

СОГЛАСОВАННАЯ ОЦЕНКА ВАЖНОСТИ ПРИЗНАКОВ И ОБЪЕКТОВ
В ЗАДАЧАХ ТАКСОНОМИИ

Ф.Т.Адьялова, М.М.Камиллов

В методе вычисления оценок большое место занимает задача оценки важности признаков [1]. Один из подходов к ее определению базируется на вычислении изменения меры $\mu(S, K_q)$ принадлежности объекта S к классу K_q при отбрасывании оцениваемого признака.

Задачу оценки объектов по важности можно ставить и решать аналогично задаче оценки важности признаков. Эта аналогия опирается на возможности рассмотрения столбцов таблицы данных (признаков) аналогично строкам (объектам) в задаче оценки признаков. Для этого вводится в рассмотрение некоторое разбиение G множества P признаков на l групп и задается некоторая мера $\mu(P, G_p)$ оценки принадлежности p -го признака ($p \in P$) к S -й группе G_p .

После этого оценка важности объекта S определяется как среднее изменение критерия качества разбиения G признаков при исключении S из множества M объектов на множестве $M \setminus S$. Таким образом, все отличие задачи оценки важности объектов от задачи определения важности признаков целиком зависит от специфики введенного критерия разбиения G . Обычно предполагается, что такой критерий определяет меру "связности" признаков внутри групп: чем они более связаны в смысле заранее выбранного коэффициента связи, тем искомый критерий больше. В этой связи формулируется новая экстремизационная задача группировки признаков, которая отличается от известных тем, что в ней существенно использована специфика признаков как элементов группировки. Эта специфика состоит в инвариантности признака к инверсии его значений. Существо подхода состоит в следующем. Строится функционал $I(G)$, с помощью которого оценивается

Группировка $G = \{G_1, \dots, G_l\}$ множества P исходных признаков на заданное число l групп. Функционал $I(G)$ строится так, что его значение тем больше, чем более "тесными" в некотором смысле оказываются признаки, оказавшиеся внутри групп разбиения G . Далее строится приближенный алгоритм его экстремизации. Подход, таким образом, оказывается вполне аналогичным подходу к автоматической классификации объектов. Различие состоит в построении такого критерия, который учитывает отличие в представлениях "сходства" между объектами и признаками.

Рассмотрим построение критерия $I(G)$ для случая обработки булевой таблицы $T_{nm} = \|t_{ij}\|_m^n$ данных.

На множестве $P = \{t_1, t_2, \dots, t_n\}$ столбцов T_{nm} определим вектор $\sigma = \{\sigma_1, \dots, \sigma_n\}$ из булевых переменных. С помощью этого вектора зададим преобразование $T_j = \sigma_j t_{.j}$ ($j = 1, n$) инверсии признака:

$$\sigma_j t_j = \begin{cases} t_{.j}, & \text{если } \sigma_j = 0, \\ \bar{t}_{.j}, & \text{если } \sigma_j = 1, \end{cases} \quad (1)$$

где черта над вектором $t_{.j}$ означает, что значения всех его компонент следует заменить на противоположные.

Пусть задано некоторое семейство Ω подмножеств объектов, каждому элементу которого $\omega \in \Omega$ поставим в соответствие характеристическую функцию сходства двух признаков:

$$r_\omega(t_{.j}, t'_{.j}) = \begin{cases} 1, & \text{если } T_{ij} = T'_{ij} \text{ для всех } i \in \omega, \\ 0 & \text{- в противном случае.} \end{cases}$$

Функция $f(t_{.j}, t'_{.j})$ близости между двумя столбцами таблицы строится как сумма характеристических функций:

$$f(t_{.j}, t'_{.j}) = \sum_{\omega \in \Omega} r_\omega(t_{.j}, t'_{.j}). \quad (2)$$

Тогда искомым критерий $I(G)$ принимает вид:

$$I(G, \sigma) = \sum_{q=1}^l \sum_{j, j' \in G_q} f(t_{.j}, t'_{.j}). \quad (3)$$

Задача группировки признаков формуруется в данном случае как задача максимизации I одновременно по G и по σ . После группировки

признаков по критерию (3) и определения способа вычисления $\mu(P, G_s, M)$, задача нахождения вектора оценок важности объектов в полной аналогии с тем, как находится важность признаков, сводится к подсчету величин:

$$\{\mu(P, G_s, M \setminus S); \forall S \in M, P \in K_s, s = 1, \dots, l\}. \quad (4)$$

Тогда оценка β_j для j -го объекта определяется как величина:

$$\beta_j = \frac{1}{S} \sum_{S=1}^l \sum_{P \in G_s} |\mu(P, G_s, M) - \mu(P, G_s, M \setminus S)|. \quad (5)$$

Пусть на таблице $T = \|t_{ij}\|_{m \times n}$ независимо получены оценки $\alpha = (\alpha_1, \dots, \alpha_n)$ и $\beta = (\beta_1, \dots, \beta_m)$ соответственно важности признаков и объектов. Каждый из этих рядов ассоциирует другие оценки:

$$\beta'_j = \sum_{i=1}^m t_{ij} \cdot \alpha_i; \quad \alpha'_i = \sum_{j=1}^n t_{ij} \cdot \beta_j. \quad (6)$$

Возникает естественная задача - найти такую пару векторов (α, β) , которая или в точности удовлетворяет соотношениям (6) или, если это невозможно, обеспечивает минимум невязки несоблюдения этих соотношений. Эту задачу мы будем называть задачей согласования оценок важности объектов и признаков. Алгоритм решения, который предлагается в настоящем докладе, состоит в построении следующей итерационной процедуры:

1) на основе матрицы T и разбиений K и G ее строк и столбцов находятся начальные векторы (α^0, β^0) ;

2) на t -м шаге итерации строятся две матрицы:

$$\left. \begin{aligned} T_{\alpha}^t &= \|t_{ij} \cdot \beta_j^{t-1}\|, \\ T_{\beta}^t &= \|t_{ij} \cdot \alpha_i^{t-1}\|. \end{aligned} \right\} \quad (7)$$

По матрице T_{α}^t вычисляется вектор α^{t+1} , а по матрице T_{β}^t - вектор β^{t+1} . Алгоритм прекращает работу после выполнения заданного числа α итераций.

В качестве искомой пары (α^*, β^*) выбирается такая, которая обеспечивает минимум критерия невязок: $I(\alpha^*, \beta^*) = \min_{\alpha, \beta} \left\{ \sum_i (\alpha_i - \sum_j t_{ij} \cdot \beta_j)^2 + \sum_j (\beta_j - \sum_i t_{ij} \alpha_i)^2 \right\}$ из заданного множества пар $\{(\alpha^0, \beta^0), (\alpha^1, \beta^1), \dots, (\alpha^{\alpha}, \beta^{\alpha})\}$.

Л и т е р а т у р а

1. ЖУРАВЛЕВ Ю.И. и др. Алгоритмы вычисления оценок и их применение. -Ташкент: ФАН, 1974.

Поступила в ред.-изд.отд.
22 марта 1984 года