

ВЫЧИСЛИТЕЛЬНЫЕ СИСТЕМЫ

Сборник трудов

Института математики СО АН СССР

1967 г.

Выпуск 28

ОБ ЭФФЕКТЕ ЗАВИСИМОСТИ МЕЖДУ ОЦЕНКАМИ
КАЧЕСТВА ПРИЗНАКОВ ОПОЗНАЮЩИХ СИСТЕМ

А.Ю. Мотузя
(Вильнюс)

В работе исследуется влияние зависимости между оценками качества систем признаков и вероятностью правильных решений при выборе наилучшей из указанных систем.

При синтезе опознавающих систем приходится решать так называемую задачу выбора качественных признаков. Формализуем эту задачу следующим образом. Пусть имеется множество $S = \{s_1, \dots, s_u\}$ (u - объем множества S) систем признаков (СП). Из этого множества требуется выбрать лучшую по заранее заданному критерию (например, по наибольшей вероятности верного опознавания речевых или зрительных образов) СП. Показатели качества (ПК) Q_1, \dots, Q_u этих СП точно неизвестны. О соотношениях этих ПК исследователь принимает решение $\pi \in R$; R - заданное множество решений) по состоятельным оценкам t_1, \dots, t_u монотонных (скажем возрастающих) функций величин Q_1, \dots, Q_u . Для подсчета каждой оценки t_i ($i=1, \dots, u$) используется выборка конечного объема N^i ($i=1, \dots, u$) реализаций $\xi_{v_i}^i$ ($v_i=1, \dots, N^i$) опознаваемых объектов (слов, геометрических фигур). Ввиду конечности объемов N^i и случайного характера объектов $\xi_{v_i}^i$ оценки t_1, \dots, t_u случайные величины, разброс которых возрастает с убыванием N^i . Это приводит к ошибкам при принятии решения π . В большинстве практических случаев требуется, чтобы вероятность правильного решения P была не ниже заданной величины P^* . Величина P возрастает с увеличением объемов выборок N^i ($i=1, \dots, u$). Но увеличение этих объемов связано с затратами. В тех случаях,

когда оценки t_1, \dots, t_u независимы, проблема определения наименьшего объема выборки $N'_{min} = \dots = N''_{min} = N'''_{min}$, при котором вероятность ошибочного решения $1 - P$ не превышает заданной величины $1 - P^*$, исследована в ряде работ [1-5]. Причем существуют практические методы решения этой задачи. Статистическую независимость оценок t_i можно достичь рандомизацией выбора реализаций ξ . При этом для вычисления каждой из оценок t_i ($i=1, \dots, u$) должны быть использованы различные выборки. Такой прием приводит к сильному завышению общего числа $N' + \dots + N'''$ реализаций опознаваемых объектов ξ . Во многих практических случаях все оценки t_1, \dots, t_u подсчитываются по одной выборке объема N опознаваемых объектов

$$\xi'_v = \dots = \xi''_v = \xi_v \quad (v=1, \dots, N) \quad . В этом случае величины t_1, \dots, t_u статистически зависимы.$$

В настоящей работе исследуем, как влияет зависимость между оценками t_1, \dots, t_u на вероятность правильного решения P . Исследования проведем в случае гауссовского распределения величин t_1, \dots, t_u . Множество решений R определим следующим образом $R = \{z_1, \dots, z_c, \dots, z_u\}$, где z — такое решение, что между ПК Q_1, \dots, Q_u существует следующее соотношение:

$$Q_i \geq \max(Q_1, \dots, Q_u) - \delta. \quad (I)$$

Здесь δ — так называемый коэффициент неточности [5].

В этом случае решение z_i будет следующей функцией оценок t_1, \dots, t_u :

$$z_i = \begin{cases} 1, & \text{если } t_i = \max(t_1, \dots, t_u); \\ 0 & \text{в другом случае.} \end{cases} \quad (2)$$

Вероятность правильного решения P тогда определяется по формулам

$$P = \sum_{i=1}^u P_i, \quad$$

где

$$P_i = P(z_i = \max(t_1, \dots, t_u)) \quad (3)$$

$$P_i = \begin{cases} P_i(z_i) = P(t_i = \max(t_1, \dots, t_u)) & \text{при } Q_i > \max(Q_1, \dots, Q_u); \\ 0 & \text{в других случаях.} \end{cases}$$

Вероятность $P_i(z_i)$ можно определить, зная условные законы плотности распределения вероятностей случайных величин t_1, \dots, t_u . Пусть эти величины распределены по нормальному закону и их математические ожидания равны показателям качества Q_1, \dots, Q_u , т.е.

$$f(\bar{t} | \bar{Q}) = (2\pi)^{-\frac{u}{2}} |\Sigma|^{-\frac{1}{2}} \exp[-(\bar{t} - \bar{Q})' \Sigma^{-1} (\bar{t} - \bar{Q})], \quad (4)$$

где

$$\Sigma = \begin{pmatrix} K_{11} & \dots & K_{1u} \\ \vdots & \ddots & \vdots \\ K_{u1} & \dots & K_{uu} \end{pmatrix},$$

Σ — матрица коэффициентов корреляции величин t_1, \dots, t_u ;

$$t = \begin{pmatrix} t_1 \\ \vdots \\ t_u \end{pmatrix}; \quad \bar{Q} = \begin{pmatrix} Q_1 \\ \vdots \\ Q_u \end{pmatrix},$$

где t и \bar{Q} — векторы оценок t_i и ПК Q_i , соответственно.

Вероятности $P_i(z_i)$, а тем самым и вероятность правильного решения P , можно определить, используя теорию многомерного статистического анализа нормально распределенных величин [6]. Она будет

$$P_i(z_i) = \int_{-\infty}^{\dots} \int_{-\infty}^{\dots} (2\pi)^{-\frac{u(u-1)}{2}} |\rho|^{-\frac{1}{2}} \exp(-\bar{x}' \rho' \bar{x}) dx_1 \dots dx_u,$$

где \bar{x} — $u-1$ -мерный случайный вектор с равным σ математическим ожиданием и корреляционной матрицей ρ ;

$$\Delta Q_{ij} = (Q_i - Q_j) \cdot [K_{ii} + K_{jj} - 2\rho_{ij}(K_{ii} \cdot K_{jj})^{1/2}]^{-1/2} \quad (6)$$

Элементы P_{kj} матрицы P определяются по формуле

$$P_{kj} = \frac{(K_{ii} + K_{jk} - K_{ik} - K_{ij})}{[(K_{ii} + K_{kk} - 2K_{ij})(K_{ii} + K_{jj} - 2K_{ij})]}^{1/2} \quad (7)$$

$$\alpha_{ij} = K_{ij}(K_{ii} \cdot K_{jj})^{-1/2},$$

где α_{ij} — коэффициент корреляции оценок t_i, t_j ($i, j = 1, \dots, U$). Из формулы (5) видно, что вероятность $P_i(t_i)$ является

возрастающей функцией аргументов ΔQ_{ij} ($j = 1, \dots, U; i \neq j$).

Величины ΔQ_{ij} , как видно из формулы (6), при $Q_i > Q_j$ являются возрастающими функциями коэффициентов корреляции

α_{ij} . Следовательно, вероятность $P_i(t_i)$ возрастает с возрастанием зависимости между оценкой i -го показателя качества t_i и остальными оценками t_j ($j = 1, \dots, U; i \neq j$). Кроме того, из формул (5) и (6) видно, что при $K_{ii} = K_{jj}$ выполняется следующее равенство:

$$\lim_{\substack{\alpha_{ij} \rightarrow 1 \\ j=1, \dots, U}} P_i(t_i) = 1. \quad (8)$$

Так как вероятность правильного выбора P является суммой вероятностей $P_i(t_i)$ (формула (2)), то она тоже возрастает с увеличением корреляции α_{ij} ($i, j = 1, \dots, U; i \neq j$) между оценками t_1, \dots, t_U функции показателей качества Q_1, \dots, Q_U .

Аналогичным путем можно показать, что вероятность P является возрастающей функцией от корреляционных коэффициентов α_{ij} и при других множествах R решений τ , встречаемых при выборе качественных признаков опознавающих систем.

Известно, что в большинстве случаев, при использовании той же самой выборки опознавающих объектов ξ_1, \dots, ξ_N для подсчета всех оценок качества t_1, \dots, t_U коэффициенты корреляций $\alpha_{ij} > 0$. В таких случаях вероятность P можно заменить её оценкой P' , которая легко подсчитывается, допуская, что оценки t_1, \dots, t_U независимы. При этом будет выполнено неравенство $P > P'$. Следовательно, если нам надо определить такой объем выборки N , при котором выполняется неравенство $P > P^*$ в этом случае, мы можем воспользоваться существующими методами [1, 2, 3, 4, 5] выбора качественных систем

признаков. При этом при вычислениях вероятность P будет замена её оценкой P' .

В заключение заметим, что полученные в настоящей работе результаты легко обобщаются при использовании предельных теорем для слабо зависимых случайных величин, для случаев биномиальных и Гамма-распределений оценок t_1, \dots, t_U .

Л и т е р а т у р а

1. Paulson E. On the comparison of several experimental categories with a control. Ann.Math.Stat., 1952, 23, 239-46.
2. Bechhofer R.E. A single sample multiple decision procedure for ranking the means of normal populations with known variances. Ann.Math.Statist., 1954, 25, 16-39.
3. Mace A.E. Sample-size determination, 1964, New York, Reinhold Publishing Corporation.
4. Chamber M.L. and Jarrett P. Use of double sampling for selecting best population. Biometrika, 1964, 51, 49-64.
5. А.Ю. Мотузя Алгоритмы выбора полезных признаков опознавающих систем. Рефераты докладов I-го Всесоюзного симпозиума по статистическим проблемам в технической кибернетике. Москва, 1967, ч. II, 9-10. Введение в многомерный статистический анализ. ГИФМЛ, Москва.
6. Т. Андерсон.

Поступила в редакцию
I/X-1967г.