

OCENA REPREZENTATYWNOŚCI POBIERANIA PRÓBEK

UKD 543.05:519.26/27:552.1:553.1

Zbiór wartości ξ zmiennej losowej niezależnych od siebie określeń charakteryzujących badaną cechę danego ośrodka nazywamy próbą. Proces określania tej próby nazywamy pobieraniem próbki. Od pojedynczych elementów próby wymagamy, aby były od siebie niezależne i aby w sposób reprezentatywny przedstawiały badany zbiór. Można to osiągnąć w warunkach, gdy badana cecha posiada ciągły rozkład gęstości prawdopodobieństwa. Rozkład ten jest monotonicznie rosnący w obrębie odcinka (a, b) . Gęstość prawdopodobieństwa $f(x)$ badanej cechy określa dystrybuanta:

$$F(x) = \int_a^x f(x) dx$$

Gdy badana cecha zależy od funkcji położenia $F(x)$, wówczas miejsca pobrania próbek powinny być roz-

łożone równomiernie w całej przestrzeni badanego obszaru. Gdy warunek równomiernego rozmieszczenia próbek nie jest spełniony, np. gdybyśmy próbki pobierali tylko z części danego rejonu możemy nie otrzymać próby reprezentatywnej. Funkcję położenia można jednak rozłożyć na sumy rozkładów $F_i(x)$ odcinkowych o s_i wagach proporcjonalnych do długości a_i odcinków [1]:

$$F(x) = \sum_{i=1}^n s_i F_i(x)$$

gdzie:

$$F_i(x) = \int_{a_i}^{a_i + \Delta a_i} f_i(x) dx \quad [1]$$

$$s_i = \frac{a_i}{A}$$

i gdzie na ogół $f_i(x) \neq f(x)$

przy czym:

- $F(x)$ — jest funkcją rozkładu badanej wartości x zbioru A ,
- a_i — jest „wielkością” elementarnego zbioru położonego wewnątrz badanego zbioru,
- $F_i(x)$ — jest funkcją rozkładu x w danym cząstkowym zbiorze.

Ogólnie $F_i(x)$ różni się od $F(x)$, więc nawet przy próbkach pobranych z rejonu zawężonego założony wzór empiryczny funkcji rozkładu nie będzie zbliżony do $F(x)$. Próbkę należy pobierać tak, aby „odległości”, zmierzone między sąsiednimi miejscami poboru próbek w obrębie danego rejonu były możliwie jednakowe, tzn. aby ich zagęszczenie w całym rejonie było jednakowe.

W ten sposób można najskuteczniej zapewnić reprezentatywność pobranych próbek, gdy ilość tych próbek jest dostatecznie duża, ponieważ każdą elementarną próbkę ma tę samą wagę, mimo że pochodzi z różnych miejsc badanego rejonu. Gdy więc wszystkie próbki jednakowe mają wspólne wagi i pochodzą z różnych miejsc badanego rejonu, a ze względu na ciągłość tych zmian próbki będą charakterystyczne dla strefy o coraz to innych właściwościach, wówczas taką próbkę możemy uznać za reprezentatywną.

Takie określenie obszarów próbek eliminuje przypadek, gdy odległości między próbkami są dowolnie małe, a zatem gdy zmienna ich wartość ma charakter ciągły. Wówczas oba punkty pobrania próbek leżą w bezpośrednim sąsiedztwie. Próbkę dublującą się charakteryzują tym samym obszarem. Obszar ten byłby w tym przypadku bez uzasadnienia w próbkę liczony podwójnie, co nie byłoby prawidłowe. Dlatego też jako wartości takich dublujących się próbek należy traktować łącznie jako wartość średnią z obu próbek dublujących się.

Rozważmy dla jakiego obszaru odnosi się to ważne stwierdzenie. Wyjdźmy z założenia, iż każdy obszar próbki ma tę samą wagę, zatem w naszych rozważaniach nie ma takiej próbki, której należałoby przypisać większą wagę. Gdy liczba próbek = N , a każda próbka jest dla swego otoczenia reprezentatywna, tzn. czy spełnia warunek reprezentatywności dla obszaru o promieniu równym r . Obszar ten przedstawia n -tą część całości badanego obszaru. Wychodząc z tego założenia, gdyby miejsca pobrania dwóch sąsiednich próbek oddalone były od siebie o odległość mniejszą jak r , to nie mogą być one uznane za samostatne jednostki próbki ogólnej. Na podstawie powyższego, zajmijmy się teraz problemem związanym z określeniem wagi kilku różnych próbek. Zagadnienie to można przeprowadzić jedną z metod:

a) Przez kilkakrotne całkowanie stochastycznej zmiennej metodą Monte-Carlo, rozwiązując całkę:

$$I = \int_{G_n} f_n(P_n) dP_n \quad [2]$$

przy czym:

- G_n — jest danym obszarem w przestrzeni n -wymiarowej, określonym zbiorem,
- $P_n = P_n(x_1, x_2, \dots, x_n)$ punktów przynależnych do obszaru G_n ,
- $dP_n = dx_1 \dots dx_n$.

Jeżeli w obszarze G_n posiadamy ilość N punktów P_i rozłożonych regularnie, stosując stochastyczne całkowanie metodą Monte-Carlo możemy z [2] napisać:

$$I = |G_n| E[f_n(P_n)] \cong |G_n| \frac{1}{N} \sum_{i=1}^N f_n(P_{ni}) \quad [3]$$

gdzie:

- E — operator wartości oczekiwanej (nadzieja matematyczna),
- $|G_n|$ — oznacza obszar G_n w przestrzeni n wymiarowej.

Przyjęto, że wielowymiarowy rozkład $f_n(P_n) = \eta_n$ w obrębie ograniczonego obszaru (a, b) ma przebieg monotoniczny. Jeżeli całkowanie przeprowadzimy dla zmiennej x_n , otrzymamy:

$$I = \int_{G_{n-1}} f_{n-1}(P_{n-1}) dP_{n-1} \quad [4]$$

gdzie:

$$f_{n-1}(P_{n-1}) = \int_{x_n = r_1(P_{n-1})}^{\mu(P_{n-1})} f_n(P_n) dx_n \quad [5]$$

Znaczenie wyrażeń G_{n-1} , P_{n-1} , dP_{n-1} oraz η_{n-1} jest analogiczne jak poprzednio. Ogólnie, jeżeli się całkuje względem zmiennej „ k ”, otrzyma się w wyniku całkowania zmienną $\eta_{(n-k)}$. W przypadku pierwszym $\eta_{(n-k)}$ posiada wartość oczekiwaną w przybliżeniu:

$$\bar{\eta}_{(n-k)} = \frac{\sum_{i=1}^N (n-k) i}{\sum_{i=1}^N |G_{(n-k) i}|} \quad [6]$$

przy czym:

$|G_{(n-k)}|$ — jest powierzchnią takich obszarów o ilości wymiarów niższym od tego, dla którego określona jest całka [2].

Próbkę złożoną z $\eta_{(n-k)}$ elementów uważać należy za reprezentatywną wtedy, jeżeli dwa sąsiednie miejsca pobrania leżą w odległości mniejszej od r . Wielkość r można oznaczyć wychodząc z założenia, że suma poszczególnych płaszczyzn w rejonie o promieniu r da wartość zgodną z wartością $|G_n|$.

Jeżeli próbka pobrana została w punkcie ($k=0$), wówczas wartość przedtęta może być obliczona na podstawie [3] ze stopniem prawdopodobieństwa przy odchyleniu granicznym ϵ . Przy pobraniu próbek o wymiarze niepunktowym, jeżeli z [6] otrzymamy wartość $\bar{\eta}_{(n-k)}$, wtedy $s-t$ określić możemy na podstawie następującego wzoru na N'

$$N' \cong \frac{|G_n|}{|G_{er}|} \quad [7]$$

Przy czym:

$|G_{er}|$ — jest wielkością obszaru o promieniu r miejsca poboru próbki o wymiarze punktu.

Wartość r należy określić na podstawie metody podanej poprzednio. Jeżeli zaś można podać wartość rozkładu $f_n(P_n)$ w obszarze G_n , wtedy można również określić standard σ . Pobraną próbkę z $\eta_{(n-k)}$ jednostek możemy przyrównać do próbki punktowej. Jeżeli zaś w wyrażeniu $\eta_{(n-k)}$ także i i k jest zmienną, wtedy średnią η_n należy obliczyć według

$$\bar{\eta}_{ni} = \frac{\eta_{(n-k) i}}{|G_{(n-k) i}|} \quad [8]$$

Próbki także będą reprezentatywne, jeśli dane części $|G_{(n-k)}|$ hiperpłaszczyzn nie są oddalone od

siebie mniej, jak o r . Tak jak poprzednio r określić można z warunku, gdy części hiperpłaszczyzn $|G_{(n-k)r}|$ wraz z otoczeniem o promieniu r dadzą w sumie $|G_n|$.

Gdy wielkość pojedynczych obszarów leżących w płaszczyznach wyższego rzędu z przestrzenią o promieniu r będzie $|G_{(n-k)r}|$, wtedy:

$$\bar{\eta}_n = \frac{\sum_{i=1}^N |G_{(n-k)r}| \bar{\eta}_{ni}}{|G_n|} \quad [9]$$

czyli pojedyncze próbki $\bar{\eta}_{ni}$ należy wstawić do rachunku z takimi wagami, którym odpowiada wielkość $|G_n|$ reprezentowana w przestrzeni próby częściowej.

Wartość r wyznaczamy podobnie, jak to już uczyniono uprzednio według formuły [7], wówczas określmy N' , a więc można również obliczyć σ_n' .

b) Pobieranie próbki wieloelementarnej.

Przyjawszy, że x oznacza pewną cechę, której dystrybuanta oznaczona jest przez $F(x)$. Jeżeli $F(x)$ w obrębie obszaru jest ściśle monotoniczne, wówczas szacowanie oczekiwanej wartości może nastąpić przy użyciu metod poprzednio opisywanych na podstawie x próbek w ilości N .

W praktyce pobieranie próbki x_i w punkcie jest niemożliwe. Z materiału zbioru o objętości V przy pobieraniu próbki zwykle wyodrębnia się ilość ΔV_i , a w laboratorium oznacza się dla tej ilości przynależną wartość średnią x_i . Jeżeli ze zbioru o rozkładzie $F(x)$ pobiera się próbki V_i o wielkości stałej, wówczas średni rozkład wartości x_i w objętości V_i będzie $F_i(x)$. Jeżeli do niego przynależna wartość σ jest σ_i , wtedy:

$$\sigma_i \leq \sigma \quad [10]$$

gdzie:

σ — jest rozrzutem $F(x)$.

Przypadek ten ważny jest dla każdej próbki o wielkości $\Delta V_i < V$. Do oznaczenia odchylenia ważne są zależności podane [2]. W dalszym ciągu należy zbadać, jak można oszacować granice błędów s dla różnych poziomów prawdopodobieństwa, gdy jest różna wielkość próbek. Przy analizie tego zagadnienia należy rozróżnić dwa przypadki:

1) gdy:

$$\sum_{i=1}^N \Delta V_i = \Delta V \quad [11]$$

czyli próbka całkowita jest wobec całości badanego zbioru nikła (gdy ΔV wynosi tylko kilka % V),

2) gdy ΔV w stosunku do V nie da się pominać.

Dla uproszczenia przyjmujemy próbki o objętości V_i jako kule o promieniu r_i , co daje:

$$\frac{4r_i^3 \pi}{3} = \Delta V_i \quad (i = 1, 2, \dots, N) \quad [12]$$

więc jeżeli wartość r jest rozumiana jak podano poprzednio, to do jej określania posłużyć może zależność:

$$\sum_{i=1}^N \frac{4(r_i + r)^3 \pi}{3} = V \quad [13]$$

Z tego wynika, że w pierwszym przypadku:

$$r \approx \sqrt[3]{\frac{3V}{4N\pi}}; \quad [14]$$

gdyż wartość $\frac{\Delta V}{V}$ jest nikła. Dlatego też wobec r pominać możemy także V_i . Tym samym próbki o różnej objętości V_i są wyrazem objętości prawie że równych, więc przy określaniu na ich podstawie wartości x_i należy przyjąć jednakowe uwagi. W drugim przypadku wartość r należy określić według [13]. Znajac już tę wartość oblicza się \bar{X} , jak następuje:

$$\bar{X} = \frac{4\pi}{3} \cdot \frac{\sum_{i=1}^N (r_i + r)^3 x_i}{V} \quad [15]$$

Znaczy to, iż każda próbka częściowa do próbek ogólnej wstawiona będzie z taką wagą, jaka odpowiada jej udziałowi w całej objętości V .

Naturalnie, obliczenia te słuszne są tylko wówczas, gdy próbki są od siebie niezależne oraz reprezentatywne. Wartości średnie można oznaczyć na podstawie znanych z literatury (2, 3) metod odnośnie do odchylenia odpowiadającego danemu stopniowi prawdopodobieństwa. Oczywiście, że wartość σ należy wypośredkować z empirycznego wzoru na rozkład. (Także w przypadku różnych wag próbek). Wzór [4] udowadnia, że przy ważonych wartościach przeciętnych, lepsze wyniki daje zwykła średnia jednakowych wag, gdyż wykazuje ona najmniejszy rozrzut. To samo można udowodnić przy założeniu, iż próbka reprezentuje zawartość danej przestrzeni o promieniu r , gdy pobrano próbki łącznego zbioru. Przy warunku [16] zależność [15] określi nam wartość r ,

$$\sum_{i=1}^N \Delta V_i = \Delta V = \text{constans} \quad [16]$$

która wtedy osiągnie swoją wartość najniższą przy określonych wartościach N i V , jeżeli próbki posiadają jednakową wielkość.

$$\left(\Delta V_i = \Delta V_j = \frac{\Delta V}{N} \quad i \neq j \right)$$

LITERATURA

1. Edwin F. Backenbach — Modern matematika mérnököknek. Műszaki Könyvkiadó, Budapest, 1960.
2. Janositz J. — Adott intervallumon változó folytonos eloszlású valószínűségi változó várható értékének vizsgálata. Nehezipari Műszaki Egyetem Közleményi (w druku).
3. Janositz J. — A várható érték analízise. Bányászati és Kohászati Lapok. Bányászati (w druku).
4. Jánossy L. — A valószínűségelmélet alapjai és néhány alkalmazása különös tekintettel mérési eredmények kiértékelésére. Tankönyvkiadó, 1965.
5. Prékopa A. — Valószínűségelmélet. Műszaki Könyvkiadó, Budapest, 1962.
6. Rényi A. — Valószínűség-számítás. Tankönyvkiadó, Budapest, 1954, 1966.
7. Srejgyer I. U. — Monte Carlo módszerek. Műszaki Könyvkiadó, Budapest, 1965.
8. Smirnov N. W., Dunin-Borkowski I. W. — Mathematische Statistik in der Technik. VEB Deutscher Verlag der Wissenschaften, Berlin, 1963.

SUMMARY

The paper deals with the problem of taking representative samples in the case when any feature of a given series analysed changes continuously. To obtain a representative picture of the whole series we must take samples within the series analysed, according to a uniform density net, i.e. we must take care to obtain samples that represent similar volume. Under such conditions, when the amount of sampling points are determined, any sample must representatively determine the features of a given series within an area (having radius r). On this basis a possibility exists to determine average weight of a sample. The weight should always be proportional to the area with the radius r determined previously.

Average value of an error ε , calculated on the basis of samples having known probability level, may also be determined. The amount of samples characterized by various weight is a quotient of the entire area and of the elementary areas with the radius r .

РЕЗЮМЕ

В статье обсуждается методика отбора представительных проб, в том случае когда один из признаков исследуемого объекта изменяется последовательно. Для того, чтобы пробы были представительные для всего объекта, они должны удовлетворять условию, чтобы места опробования располагались по всему объему с одинаковой густотой, следовательно, чтобы каждая проба представляла одинаковый объем. В условиях определенного количества мест опробования проба определяет свойства объекта на площади с радиусом r . Таким образом, при пробах величиной k ($< n$) в n -размерном пространстве можно определить средний вес пробы, когда k ($< n$) является постоянной величиной, или же когда взятые пробы различны (k переменное). Вес всегда должен быть прямопропорционален ранее определенной площади с радиусом r .

Можно также определить величину средней ошибки E на основании проб с известным уровнем вероятности. Количество проб разного веса является кратностью всей площади и элементарных участков с радиусом r .