probabilitatea succesului egală cu  $\psi = 1 - p + (2p - 1)\pi$  în fiecare experiment, cel mai probabil număr de succese  $n_0$  este partea întregă a numărului fracționar  $(N - 1)(1 - p + (2p - 1)\pi)$ ; în caz contrar  $n_0$  are două valori:

$$(N - 1)(1 - p + (2p - 1)\pi); (N - 1)(1 - p + (2p - 1)\pi) + 1.$$

Aceasta ne poate permite să estimăm probabilitatea  $\pi$ . Într-adevăr, fie că am efectuat un sondaj cu  $N$  respondenți și  $n^*$  dintre aceștia au răspuns „da”. Cunoscând numărul  $n^*$  vom calcula succesiv valorile numărului  $n_0$  pentru  $\pi = h, 2h, \dots$ , eventual până la  $\left[\frac{1}{h}\right] \cdot h$ , oprindu-ne la valoarea  $\pi = kh$ , pentru care  $n_0 = n^*$ . În consecință, considerăm  $\pi \approx kh$ .

### 3. Aplicații și concluzii

Metodele descrise mai sus au fost testate pe un eșantion de 80 de studenți de la USM. Au fost abordate trei fenomene: 1) mituirea profesorilor; 2) dezamăgirea în alegerea specialității (facultății); 3) copierea la examen. În toate cazurile s-a luat  $p = \frac{1}{3}$ , adică la prima din cele două întrebări formulate în fișe a răspuns, în medie, fiecare al treilea dintre cei chestionați. Numărul răspunsurilor de „da” în legătură cu cele trei fenomene este, respectiv: 1)  $n_1 = 45$ ; 2)  $n_2 = 37$ ; 3)  $n_3 = 85$ .

Se verifică cu ușurință că aceste rezultate satisfac condițiile care permit să aplicăm metodele de estimare a probabilității  $\pi$ , expuse mai sus. Astfel, 4 din aceste metode au condus la următoarele valori aproximative (sau intervale de încredere pentru  $\pi$ ):

Metoda 1.  $\pi_1 \approx 0,3125$ ;  $\pi_2 \approx 0,6125$ ;  $\pi_3 \approx 0,6875$ .

Metoda 2.  $\pi_1 \in (0,0039; 0,636)$ ;  $\pi_2 \in (0,2887; 0,9285)$ ;  $\pi_3 \in (0,3667; 0,9884)$ .

Metoda 3.  $\pi_1 \approx 0,30$ ;  $\pi_2 \approx 0,61$ ;  $\pi_3 \approx 0,70$ .

Metoda 5.  $\pi_1 \in (0,297; 0,353)$ ;  $\pi_2 \in (0,593; 0,629)$ ;  $\pi_3 \in (0,67; 0,70)$ .

Trebuie să recunoaștem că volumul eșantionului este destul de modest și sondajul nu s-a efectuat în toate grupele concomitent. Rolul sondajului este unul mai mult ilustrativ. Noi nu pretindem că rezultatele lui reflectă exact situația reală în problemele respective. Totodată, trebuie să menționăm faptul că referitor la mită în învățământ rezultate similare cu ale noastre am întâlnit în unele materiale oficiale.

### Bibliografie:

1. NATHAN, G. Bibliographie de la méthode des réponses randomisées (1965-1987). En: *Statistique Canada: Techniques d'enquête*, 1988, no.14, p.351-365.
2. POȘTARU, A., PRODAN, N., BENDERSCHI, O. Asupra metodei răspunsurilor randomizate. În: *Materialele Conferinței internaționale „Modelarea Matematică, Optimizare și Tehnologii Informaționale”*. Ediția a IV-a. Chișinău: Evrica, 2014. Vol.I, p.184-190. ISBN 978-9975-62-365-0
3. WARNER, S. Randomizes response: a survey technique for eliminating evasive answer bias. In: *J. Am. Statist. Assoc.*, 1965, no.60, p.63-69.

Prezentat la 01.10.2015