

УДК 621.391:512.2

**О НЕКОТОРЫХ ЗАДАЧАХ ТАКСОНОМИИ (ГРУППИРОВКИ) ОБЪЕКТОВ**

В.И. Котиков

В данной статье предложены методы оптимизации дисперсионного и дискриминантного критериев качества таксономии, а также рассмотрена постановка ряда других задач как задач таксономии. Приведены алгоритмы минимизации исходного описания в задачах таксономии.

**§ I. Постановка задачи.**

**Оптимизация дисперсионного критерия**

Допустим, в  $\mu$ -мерном исходном пространстве признаков  $X = \{X_1, \dots, X_p\}$  дано описание  $N$  объектов  $\{a_1, \dots, a_N\}$ . Необходимо разбить  $N$  объектов на  $K$  непересекающихся таксонов (групп, подмножеств) так, чтобы оптимизировать при этом некоторый критерий. Это обычная постановка задачи таксономии объектов [1,2].

Предположим, что вычислена матрица взаимных расстояний  $\{p_{ij}\}$  между объектами в пространстве  $X$ . Тогда задача заключается в минимизации величины

$$D = \sum_{\ell=1}^K \sum_{i=1}^{N_\ell} (p_{i0}(\ell))^2, \quad (1)$$

где  $N_\ell$  - число объектов в  $\ell$ -м таксоне, а  $p_{i0}(\ell)$  - расстояние  $i$ -го объекта  $\ell$ -го таксона до центра "тяжести" этого же таксона.

В работе [3] данная задача сведена к задаче квадратичного, а также целочисленного линейного программирования. Здесь же будет дано более простое решение.

Будем минимизировать несколько иной дисперсионный критерий:

$$D^* = \sum_{\ell=1}^K \sum_{i=1}^{N_\ell} (p_{qi}(\ell))^2, \quad (2)$$

где  $p_{qi}(\ell)$  - расстояние  $q$ -го объекта  $\ell$ -го таксона до  $i$ -го объекта того же таксона; причем  $q$ -й объект это такой объект  $\ell$ -го таксона, который наиболее близок к центру "тяжести" того же таксона.

Очевидно, что если в каждом из  $K$  таксонов будет хотя бы один объект, достаточно близкий к центру "тяжести" своего таксона, то решения задачи по критерию  $D$  и критерию  $D^*$  будут близки.

Обозначим  $p_{ij}^2 = t_{ij}$ . Введем целочисленные переменные (неизвестные):

$$x_{ij} = \{0,1\}; \quad i=1, \dots, N; \quad j=1, \dots, N+1 \quad (3)$$

Причем будем полагать "затраты", равными:

$$\left. \begin{aligned} t_{i(N+1)} &= 0, \quad i=1, \dots, N; \\ t_{ij} &= p_{ij}^2; \quad i, j=1, \dots, N. \end{aligned} \right\} \quad (4)$$

Столбец  $(N+1)$  матрицы  $\{x_{ij}\}$  является фиктивным и служит для контроля числа получившихся таксонов.

Определим функцию

$$Q^* = \sum_{i=1}^N \sum_{j=1}^{N+1} t_{ij} x_{ij} \quad (5)$$

при ограничениях на неизвестные:

$$\sum_{i=1}^N x_{ij} = b_j; \quad j=1, \dots, N+1; \quad (6)$$

где  $b_j = 1; j=1, \dots, N$  и  $b_{N+1} = K$ ;

$$\sum_{j=1}^N x_{ij} - N x_{i(N+1)} \leq 0; \quad i=1, \dots, N. \quad (7)$$

Число ненулевых строк в получившейся после решения матрице  $\{x_{ij}\}$  будет равно числу таксонов. Это гарантируется неравенствами (7), а также равенством

$$\sum_{i=1}^N x_{i(N+1)} = b_{N+1} = K.$$

Равенства (6) гарантируют попадание каждого объекта только в один таксон.

Если  $q$  - я строка в матрице  $\{x_{ij}\}$  окажется ненулевой и  $\{x_{q1}=1, x_{q2}=1, x_{qe}=1\}$ , то это означает, что объекты  $a_q, a_{q1}, a_{q2}, a_e$  и  $a_e$  входят в один таксон, к "центру" которого наиболее близким является объект  $a_q$ .

ЛЕММА. Объект  $a_q$  является на  $\epsilon$  более близким к центру "тяжести"  $e$  - го таксона.

ДОКАЗАТЕЛЬСТВО. Дисперсия

$$\sigma = (1/N_e) \sum_{i=1}^{N_e} (x_i - \mu)^2$$

достигает минимума при

$$\mu = \bar{x} = (1/N_e) \sum_{i=1}^{N_e} x_i;$$

это следует из уравнения  $(\partial\sigma/\partial\mu) = 0$ , то есть если  $\mu = \bar{x}$ , то

$$\sum_{i=1}^{N_e} p_{iq}^2 \geq \sum_{i=1}^{N_e} p_{i\bar{x}}^2$$

Так как функция (5) подвергается минимизации, то утверждение леммы справедливо.

Таким образом, нами доказана следующая

ТЕОРЕМА I. Минимизация целевой функции (5) при ограничениях (3), (6), (7) и равенствах (4) экви-

валентна минимизации критерия  $D^*$ .

Данная задача есть задача целочисленного линейного программирования [4].

Для случая большого  $N$  сформируем алгоритм, позволяющий достигать локальный минимум функции  $D^*$ . Структура алгоритма ("процедуры  $\gamma$ ") аналогична структуре многих итерационных методов таксономии [6].

Определим "процедуру  $\gamma$ " следующим образом: допустим, на  $i$ -м шаге процедуры осуществлено какое-то разбиение множества  $N$  объектов на  $K$  подмножеств и в каждом из этих подмножеств найден объект (эталон)  $a_q$ , сумма квадратов расстояний которого до всех остальных объектов данного подмножества минимальна. Тогда на  $(i+1)$ -м шаге алгоритма происходит новое перераспределение всех  $N$  объектов по  $K$  таксонам: объект относится к тому подмножеству, к эталону которого  $a_q$  он наиболее близок. Для вновь полученных  $K$  подмножеств по вышеописанному правилу находятся новые эталоны  $a_q$ . Далее осуществляется  $(i+2)$ -й шаг. Из самого алгоритма непосредственно вытекает

ТЕОРЕМА 2. Если  $D_i^*$  и  $D_{i+1}^*$  есть значения критерия качества таксономии  $D^*$  на  $i$ -м и  $(i+1)$ -м шагах "процедуры  $\gamma$ ", то справедливо соотношение  $D_{i+1}^* \leq D_i^*$ .

Хотя критерий  $D^*$  и "слабее" в статистическом смысле критерия  $D$ , но имеет более широкую область применения, так как не требует, чтобы величины  $\{p_{ij}\}$  обязательно удовлетворяли всем аксиомам метрического пространства: величина  $p_{ij}$  может просто численно выразить степень "различия" между объектами  $a_i$  и  $a_j$ . Кроме того, в некоторых задачах понятие центра "тяжести" (оценки математического ожидания) таксона физически просто не реализуемо.

Рассмотрим некоторые иные задачи, которые могут быть решены с помощью описанного метода.

Группировка признаков  $\{x_1, \dots, x_p\}$ . Если здесь за объекты принять сами признаки, а в качестве расстояний между ними  $\{p_{ij}\}$  взять  $\{(1 - |r_{ij}|)\}$ , где  $|r_{ij}|$  - модуль коэффициента корреляции, то мы можем осу-

пествить разбиение всей совокупности исходных признаков  $X$  на  $K$  групп, в каждой из которых они между собой коррелируют достаточно сильно, а между группами слабо.

#### Минимизация описания объектов

Допустим, что каждый из  $\rho$  исходных признаков  $\{X_1, \dots, X_\rho\}$  имеет  $n_e$  числовых градаций. Необходимо из имеющихся  $N$  градаций вывал ( $\sum_{e=1}^{\rho} n_e = N$ ) выбрать те  $K$  градаций ( $K < N$ ), при которых был бы минимум среднеквадратичного отклонения описания объектов в новом пространстве  $\tilde{X}$  от их описания в исходном пространстве  $X$ . Пусть всего  $M$  объектов. Если объект  $a_i$  имел по признаку  $X_j$  значение  $x_{ij}$ , то после удаления этой градации он примет значение  $\tilde{x}_{ij}$  ближайшей градации из неудаленных у  $X_j$ . Этого требует критерий минимума среднеквадратичного отклонения.

Если значения  $k$ -й градации  $x_j^k$  признака  $X_j$  имело  $m_k$  из исходных объектов, а значение  $l$ -й градации  $x_j^l$  -  $m_l$  объектов, то в качестве квадрата расстояния нужно взять:

$$\rho_{kl}^2 = m_l (x_j^k - x_j^l)^2 \text{ и } \rho_{lk}^2 = m_k (x_j^k - x_j^l)^2.$$

Расстояния между градациями различных признаков нужно положить равными  $\infty$ .

Справедливость подобной постановки вытекает из равенства

$$\sum_{i=1}^M \sum_{j=1}^{\rho} (x_{ij} - \tilde{x}_{ij})^2 = \sum_{j=1}^{\rho} \sum_{i=1}^M (x_{ij} - \tilde{x}_{ij})^2.$$

Ненулевые строки матрицы найденных неизвестных  $\{\rho_{ke}\}$  будут соответствовать тем  $K$  из исходных градаций пространства  $X$ , которые нужно оставить.

#### Квантование сигналов по уровню

Допустим, имеются  $M$  сигналов  $y = f(t)$ , заданных своими отсчетами по оси времени. Амплитудное значение при каждом отсчете времени может принимать одно из  $N$  значений. Нужно из  $N$  выбрать такие  $K$  градаций ( $K < N$ ), чтобы среднеквадратичное отклонение семейства  $M$  сигналов от их исходных было

бы минимальным. Пусть число всех отсчетов по всем сигналам, имеющих значение градации  $y_k$ , равно  $m_k$ , а для градации  $y_l$  -  $m_l$ , соответственно. Тогда положим

$$\rho_{kl}^2 = m_l (y_k - y_l)^2 \text{ и } \rho_{lk}^2 = m_k (y_k - y_l)^2.$$

Аналогичным образом нетрудно убедиться, что решение этой задачи как задачи таксономии будет минимизировать вышеупомянутый критерий. Не нулевые строки матрицы  $\{\rho_{ke}\}$  будут соответствовать тем  $K$  из  $N$  исходных градаций уровня  $y$ , которые нужно оставить.

#### § 2. Оптимизация дискриминантного критерия

Как и ранее, будем предполагать, что  $N$  объектов заданы своим описанием в  $\rho$ -мерном признаковом пространстве  $\{X_1, \dots, X_\rho\}$ . Необходимо разбить  $N$  объектов на два класса  $A$  и  $B$  так, чтобы максимизировать значение критерия Фишера

$$F = \frac{(\bar{y}_a - \bar{y}_b)^2}{\sigma_a^2 + \sigma_b^2} \quad (8)$$

по наилучшему дискриминантному вектору  $y = \alpha_1 X_1 + \dots + \alpha_\rho X_\rho$ . Величины  $\bar{y}_a$  и  $\bar{y}_b$  есть центры "тяжести" проекции выборки классов  $A$  и  $B$  на вектор  $y$ , а  $\sigma_a$  и  $\sigma_b$  - оценки соответствующих дисперсий по вектору  $y$ .

Допустим, имеется произвольный вектор  $y = \alpha X$ . Значения проекций  $N$  исходных объектов (точек) на  $y$  равны  $y_1, y_2, \dots, y_N$ , соответственно.

Сначала рассмотрим одномерную задачу установления такого порога  $y^0$  (границы между классами) на фиксированном векторе  $y$ , который бы доставлял  $F = F_{max}$ . Данную задачу будем решать методом "ветвей и границ" [4].

Относительно любого порога  $y^i$  расположение классов будем считать следующим:

$$y_l < y^i; y_l \in A \text{ и } y_k \geq y^i; y_k \in B.$$

Допустим, для произвольных порогов  $y^i$  и  $y^j$  у нас вычислены следующие характеристики  $\langle F^i, \bar{y}_a^i, \bar{y}_b^i, \sigma_a^i, \sigma_b^i, N_a^i \rangle$  и  $\langle F^j, \bar{y}_a^j, \bar{y}_b^j, \sigma_a^j, \sigma_b^j, N_a^j \rangle$ , соответственно.  $N_a^i$  - объем выборки

класса  $A$  при границе  $y^i$  на векторе  $y$ ,  $N_0^i = N - N_0^i$ . Пусть  $F^i < F^j$  и  $y^i < y^j$ .

Вычислим верхнюю границу значения  $F_{ij}$  для порога на отрезке  $(y^i, y^j)$ .

Между порогами  $y^j$  и  $y^i$  находится  $(N_0^j - N_0^i)$  точек. Нетрудно убедиться в том, что величина  $F$  достигнет максимума на отрезке  $(y^i, y^j)$ , если все  $(N_0^j - N_0^i)$  этих точек будут иметь одну и ту же координату, почти совпадающую с  $y^i$ .

Тогда верхняя граница значения  $F_{ij}$  для порога на отрезке  $(y^i, y^j)$  вычисляется так:

$$\hat{F}_{ij} = (\hat{\mu}_a - \hat{\mu}_b)^2 / (\hat{\sigma}_a^2 + \hat{\sigma}_b^2),$$

где:

$$\hat{\mu}_a = \frac{1}{N_a^j} (N_a^i \bar{y}_a^i + (N_a^j - N_a^i) y^i);$$

$$\hat{\mu}_b = \bar{y}_b^j; \quad \hat{\sigma}_a^2 = \hat{\sigma}_{aH}^2 + \hat{\sigma}_{aC}^2; \quad \hat{\sigma}_b^2 = \sigma_c^2;$$

$$\hat{\sigma}_{aH}^2 = (1/N_a^j) (N_a^i (\bar{y}_a^i - \mu_a)^2 + (N_a^j - N_a^i) (y^i - \mu_a)^2);$$

$$\hat{\sigma}_{aC}^2 = (1/N_a^j) (\sigma_a^i N_a^i + 0);$$

$$\mu_a = (\bar{y}_a^i + y^i) / 2.$$

Если  $(N_a^j - N_a^i) \leq 1$ , то  $\hat{F}_{ij} = \max(F^i, F^j)$ . Координату нового порога  $y^q$  на отрезке  $(y^i, y^j)$  можно выбрать пропорционально величинам  $F^i$  и  $F^j$ . Поскольку получено выражение для определения верхней границы значения критерия на отрезке и дано правило дробления отрезков, то возможно применение метода "ветвей и границ". В силу конечности выборки решение оптимизационной одномерной задачи может быть получено точно.

Перейдем теперь к решению задачи в многомерном пространстве.

Допустим, для произвольного вектора  $y_i$  мы решили одномерную задачу таксономии, получив множества  $\{A_i, B_i\}$  и значение критерия  $F_i$ . Для множеств  $A_i$  и  $B_i$  в простран-

стве  $X$  найдем дискриминантный вектор  $y_{i+1}$ , максимизирующий критерий  $F$ . Вектор  $y_{i+1}$  целесообразно находить методом, описанным в работе [5].

Для полученного вектора  $y_{i+1}$  решим одномерную задачу таксономии и получим  $F_{i+1} \leftarrow \{A_{i+1}, B_{i+1}\}$ .

Легко может быть доказана

ТЕОРЕМА 3. Если выполняется неравенство  $y_{i+1} \neq y_i$ , то справедливо соотношение  $F_{i+1} > F_i$ .

Таким образом, имеем сходимость к экстремуму функционала  $F$ .

На основе вышеописанной процедуры и процедуры работы [4] мы несложным образом можем осуществить таксономию на два класса  $N$  исходных объектов по вектору  $y$ , ортогональному первому, и так далее до тех пор, пока значение качества таксономии  $F$  по вектору  $y$ , ортогональному всем предыдущим, не станет достаточно плохим. Это можно использовать для общей таксономии  $N$  объектов на большее число классов.

### § 3. Минимизация исходной системы признаков

Необходимо из исходной системы  $p$  признаков  $\{X_i\}$  выбрать такие  $m$  признаков  $\{X_j\}$  ( $m < p$ ), в  $m$ -мерном пространстве которых качество группировки объектов на  $K$  таксонов было бы наилучшим. Это так называемая задача комбинированного типа "SX" [6]. Рассмотрим ее решение для двух частных случаев.

#### Дисперсионный критерий

Качество таксономии здесь будем измерять величиной критерия  $D^*$ . Квадрат расстояния между  $i$ -м и  $j$ -м объектами, например, в евклидовой метрике измеряется таким образом:

$$\rho_{ij}^2 = \sum_{l=1}^p (\rho_{il}^2 - \rho_{jl}^2)^2 = \sum_{l=1}^p \rho_{ijl}^2.$$

Если  $t_{ij} = \rho_{ij}^2$ , то  $t_{ij} = \sum_{l=1}^p t_{ijl}$ , где  $t_{ijl} = \rho_{ijl}^2$ .

Введем целочисленные переменные  $z_{ijl} = \{0, 1\}$ , где  $z_{ijl} =$

$= 1, \dots, (N+1)$ ;  $j = 1, \dots, (N+1)$ ;  $l = 1, \dots, (p+1)$ .

Неизвестные  $z_{(N+1)j\ell}$  введены для контроля числа признаков  $\{X_j\}$ .

Теперь наша цель заключается в минимизации функционала

$$Q^* = \sum_{l=1}^N \sum_{j=1}^N \sum_{\ell=1}^p t_{ij\ell} z_{ij\ell} \quad (9)$$

при ограничениях:

$$\left. \begin{aligned} \sum_{l=1}^N \sum_{j=1}^{N+1} z_{ij\ell} - (N+K)z_{(N+1)\ell(p+1)} &= 0; \\ \ell &= 1, \dots, p; \\ \sum_{l=1}^p z_{(N+1)\ell(p+1)} &= m; \end{aligned} \right\} \quad (10)$$

$$\left. \begin{aligned} \sum_{l=1}^N \sum_{j=1}^p z_{ij\ell} &= b_j; \quad j=1, \dots, N+1; \\ b_j &= m \quad \text{при } j=1, \dots, N; \quad b_{(N+1)} = mK \end{aligned} \right\} \quad (11)$$

$$\sum_{j=1}^N \sum_{l=1}^p z_{ij\ell} - N \sum_{\ell=1}^p z_{i(N+1)\ell} \leq 0; \quad i=1, \dots, N. \quad (12)$$

Трактовка переменных и ограничений аналогична задаче для критерия  $D^*$  второго параграфа.

Задача  $\langle (9), (10), (11), (12) \rangle$  есть задача целочисленного линейного программирования [3].

Для случая больших  $N$  и  $p$  может быть предложен такой алгоритм. Сначала задача таксономии решается в исходном  $p$ -мерном пространстве  $X$ . Заметим, что функционал (9) можно записать и в таком виде:

$$Q^* = \sum_{\ell=1}^p Q_{\ell}^*, \quad \text{где} \quad Q_{\ell}^* = \sum_{i=1}^N \sum_{j=1}^N t_{ij\ell} z_{ij\ell}$$

Далее, исходя из полученных таксонов, для каждого из  $\{X_{\ell}\}$  определим величину  $Q_{\ell}^*$ . Так как величины  $Q_{\ell}^*$  входят в  $Q^*$  чисто аддитивно, то "наилучшие"  $m$  признаков  $\{X_j\}$  будут те, которые обладают меньшими величинами  $\{Q_j^*\}$ . Таким образом, здесь наилучшей подсистемой из  $m$  признаков будет та, которая вносит "наилучший" вклад в оптимизацию критерия качества

$Q^*(D^*)$  таксономии во всем исходном признаковом пространстве.

### Дискриминантный критерий

Допустим, в  $p$ -мерном пространстве  $\{X_{\ell}\}$  на основе методики, изложенной в предыдущем параграфе, получены разбиения объектов на множества  $A$  и  $B$  по  $n$  ортогональным дискриминантным функциям  $\{y_1, \dots, y_n\}$ . Величина  $n$  выбирается такой, что качество таксономии  $F$  для  $y_{(n+1)}$  становится уже достаточно "плохим".

Этот факт можно трактовать как существование в исходном пространстве  $n$  возможных "различных" вариантов таксономии объектов на два класса  $A$  и  $B$ . Так как каждый  $y_{\ell}$  из полученных векторов  $\{y_1, \dots, y_n\}$  обладает качеством таксономии по нему, равным  $F_{\ell}$ , и является "наилучшим" собственным вектором некоторой матрицы  $H(y_{\ell})$ , то мы для выбора оптимальной  $m$ -мерной подсистемы  $\{X_j\}$  можем воспользоваться алгоритмом, изложенным в работе [4].

В заключение автор выражает благодарность Н.Г. Загоруйко за обсуждение данной работы.

### Л и т е р а т у р а

1. ДОРОЖНИК А.А. Алгоритмы автоматической классификации (обзор). - "Автоматика и телемеханика", М., 1971, № 12.
2. БИКИНА В.Н., ЗАГОРУЙКО Н.Г. Количественные критерии качества таксономии и их использование в процессе принятия решений. - "Вычислительные системы", Новосибирск, "Наука", 1969, вып. 36.
3. КОТИКОВ В.И. Об одном методе таксономии множества объектов и параметров. - "Вычислительные системы", Новосибирск, 1971, вып. 46.
4. КОРБУТ А.А., ФИНКЕЛЬШТЕЙН В.Д. Дискретное программирование. М., "Наука", 1969.
5. КОТИКОВ В.И. Оптимизация критерия Фишера-Уилкса у сокращение исходной системы описания в задачах распознавания образов. - Настольный сборник, стр. 136-142.
6. ЗАГОРУЙКО Н.Г. Одновременный поиск эффективной системы признаков и наилучшего варианта таксономии (алгоритм " "). - "Вычислительные системы", Новосибирск, "Наука", 1969, вып. 36.

Получена в ред.-изд.отд.

1. I. 1972 г.