

Член-кореспондент НАН України С. І. Ляшко, В. В. Алексеєнко,
Д. А. Ключин

Непараметричний критерій однорідності двох вибірок на основі статистики включення

Розглянуто нову непараметричну статистику, основу на порядкових статистиках, що відображає близькість двох вибірок. Наведено алгоритм обчислення статистики. Побудовано критерій однорідності двох вибірок із довільним визначенням рівнем значущості.

Нехай $x = (x_1, x_2, \dots, x_n)$ і $y = (y_1, y_2, \dots, y_m)$ — вибірки, отримані простим випадковим вибором із генеральних сукупностей G та H відповідно, що мають абсолютно неперервні функції розподілу $F(u)$ і $G(u)$. Одна з основних задач статистичного розпізнавання образів — перевірка гіпотези H_0 про те, що вибірки x і y є однорідними, тобто $F(u) = G(u)$.

У роботі [1] розглядається статистика, що є сумою абсолютних відхилень між наближеними функціями щільності еталонної вибірки та еталонної вибірки, доповненої елементами тестової. Статистика має досить складну будову і не може застосовуватись для класифікації тестової вибірки серед класів, представлених еталонними вибірками різного обсягу.

Інша статистика, наведена у [2], будується як квадратичне відхилення між частотою і ймовірністю потрапляння елементів тестової вибірки в інтервали, визначені еталонною вибіркою. Для цієї статистики наводиться функція розподілу, але не наведений алгоритм обчислення, що ускладнює використання результатів на практиці.

На базі ідей цих двох робіт будується нова непараметрична статистика, що дозволяє розв'язувати задачу перевірки гіпотези H_0 .

Статистика включення. Нехай $x_{(1)} < x_{(2)} < \dots < x_{(n)}$ — варіаційний ряд вибірки x_1, x_2, \dots, x_n і $y_{(1)} < y_{(2)} < \dots < y_{(m)}$ — варіаційний ряд вибірки y_1, y_2, \dots, y_m .

Згідно з гіпотезою Хілла,

$$P(x_k \in (x_{(i)}, x_{(j)})) = \frac{j-i}{n+1}. \quad (1)$$

Введемо позначення:

$$I_0 = (-\infty; x_{(1)}), \quad I_n = (x_{(n)}, \infty), \quad I_j = (x_{(j)}, x_{(j+1)}), \quad j = \overline{1, n-1}. \quad (2)$$

Тоді, за формулою (1), має місце $p = P(x_k \in I_i) = 1/(n+1)$.

Позначимо l_i кількість y_j , таких, що $y_j \in I_i$, і $f_i = l_i/m$ — частоту потрапляння елементів вибірки y_j в інтервали I_i .

Введемо статистику включення, яка є абсолютним відхиленням між частотою і ймовірністю потрапляння елементів вибірки y_j в інтервали I_i :

$$\eta = \sum_{i=0}^n |p - f_i|. \quad (3)$$

Введемо такі позначення:

$$\bar{b} = \begin{cases} \left[\frac{m}{n+1} \right] + 1, & \frac{m}{n+1} > \left[\frac{m}{n+1} \right] \\ \frac{m}{n+1}, & \frac{m}{n+1} = \left[\frac{m}{n+1} \right] \end{cases};$$

$$\bar{h} - \text{кількість } i \text{ таких, що } l_i \geq \bar{b}, \quad \bar{w} = \sum_{l_i \geq \bar{b}} l_i - \bar{b}, \quad \delta_+ = \bar{b} - \frac{m}{n+1}; \quad \underline{b} = \bar{b} - 1;$$

$$\underline{h} - \text{кількість } i \text{ таких, що } l_i \leq \underline{b}, \quad \underline{w} = \sum_{l_i \leq \underline{b}} \underline{b} - l_i, \quad \delta_- = \frac{m}{n+1} - \underline{b}.$$

Лема 1. \underline{h} , \underline{w} , δ_- можна визначити через \bar{h} , \bar{w} , δ_+ .

Доведення. $\underline{h} = n - \bar{h}$, де n — довжина вектора \bar{l} ,

$$\begin{aligned} \underline{w} &= \sum_{l_i \leq \underline{b}} \underline{b} - l_i = \sum_{l_i \leq \underline{b}} \underline{b} - \sum_{l_i \leq \underline{b}} l_i = \underline{hb} - \left(m - \sum_{l_i \geq \bar{b}} l_i \right) = \\ &= (n - \bar{h})(\bar{b} - 1) - \left(m + \sum_{l_i \geq \bar{b}} (\bar{b} - l_i) - \sum_{l_i \geq \bar{b}} \bar{b} \right) = \\ &= n\bar{b} - n - \bar{hb} + \bar{h} - (m + \bar{w} - \bar{hb}) = n\bar{b} + \bar{h} - n - m - \bar{w}. \end{aligned}$$

Отже,

$$\delta_- = \frac{m}{n+1} - \underline{b} = \frac{m}{n+1} - \bar{b} + 1 = 1 - \delta_+.$$

Лема 2. Пара (\bar{h}, \bar{w}) однозначно визначає статистику η для визначених m і n .

Доведення.

$$\begin{aligned} \eta &= \sum_{i=0}^n |p - f_i| = \sum_{l_i \geq \bar{b}} \frac{\left| \frac{m}{n+1} - l_i \right|}{m} + \sum_{l_i \leq \underline{b}} \frac{\left| \frac{m}{n+1} - l_i \right|}{m} = \frac{\sum_{l_i \geq \bar{b}} \left(l_i - \frac{m}{n+1} \right) + \sum_{l_i \leq \underline{b}} \left(\frac{m}{n+1} - l_i \right)}{m} = \\ &= \frac{\sum_{l_i \geq \bar{b}} (l_i - \bar{b} + \delta_+) + \sum_{l_i \leq \underline{b}} (\underline{b} - l_i + \delta_-)}{m} = \frac{\bar{w} + \bar{h}\delta_+ + \underline{w} + \underline{h}\delta_-}{m}. \end{aligned}$$

За лемою 1 маємо

$$\eta = \frac{\bar{w} + \bar{h}\delta_+ + n\bar{b} + \bar{h} - n - m - \bar{w} + (n - \bar{h})\delta_-}{m} = \eta(\bar{h}, \bar{w}).$$

Твердження лєми справедливе, оскільки \bar{b} , δ_+ і δ_- обчислюються через m і n .

Розподіл статистики включення. Припустимо, що справедлива гіпотеза H_0 , тобто $F(u) = G(u)$. Введемо в розгляд об'єднану вибірку $z_1, z_2, \dots, z_{m+n} = x_1, x_2, \dots, x_n, y_1, y_2, \dots, y_m$. Нехай $z_{(1)}, z_{(2)}, \dots, z_{(m+n)}$ — відповідний варіаційний ряд.

Оскільки елементи вибірки y_j належать генеральній сукупності G і відповідному варіаційному ряду, то ми можемо вважати, що всі можливі позиції, що займають елементи y_j у варіаційному ряді, є рівноімовірними. Тобто будь-які m елементів вибірки z_1, z_2, \dots, z_{m+n} можуть належати тестовій вибірці. Набір позицій елементів тестової вибірки у варіаційному ряді однозначно визначає вектор \bar{l} частот потрапляння елементів вибірки y_j в інтервали I_i . Всього різних розподілів \bar{l} існує C_{m+n}^m — кількість способів обрати у варіаційному ряді m позицій, що відповідають елементам вибірки y_j .

Позначимо як $S(n, k, l)$ кількість способів розмістити k кульок по n урнах так, щоб в кожній було не більше, ніж l кульок. Ця величина може бути обчислена за допомогою таких рекурентних формул:

$$S(0, k, l) = 0, \quad S(n, 0, l) = 1, \quad S(n, 1, l) = n, \quad S(1, k, l) = 1, \quad k \leq l,$$

$$S(n, k, l) = S(n-1, k, l) + S(n-1, k-1, l) + \dots + S(n-1, k-l, l).$$

Теорема 1. Для визначених характеристик (\bar{h}, \bar{w}) кількість різних векторів \bar{l} , для яких характеристики будуть набувати цих значень, можна обчислити за формулою:

$$N(\bar{h}, \bar{w}) = C_{n+1}^{\bar{h}} C_{\bar{h}+\bar{w}-1}^{\bar{w}} S(n+1-\bar{h}, m-\bar{h}\bar{w}-\bar{w}, \bar{b}-1). \quad (4)$$

Доведення. Кількість способів обрати \bar{h} позицій серед $n+1$ дорівнює $C_{n+1}^{\bar{h}}$. Згідно з позначеннями, \bar{h} — кількість позицій вектора \bar{l} , вибіркові значення в яких перевищують або дорівнюють \bar{b} , \bar{w} — величина, на яку за цими позиціями сума елементів в інтервалі більша за \bar{b} . Кількість векторів, в яких на \bar{h} позиціях знаходяться числа, що перевищують або дорівнюють \bar{b} так, що $\sum_i (l_i - \bar{b}) = \bar{w}$ — це те саме, що кількість способів спочатку розташувати \bar{b} кульок у \bar{h} урнах, потім ще додати \bar{w} кульок до цих урн. Кількість способів — $C_{\bar{h}+\bar{w}-1}^{\bar{w}}$. На інших $n+1-\bar{h}$ позиціях залишається всього $m-\bar{h}\bar{w}-\bar{w}$ елементів, причому кількість елементів в кожному інтервалі не більше ніж $\bar{b}-1$. Тобто кількість варіантів дорівнює $S(n+1-\bar{h}, m-\bar{h}\bar{w}-\bar{w}, \bar{b}-1)$. Отже, загальна кількість векторів дорівнює $C_{n+1}^{\bar{h}} C_{\bar{h}+\bar{w}-1}^{\bar{w}} S(n+1-\bar{h}, m-\bar{h}\bar{w}-\bar{w}, \bar{b}-1)$.

Таким чином, показано, що для визначених m і n імовірність того, що розподілу інтервалів \bar{l} буде відповідати пара (\bar{h}, \bar{w}) , дорівнює

$$p(\bar{h}, \bar{w}) = \frac{C_{n+1}^{\bar{h}} C_{\bar{h}+\bar{w}-1}^{\bar{w}} S(n+1-\bar{h}, m-\bar{h}\bar{w}-\bar{w}, \bar{b}-1)}{C_{m+n}^m}. \quad (5)$$

Критерії однорідності. На основі викладеного можна побудувати такий критерій для перевірки гіпотези H_0 :

1) для даних m та n (розмірностей тестової та еталонної вибірки) обчислити ймовірності для всіх можливих пар (h, w) . Зазначмо, що ці ймовірності можна обчислити заздалегідь і занести в таблицю, адже вони залежать тільки від m та n ;

2) визначити пару характеристик (h_0, w_0) для заданих тестової і еталонної вибірки і відповідну ймовірність $p(h_0, w_0)$;

3) обчислити α — суму ймовірностей менш імовірних пар характеристик, ніж (h_0, w_0) . Фактично α можемо вважати досяжним рівнем значущості, тобто найменшим рівнем значущості, при якому гіпотеза H_0 відхиляється для таких значень характеристик;

Таблиця 1. Розподіли еталонних вибірок

n	μ	σ
100	0	1
150	0	2
200	1	1
250	-1	4

Таблиця 2. Результати експерименту

Класифікація	Експеримент 1				Експеримент 2				Експеримент 3			
	Розподіл	1	2	3	4	1	2	3	4	1	2	3
$K_{\text{успіх}}$	98	84	69	84	86	93	94	92	98	69	98	94
$K_{\text{помилка}}$	0	1	1	0	0	0	0	1	0	4	0	0
$K_{\text{невизн}}$	2	15	30	16	14	7	6	7	2	27	2	6
$AK_{\text{успіх}}$	99	96	93	99	98	100	100	99	99	94	100	99
$AK_{\text{помилка}}$	1	4	7	1	2	0	0	1	1	6	0	1

4) нехай бажаний рівень значущості α_0 . Тоді критерієм прийняття гіпотези H_0 буде $\alpha > \alpha_0$.

Розглянемо задачу класифікації. Нехай задано k класів, кожен з яких представлений відповідною еталонною вибіркою $x^{(i)}$, $i = 1, \dots, k$, і тестова вибірка y . Треба визначити, якому класу належить тестова вибірка. Для кожної еталонної вибірки і тестової вибірки за критерієм і заданим рівнем значущості перевіримо справедливість гіпотези H_0 . Якщо гіпотеза H_0 приймається згідно з критерієм рівно для однієї еталонної вибірки, то вважаємо, що тестову вибірку вдалося класифікувати і вона належить класу цієї еталонної вибірки, інакше класифікацію виконати не вдалося.

Побудована згідно з алгоритмом величина $1 - \alpha$ може бути використана як міра близькості двох вибірок.

Якщо будь-що необхідно обрати якийсь із класів, то можна побудувати альтернативний критерій класифікації — будемо вважати, що тестова вибірка належить тому класу, для еталонної вибірки якого величина $1 - \alpha$ є мінімальною. В такому випадку невизначеність не виникає.

Практична перевірка критеріїв однорідності. Методика перевірки критеріїв класифікації має такий вигляд. Генеруємо чотири еталонні вибірки, які будуть представляти класи класифікації (табл. 1).

Для кожного типу нормального розподілу генеруємо 100 тестових вибірок по 300 елементів в кожній. Кожну тестову вибірку класифікуємо за критерієм (К) і альтернативним критерієм (АК). Експеримент повторимо тричі. Результати наведені у табл. 2.

1. *Алексеев В. В.* Статистика включення // Журн. обчисл. і прикл. математики. – 2012. – № 1(107). – 105–111 с.
2. *Vairamov I. G., Ozkaya N.* On the nonparametric test for two sample problems based on spacing // J. of App. Stat. Science. – 2000. – 10, No 1. – P. 57–68.
3. *Hill B.* Posteriori distribution of percentiles: Bayes' theorem for sampling from a population // J. Amer. Statist. Assoc. – 1968. – 63, No 322. – P. 677–691.

Член-корреспондент НАН Украины С. И. Ляшко, В. В. Алексеенко,
Д. А. Ключин

Непараметрический критерий однородности двух выборок на основе статистики включения

Рассмотрена новая непараметрическая статистика, основанная на порядковых статистиках и отражающая близость двух выборок. Приведен алгоритм вычисления статистики. Построен критерий однородности двух выборок с произвольно заданным уровнем значимости.

Corresponding Member of the NAS of Ukraine S. I. Lyashko, V. V. Alexeyenko,
D. A. Klyushin

Nonparametrical test for homogeneity of two samples based on the inclusion statistics

A new nonparametrical statistics based on order statistics reflecting the proximity of two samples is investigated. An algorithm for calculation of this statistics is proposed. The test for homogeneity of two samples with an arbitrary significance level is developed.