

УДК 519.95

ИСПОЛЬЗОВАНИЕ ДЛИН ТУПИКОВЫХ ТЕСТОВ  
ПРИ ОБРАБОТКЕ ТАБЛИЦ

А.Н. Дмитриев

Цель настоящей заметки — дополнить тестовый метод [2] обработки таблиц бинарных символов учетом длин тупиковых тестов. Таблицы содержат информацию о реальных объектах и составлены из нулей и единиц, указывающих на отсутствие или наличие признаков, характеризующих эти объекты. Обработка таблиц употребляется при распознавании образов, классификации, упорядочивании объектов по заданному критерию и в других подобных задачах [3-5]. Длины тупиковых тестов (как было обнаружено в конкретных практических задачах) позволяют строить новые величины, уточняющие и дополняющие решение задач указанного профиля. Содержание заметки состоит из небольшого введения, числовых характеристик таблиц, строк, столбцов, связанных с длинами тупиковых тестов и некоторых рекомендаций по использованию введенных величин.

I. Напомним, что тестом для таблицы  $T$  называется такая ее часть  $t$ , которая получается из  $T$  удалением тех или иных столбцов, и в которой все строки различны; тупиковым называется такой тест, в котором никакая его часть не является тестом [1]; длиной  $l(t)$  теста  $t$  называется число содержащихся в нем столбцов.

Например, в таблице

$$T = \begin{array}{|c|c|} \hline 000\dots 00 \\ 100\dots 01 \\ 010\dots 10 \\ \hline 001\dots 11 \\ \hline \end{array} \begin{array}{l} \\ \\ \\ t_1 \quad t_2 \end{array}$$

ее части  $t_1$  и  $t_2$  образуют тупиковые тесты, причем  $l(t_1) = 3$ , а  $l(t_2) = 2$ . В силу того, что в тупиковом тесте все строки различны, его можно интерпретировать как носителя распознавания строк таблицы, а в силу несжимаемости такого теста его можно рассматривать как первичный элементарный носитель распознавания. Для данной таблицы  $T$  и теста  $t$  разность  $l(T) - l(t)$  можно принимать за меру его выразительности, то есть чем больше  $l(T) - l(t)$ , тем больше "экономичность" распознавания строк  $T$  с помощью  $t$  (по сравнению с другими тестами той же таблицы). Эти соображения составляют основу построений во втором пункте.

2. Пусть  $T = \{a_{ij}\}$  — бинарная (т.е.  $a_{ij} = 0, 1$ ) таблица из  $m$  строк ( $j = 1, 2, \dots, m$ ) и  $n$ -столбцов ( $i = 1, 2, \dots, n$ ). В задачах указанного типа основную роль играют следующие величины, разработанные на тестовой основе [2,3]:

а)  $P_i = \frac{\kappa_i}{\kappa}$  — информационный вес  $i$ -го столбца, где

$\kappa$  — общее число тупиковых тестов;  $\kappa_i$  — нагрузка  $i$ -го столбца, или число вхождений  $i$ -го столбца в тупиковые тесты  $t$  таблицы  $T$ .

б)  $P(\sigma) = \sum_{i=1}^n P_i \sigma_i$  — информационный вес строки  $\sigma = (\sigma_1, \dots, \sigma_n)$  таблицы  $T$ .

Как видно из определения для  $P_i$ , информация о длинах тупиковых тестов  $t$  в операции вычисления  $P_i$  не участвует, поскольку все  $t$  в  $T$  считаются равноправными. Указанная выше основа (п. а) и п. б)) сохраняется, но в дополнение к ней теперь предлагается некоторое расчленение информационного веса столбцов, учитывающее длины тупиковых тестов, а именно: для  $i$ -го столбца таблицы  $T$  ( $i = 1, \dots, n$ ) определим

$$R_i = P_i \left(1 - \frac{l_i}{m}\right), \quad Q_i = P_i \cdot \frac{l_i}{m},$$

где  $\bar{L}_i$  — средняя длина тупикового теста, содержащего  $i$ -й столбец, а  $m$  — число строк. Как и в п. б), для строки  $\sigma = (\sigma_1, \dots, \sigma_n)$  определим числа:

$$R(\sigma) = \sum_{i=1}^n R_i \sigma_i \quad \text{и} \quad O(\sigma) = \sum_{i=1}^n O_i \sigma_i.$$

Числа  $R_i$ ,  $R$  и  $O_i$ ,  $O$  назовем соответственно различающими и отождествляющими весами столбца или строки. По определению  $R_i + O_i = P_i$ ,  $R(\sigma) + O(\sigma) = P(\sigma)$ .

В связи с этим расчленением и практическими нуждами возникает потребность в изучении распределения длин тупиковых тестов и в отыскании класса таблиц с заданным распределением тестов — таблиц, у которых преимущественно "короткие" тесты и т.п.; большое прикладное значение имеет сравнение максимальной длины тупикового теста таблицы со средней длиной тупиковых тестов той же таблицы. Отметим попутно, что последняя равна  $\sum_{i=1}^n P_i$ .

3. Прикладное значение величин различающих и отождествляющих информационных весов состоит в том, что с их помощью более тонко фракционируются на группы столбцы и строки таблицы  $T$ .

Употребляя  $R_i$  и  $O_i$  при анализе столбцов, интерпретатор из общей совокупности признаков может выделить группу признаков, которая либо отождествляет, либо различает объекты в  $T$ . Это подразделение минимизирует число признаков в соответствии с целями обработки таблиц и обнаруживает качественный состав признаков.

Употребление величин  $R$  и  $O$  для фракционирования строк приводит к подразделению такого вида:

а) Оценка строк по величине  $R$  предпочтительна для процедуры упорядочивания объектов по заданному критерию внутри данной таблицы.

б) Оценка строк по величине  $O$  предпочтительна для процедур обнаружения родственных групп исследуемых объектов одной или нескольких таблиц, то есть в задачах компоновки классов.

Данные, получаемые с помощью указанных величин, позволяют естествоиспытателю существенно корректировать этапы решения задач, связанных с обработкой нескольких таблиц, содержательно приближая их к цели исследования.

1. И.А. ЧЕРГИС, С.В. ЯБЛОНСКИЙ. Логические способы контроля электрических схем. — Труды Математического института им. В.А. Стеклова, 1958, т. 51, стр. 270-360.

2. А.Н. ДМИТРИЕВ, Ю.И. ЖУРАВЛЕВ, Ф.П. КРЕНДЕЛЕВ. О математических принципах классификации предметов и явлений. — Дискретный анализ, Новосибирск, 1966, вып. 7, стр. 3-15.

3. А.Н. ДМИТРИЕВ. Некоторые табличные числа. — Дискретный анализ, Новосибирск, 1968, вып. 12, стр. 22-26.

4. Сб. "Алгоритмы и программы вычислительной диагностики психических заболеваний". Новосибирск, 1969, стр. 51-76.

5. Сб. "Алгоритмы и программы решения геологических задач на ЭЦВМ "Минск-2" и "БЭСМ-3М", вып. 2, Алма-Ата, 1969, стр. 30-53.

Поступила в редакцию  
23.9.1970 г.