

В. А. СЛЕПЯН

**О ЧИСЛЕ ТУПИКОВЫХ ТЕСТОВ И О МЕРАХ ИНФОРМАТИВНОСТИ  
СТОЛБЦА ДЛЯ ПОЧТИ ВСЕХ БИНАРНЫХ ТАБЛИЦ**

(Представлено академиком С. Л. Соболевым 2 VII 1969)

1. Пусть  $T_{ln}$  — бинарная таблица, имеющая  $l$  строк и  $n$  столбцов, а  $T_{lk}$  — таблица, составленная из  $k$  ( $k \leq n$ ) столбцов таблицы  $T_{ln}$ .

Определение 1. Таблицу  $T_{lk}$  будем называть тестом таблицы  $T_{ln}$ , если она состоит из разных строк.

Определение 2. Тест будем называть тупиковым, если после удаления из него любого столбца, он перестает быть тестом. Понятие теста и тупикового теста дано в работе (1).

Таблице  $T_{ln}$  поставим в соответствие величины  $N(T_{ln})$  и  $N^i(T_{ln})$  — число всех тупиковых тестов и число тупиковых тестов, в которые входит  $i$ -й столбец, соответственно.

В тестовых алгоритмах распознавания образов (2) за меру информативности столбца принимается информационный вес столбца

$$I = N^i(T_{ln}) / N(T_{ln}). \quad (1)$$

Для так называемых алгоритмов «голосования» Ю. И. Журавлевым в качестве меры информативности предложена величина

$$J = \sum_{i=1}^{l-1} \sum_{j=i+1}^l (C_{n-\rho_{ij}}^k - C_{n-1-\tilde{\rho}_{ij}}^k) / \sum_{i=1}^{l-1} \sum_{j=i+1}^l C_{n-\rho_{ij}}^k, \quad (2)$$

где  $\rho_{ij}$  — расстояние Хэмминга между строками  $i$  и  $j$  в таблице  $T_{ln}$ , а  $\tilde{\rho}_{ij}$  — то же расстояние после удаления одного из столбцов.

2. Будем теперь считать  $T_{ln}$  случайной таблицей, а именно таблицей, каждый элемент которой принимает значения 0 и 1 с вероятностью  $1/2$ . Тогда  $N(T_{ln})$ ,  $N^i(T_{ln})$ ,  $I$  и  $J$  являются случайными величинами.

В настоящей заметке для класса таблиц, удовлетворяющих соотношению

$$n \leq 2^{l^\alpha} \quad (\alpha < 1/2) \quad (3)$$

находится асимптотика среднего ( $n, l \rightarrow \infty$ ) случайной величины  $N(T_{ln})$ , а для класса таблиц, удовлетворяющих условию

$$\log l / \alpha \leq n \leq l^\beta \quad (\beta < 2, \alpha = 1/2 - \lambda_0 / \log l, \lambda_0 \rightarrow \infty, \lambda_0 / \log l \rightarrow 0), \quad (4)$$

доказывается, что

$$N/\bar{N} \xrightarrow{p} 1, \quad \ln/2 \log l \xrightarrow{p} 1 \quad (l, n \rightarrow \infty).$$

Для таблиц  $T_{ln}$  любых размеров, если только выполнено одно из соотношений  $k \leq \sqrt[3]{n}$ ,  $l^2 2^{-k} \rightarrow \infty$ , справедлива формула

$$Jn/k \xrightarrow{p} 1 \quad (l, n \rightarrow \infty).$$

Здесь символ  $\xrightarrow{p}$  обозначает, как обычно, сходимость по вероятности.



3. Оценим вероятности тупикового теста

Рассмотрим таблицу  $T_{lk}$ . Дадим некоторые определения.

О п р е д е л е н и е 3. Две неравные строки таблицы  $T_{lk}$  образуют соединение на  $i$ -м столбце, если после удаления  $i$ -го столбца они становятся равными.

Будем говорить, что таблица  $T_{lk}$  имеет соединение на  $i$ -м столбце, если хотя бы одна пара строк  $T_{lk}$  образует соединение на  $i$ -м столбце; таблица  $T_{lk}$  имеет соединение на  $m$  столбцах, если она имеет соединение на каждом из  $m$  столбцов.

Введем обозначения:  $T_{lk}^*$ ,  $T_{lk}'$  — таблица  $T_{lk}$  тест, тупиковый тест соответственно;  $R_l^0$ ,  $R_l^1$  — множество тестов, имеющих соединение на первом столбце, соответственно;  $r_l^i$  — множество тестов, имеющих  $i$  и только  $i$  соединений на первом столбце;  $\mu(A)$  — мощность множества  $A$ .

Пусть  $Q_0$  и  $Q$  — некоторые множества, элементами которых являются таблицы  $T_{lk}$ ,  $T_{(l-1)k}$  соответственно, а  $\alpha_0$  и  $\alpha_1$  — целые числа. Под записью  $\alpha_0 Q_0 \rightarrow \alpha_1 Q_1$  условимся понимать следующее: из  $\alpha_0$  множеств  $Q_0$  путем удаления по одной строке из каждой таблицы можно построить  $\alpha_1$  множеств  $Q_1$ . Для множеств  $r_l^i$ ,  $r_{l-1}^i$ ,  $r_{l-1}^{i-1}$  нетрудно доказать справедливость соотношений

$$2ir_l^i \rightarrow (l - 2i + 1)r_{l-1}^{i-1}, \quad (5)$$

$$(l - 2i)r_l^i \rightarrow 2(2^{k-1} - l + i + 1)r_{l-1}^i. \quad (6)$$

Введенные обозначения позволяют записать равенства

$$R_{l-1}^0 = \sum_{i=0}^{[(l-1)/2]} r_{l-1}^i, \quad R_l^1 = \sum_{i=0}^{[l/2]} r_l^i. \quad (7)$$

На основании (5) и (6) всегда можно найти такое целое число  $A_i > 0$ , что справедливо представление

$$A_i(\alpha_i + \beta_i)r_l^i \rightarrow m_i r_{l-1}^{i-1} \cup n_i r_{l-1}^i \quad (i = 0, 1, \dots, [l/2]),$$

где  $\alpha_i, \beta_i \geq 0$ ;  $\alpha_i + \beta_i = 1$ ;  $m_i, n_i$  — целые положительные числа.

Доказывается, кроме того, что  $\alpha_i, \beta_i$  удовлетворяют уравнению

$$\alpha_i \mu(r_l^i) + \beta_{i-1} \mu(r_{l-1}^{i-1}) = \frac{\mu(R_l^1)}{\mu(R_{l-1}^0)} \mu(r_{l-1}^{i-1}).$$

Отсюда и из формул (7) следует

Л е м м а 1. Существуют такие целые положительные числа  $A_i, B_i$ , что

$$A_i R_l^1 \rightarrow B_i R_{l-1}^0.$$

Следствием леммы 1 является

Л е м м а 2. Вероятность тупикового теста удовлетворяет неравенству

$$p(T_{lk}') \geq p(T_{lk}^*) \prod_{i=1}^k p_i, \quad p_i = 1 - \frac{2^{l-i+1} C_{2^{k-1}}^{l-i+1}}{C_{2^k}^{l-i+1}}.$$

Пусть таблица  $T_{lk}$  образует тест. Зафиксируем  $m$  столбцов.

О п р е д е л е н и е 4. Назовем гирляндой группу из  $i + 1$  строк, образующих между собой соединения на  $i$  столбцах, принадлежащих  $m$  указанным столбцам.

О п р е д е л е н и е 5. Максимальной гирляндой на  $j$ -м столбце назовем гирлянду, образующую соединения на максимальном числе столбцов из  $m$  фиксированных, включая  $j$ -й.

Обозначим через  $X_l^m$  тест  $T_{lk}^*$ , имеющий соединение на  $m$  фиксированных столбцов, в число которых входит первый, а через  $Y_{i \max}$  — тест  $T_{lk}^*$ , содержащий максимальную гирлянду из  $i + 1$  строк на первом столбце.



Тогда

$$p(X_i^m/R_i^0) = \sum_{i=1}^m p(Y_{i \max} X_i^{m-i}/R_i^0) \quad (X_i^0 \equiv 1). \quad (8)$$

Из утверждения леммы 1 и того факта, что вероятность соединения на  $m$  столбцах не увеличивается при уменьшении числа строк таблицы, оцениваются сверху вероятности, входящие в правую часть (8). В результате имеем

$$p(X_i^m/R_i^0) \leq p_1 p(X_i^{m-1}/R_i^0) + \sum_{i=2}^m 2^{k-1} C_{m-1}^{i-1} i! (i-1)! (l2^{-k})^{i+1} p(X_i^{m-i}/R_i^0).$$

Далее методом математической индукции, доказывается

Лемма 3. Вероятность тупикового теста удовлетворяет неравенству  $p(T_{lk}') \leq p(T_{lk}^*) p_1^k (1 + \delta)^k$ , где

$$\delta = \frac{l^3 k}{2^{2k} p_1^2} \frac{1 - (lk^2/2^k p_1)^k}{1 - lk^2/2^k p_1}.$$

Следствием лемм 2 и 3 является

Теорема 1. Если  $k^2 l^{-1} \rightarrow 0$ , то  $p(T_{lk}') \sim p(T_{lk}^*) (1 - \exp(-l2^{-k-1}))$ .

Для тех  $k$ , при которых условие теоремы 1 не выполнено, приведем следующую верхнюю оценку вероятности тупикового теста:

$$p(T_{l(k+1)}') \leq p(\bar{T}_{lk}^*) p(T_{lk}') / p(T_{lk}^*), \quad (9)$$

4. Среднее число тупиковых тестов может быть представлено суммой

$$\bar{N} = \sum_{k=\lceil \log l \rceil}^{\min(n, l-1)} \bar{N}_{k_0} \quad \bar{N}_k = C_n^k p(T_{lk}'). \quad (10)$$

где  $\bar{N}_k$  — среднее число тупиковых тестов длины  $k$  (длина теста — число столбцов в нем). С использованием теоремы 1 и формул (9) и (10) доказывается

Теорема 2. 1) Если  $n \leq l^2 (\log l)^\gamma$  ( $\gamma < 3$ ), то  $\bar{N} \sim \bar{N}_{k_0} + \bar{N}_k$ .

2) Если  $n \leq 2^{l^\alpha}$  ( $\alpha < 1/2$ ), то

$$\bar{N} \sim \bar{N}_{k_0} \left( 1 + \sum_{i=1}^{\lceil l^{\alpha_0} \rceil - k_0} (\eta_{k_0}^{-1} a^{(i-1)/2})^i + \sum_{i=1}^{k_0 - \lceil \log l \rceil} (\eta_{k_0-1} a^{(i-1)/2})^i \right),$$

где  $\bar{N}_k \sim \bar{N}_k = C_n^k p(T_{lk}^*) (1-x)^k$ ;  $x = \exp(-l2^{-k-1})$ ;  $\eta_k = \bar{N}_k / \bar{N}_{k+1}$ ;  $k_0$  — корень уравнения  $\eta_k = 1$ ;  $k_1 = k_0 - 1$  или  $k_0 + 1$ ;  $\eta_{k_0} \geq 1$ ,  $\eta_{k_0-1} \leq 1$ ;  $0 < a < a_0 < 1/2$ ;  $0 \leq a \leq 1/4$ .

В частности: 1) при  $2 \log l \leq n \leq l^\beta$  ( $\beta < 2$ )  $k_0 = \lceil 2 \log l - \log(-\ln(n^{1/2} \log l - 1)) - 2 \rceil$ ; 2) при  $l^\varphi \leq n \leq 2^{l^\alpha}$  ( $\varphi > 2$ )  $a = 1/4$ ,  $k_0 = \lceil 1/2 \log(nl^2) - 1/2 \log \log(nl^2) - 1/2 \rceil$ .

Рассмотрим множество таблиц  $Q_n = \{T_{ln}\}$ , где  $2 \log l \leq n \leq l^\beta$ ,  $l \in \in (l_1, (1+a)l_1)$ ,  $a > 0$ .

Теорема 3. Почти во всех таблицах множества  $Q_n$

$$\bar{N} \sim \bar{N}_{k_0} \quad (l \rightarrow \infty).$$

5. В таблице  $T_{ln}$  пронумеруем столбцы (1, 2, ..., n). Рассмотрим две таблицы, составленные из  $k$  столбцов таблицы  $T_{ln}$ . Пусть в каждую из них входит  $m$  и только  $m$  столбцов с одинаковыми номерами. Обозначим через  $p(m)$ ,  $p'(m)$  вероятность того, что обе таблицы одновременно являются тестами, тупиковыми тестами соответственно. Относительно их справедлива



Лемма 4. 1)  $p(m)$  является неубывающей функцией  $m$ .

2)  $p'(m) \leq p(m) \leq \exp(-l^2 2^{-k} + \beta_m)$ , где

$$\beta_m = \frac{l^2}{2^{2k-m+1}} + \frac{l}{2^{m+1}} + \frac{l^3}{2^{2m+1}} \left( \frac{1}{1-l^2^{-m}} + \frac{2 \exp(l \cdot 2^{-m+1})}{1-q} \right) \quad (q < 1).$$

3)  $p'(m) \leq p^* \sim p(T_{lk}^*) (1-x)^k (1-x^{2^m})^{k-m}$  ( $k^2 l^{-1} \rightarrow 0, l \rightarrow \infty$ ).

Обозначим через  $DN_k$  дисперсию числа тушиковых тестов длины  $k$ . Тогда (3) справедлива формула

$$V_k = \frac{DN_k}{\bar{N}_k^2} = \sum_{m=\max(0, 2k-n)}^k \frac{C_k^m C_{n-k}^{k-m}}{C_n^k} \frac{p'(m)}{p^2(T_{lk}^*)} - 1. \quad (11)$$

Используя лемму 4, можно показать, что справедлива

Лемма 5. Для класса таблиц (4)  $V_{k_0} \rightarrow 0$  ( $l, n \rightarrow \infty$ ).

Далее будем рассматривать только те таблицы класса (4), для которых  $\bar{N} \sim \bar{N}_{k_0}$  ( $l \rightarrow \infty$ ). Из неравенств Чебышева (4) и леммы 5 следует

Теорема 4.  $N/\bar{N} \xrightarrow{p} 1$  ( $l, n \rightarrow \infty$ ).

Для случайной величины  $N^i(T_{ln})$  можно доказать теорему, аналогичную теореме 4, и показать, что  $\bar{N}^i \sim 2 \log l \bar{N} / n$  ( $n, l \rightarrow \infty$ ). Отсюда и из теоремы 4 следует

Теорема 5.  $ln/2 \log l \xrightarrow{p} 1$  ( $n, l \rightarrow \infty$ ).

6. Рассмотрим таблицу  $T_{ln}$ , размеры которой произвольны. Относительно случайной величины  $J$  справедлива

Теорема 6. Если  $k \leq \sqrt[3]{n}$  или  $l^2 2^{-k} \rightarrow \infty$ , то  $Jn/k \xrightarrow{p} 1$  ( $l, n \rightarrow \infty$ ).

Институт математики  
Сибирского отделения Академии наук СССР  
Новосибирск

Поступило  
2 VII 1969

#### ЦИТИРОВАННАЯ ЛИТЕРАТУРА

- <sup>1</sup> И. А. Чегис, С. В. Яблонский, Тр. Матем. инст. им. В. А. Стеклова АН СССР, 51, 270 (1958). <sup>2</sup> А. Н. Дмитриев, Ю. И. Журавлев, Ф. П. Кренделев, Дискретный анализ, в. 7, 3 (1966). <sup>3</sup> В. А. Слепян, Дискретный анализ, в. 12, 50 (1968). <sup>4</sup> Б. В. Гнеденко, Курс теории вероятностей, М., 1961.