

ОБ ОДНОМ МЕТОДЕ ПОСТРОЕНИЯ ОПТИМАЛЬНЫХ КЛАССИФИКАЦИЙ

А. И. ЗЕНКИН, А. А. ЗЕНКИН

Вычислительный центр АН СССР, Москва, СССР

Рассматривается задача классификации, являющаяся одной из центральных в теории распознавания и прогнозирования. В работе предложен метод разбиений произвольных множеств объектов на классы в предположении, что априорное описание этих классов и их обучающие выборки — отсутствуют. Полученная классификация представляет собой решение некоторой задачи минимизации определенного типа функционалов, характеризующих качество классификации. Метод реализован в виде программного комплекса на языке АЛГОЛ-60 для ЭВМ БЭСМ-6.

§ 1. Вычисление расстояний между объектами

Пусть R — множество всех действительных чисел и T — прямоугольная матрица (таблица чисел из R) фиксированной размерности $M \times N$:

$$(1) \quad T = \{t_{ij}\}, \quad t_{ij} \in R, \quad i = 1, 2, \dots, M, \quad j = 1, 2, \dots, N.$$

С таблицей T однозначно связаны упорядоченная последовательность её строк

$$(2) \quad S = \{s_1, s_2, \dots, s_M\}$$

такая, что

$$(3) \quad \forall i ((s_i \in S) \rightarrow (s_i = (t_{i1}, t_{i2}, \dots, t_{iN}))),$$

и последовательность её столбцов

$$(2a) \quad S^+ = \{s_1^+, s_2^+, \dots, s_N^+\}$$

такая, что

$$(3a) \quad \forall j ((s_j^+ \in S^+) \rightarrow (s_j^+ = (t_{1j}, t_{2j}, \dots, t_{Mj}))).$$

Строки (2) таблицы T будем также называть *объектами* (или *векторами*), а элементы строк (3) — *признаками объекта* (или *компонентами вектора*).

Таким образом, задание числовой таблицы T (1) равносильно заданию упорядоченного множества S из M объектов (2), каждый из которых характеризуется набором из N признаков (3).

Для характеристики „близости”, „сходства” или расстояния между строками исходной таблицы T введём ряд оценок.

1. *Оценки расстояний между строками.* Введём в рассмотрение N -мерный вектор весов признаков

$$(4) \quad \bar{w} = \bar{w}(w_1, w_2, \dots, w_N), \quad 0 \leq w_j \leq 1, j = 1, 2, \dots, N,$$

так что w_j — вес j -го признака таблицы T .

Пусть $s_p, s_q \in S$ — две произвольные строки из (2). Исходя из „физического” смысла рассматриваемого ниже класса задач, за меру близости этих двух строк по j -му признаку удобно взять некоторую функцию от абсолютной величины разности значений этого признака для объектов s_p и s_q :

$$(5) \quad r_j(p, q) = w_j f(|t_{pj} - t_{qj}|).$$

Следуя методу вычисления оценок, описанному в [1], введём числовой N -мерный вектор

$$(6) \quad \bar{\varepsilon} = \bar{\varepsilon}(\varepsilon_1, \varepsilon_2, \dots, \varepsilon_N), \quad 0 \leq \varepsilon_j \leq 1, j = 1, 2, \dots, N,$$

— вектор пороговых оценок по признакам, и определим явно функцию (5) с помощью следующего условия

$$(7) \quad r_j(p, q) = \begin{cases} 0, & \text{если } |t_{pj} - t_{qj}| > \varepsilon_j, \\ w_j, & \text{если } |t_{pj} - t_{qj}| \leq \varepsilon_j. \end{cases}$$

Расстояние $R(p, q)$ между строками s_p и s_q определим теперь следующим образом

$$(8) \quad R(p, q) = \frac{1}{N} \sum_{j=1}^N r_j(p, q).$$

Из (8) следует, что для любых $s_p, s_q \in S$ имеет место неравенство

$$(8a) \quad 0 \leq R(p, q) \leq 1.$$

Введём пороговую оценку „близости” строк

$$(8b) \quad 0 \leq \delta_s \leq 1$$

и на множестве всех пар объектов из (2) определим бинарную функцию „близости” для строк по порогу δ_s :

$$(9) \quad \Gamma(p, q) = \begin{cases} 0, & \text{если } R(p, q) < \delta_s, \\ 1, & \text{если } R(p, q) \geq \delta_s. \end{cases}$$

Бинарная функция (9) имеет следующую содержательную интерпретацию: если усреднённое количество „совпадающих” по ε -порогу признаков больше порога δ_s , то строка s_p „голосует” за „сходство” со строкой s_q (или, что

то же, строка s_q „голосует” за „сходство” со строкой s_p); в противном случае количество „голосов” за „сходство” строк s_p и s_q по порогу δ_s равно нулю.

2. *Оценка расстояния от строки s_p до множества K .* Пусть $K \subset S$ — некоторое подмножество строк из (2)

$$K = \{s_1', s_2', \dots, s_{M'}'\}, \quad M' \leq M.$$

Задача расстояние от произвольного объекта $s_p \in S$ до заданного множества объектов K примем следующую величину:

$$(10) \quad R(p; K) = \frac{1}{M'} \sum_{s_q \in K} \Gamma(p, q),$$

которую можно интерпретировать как усредненное по числу объектов в K количество „голосов” в пользу принадлежности рассматриваемого объекта s_p к множеству K .

Введём пороговую оценку

$$(10a) \quad 0 \leq \delta_K \leq 1$$

и на множестве всех подмножеств объектов из (2) определим бинарную функцию „близости” для произвольной строки $s_p \in S$ и произвольного подмножества $K \subset S$ по порогу δ_K :

$$(11) \quad \Gamma(p; K) = \begin{cases} 0, & \text{если } R(p; K) < \delta_K, \\ 1, & \text{если } R(p; K) \geq \delta_K. \end{cases}$$

Введём теперь обобщенный вектор пороговых оценок

$$\bar{u} = \bar{u}(\varepsilon_1, \varepsilon_2, \dots, \varepsilon_N, \delta_s, \delta_K) = \bar{u}(\bar{\varepsilon}, \delta_s, \delta_K)$$

и следующее определение принадлежности объекта s_p множеству K объектов из S .

Определение 1. Объект $s_p \in S$ принадлежит множеству объектов K по фиксированному порогу \bar{u} , если $\Gamma(p; K) = 1$; в противном случае $s_p \notin K$.

§ 2. Статическая классификация Задача классификации без учителя

Фиксируем вектор $\bar{u}(\bar{\varepsilon}, \delta_s, \delta_K)$ пороговых оценок и рассмотрим задачу о разбиении множества S последовательности строк (2) таблицы T на классы (число которых не задаётся априори).

Определение 2. Классификацией K упорядоченного множества объектов S по заданному вектору $\bar{u}(\bar{\varepsilon}, \delta_s, \delta_K)$ пороговых оценок назовём упорядоченную (вообще говоря, произвольно) последовательность

$$(12) \quad K = \{K_1, K_2, \dots, K_L\}$$

подмножеств множества S , удовлетворяющую следующим условиям:

(1) $\forall i, j ((i \neq j) \rightarrow (K_i \cap K_j = \emptyset))$.

$$(2) \bigcup_{i=1}^L K_i = S.$$

(3) С точностью до порядка в (12) K не зависит от порядка строк в исходной таблице T .

(4) $\forall s_p \in S ((s_p \in K_i) \rightarrow \neg \exists j \neq i (R(p; K_j) > R(p; K_i)))$.

Условия (1) и (2) вытекают из определения понятия разбиения множества на классы; условие (3) исключает зависимость классификации от порядка (как правило, случайного) ввода исходной информации и означает следующее: если S' — любая последовательность, отличающаяся от S только порядком элементов, и $K' = \{K'_1, K'_2, \dots, K'_L\}$ — классификация S' по тому же вектору $\bar{u}(\epsilon, \delta_s, \delta_K)$, то $L' = L$ и существует такое 1-1 соответствие φ между индексами последовательностей K и K' , что

$$\forall i ((i = \varphi(i')) \rightarrow (K_i \equiv K_{i'})), \quad i = 1, 2, \dots, L.$$

Наконец, условие (4) можно содержательно интерпретировать как устойчивость классификации K по принадлежности любой строки „своему” классу или, другими словами, для классификации K понятие „объект $s_p \in K_j$ ” и „класс K_j является ближайшим (в смысле (11)) к объекту s_p ” должны совпадать.

Классификацию, удовлетворяющую введенному определению, мы будем строить в два этапа. Любую строку $s_p \in S$ будем называть в процессе классификации „отнесенной”, если она уже отнесена к какому-либо классу, и „неотнесенной” — в противном случае. Очевидно, что до начала классификации любая строка $s_p \in S$ является неотнесенной, а после окончания классификации — отнесенной.

На I этапе строим так называемую предварительную классификацию

$$(13) \quad \tilde{K} = \{\tilde{K}_1, \tilde{K}_2, \dots, \tilde{K}_{L'},\}$$

согласно следующему простейшему алгоритму.

Заносим в I класс строку s_1 и все строки $s_q \in S$, удовлетворяющие условию (9):

$$(14) \quad \Gamma(1; q) = 1.$$

Очевидно, что условие (14) однозначно определяет некоторый класс \tilde{K}_1 . Если после этого в (12) имеются неотнесенные строки, то неотнесенную строку с наименьшим индексом, — пусть это будет, например, s_{l_2} , — и все неотнесенные строки $s_q \in S$, удовлетворяющие условию:

$$\Gamma(l_2, q) = 1,$$

заносим во II класс и формируем \tilde{K}_2 . Эта процедура повторяется до тех пор, пока в (2) не останется неотнесенных строк. В результате мы получим предварительную классификацию \tilde{K} .

Очевидно, что классификация \tilde{K} существенно зависит от порядка элементов в (2) и, следовательно, не удовлетворяет условию (3) определения 2. Для исключения зависимости классификации от порядка строк в исходной таблице T проводится II этап классификации представляющий собой следующую итерационную процедуру „перераспределения” элементов последовательности (2) между классами (13).

Для каждого элемента $s_p \in S$ определяется

$$(15) \quad \max \{R(s_p; \tilde{K}_1), R(s_p; \tilde{K}_2), \dots, R(s_p; \tilde{K}_{L'})\}.$$

Пусть этот максимум достигается на \tilde{K}_j . Если $R(s_p; K_j) \geq \delta_K$, то строку s_p относим к классу \tilde{K}_j . При этом, если строка s_p принадлежала классу $\tilde{K}_{j'}$, и $j' \neq j$, то последний может оказаться пустым и будет в таком случае элиминирован. Эта процедура повторяется для всех элементов последовательности (2) до тех пор, пока не будет получена (точная) классификация

$$(16) \quad K = \{K_1, K_2, \dots, K_L\}, \quad L \leq L',$$

обладающая тем свойством, что

$$(17) \quad \forall s_p \in S ((s_p \in K_j) \rightarrow \max_{K_i \in K} \{R(s_p; K_i)\} = R(s_p; K_j)).$$

Поскольку в данной работе излагается лишь прикладной аспект разработанного метода классификации, заметим, что описанная процедура сходится очень быстро (для весьма различной по своему характеру исходной информации число итераций, как правило, не превышало 5).

§ 3. Динамическая классификация

Задача поиска оптимальной классификации

До сих пор рассматривалась классификация при фиксированном векторе $\bar{u}(\epsilon, \delta_s, \delta_K)$ пороговых оценок и было показано, что в этих условиях удается получить классификацию K , оптимальную в том смысле, что она не зависит от способа упорядочения классифицируемых объектов и является устойчивой с точки зрения принадлежности объектов определенным классом. Очевидно, что в общем случае

$$K = K(\bar{u})$$

и для различных \bar{u} мы будем получать различные классификации. Поэтому возникает задача поиска классификаций, оптимальных в более широком смысле.

Назовём допустимым разбиение K множества объектов S на подмножества

$$(18) \quad K = \{K_1, K_2, \dots, K_L\},$$

если (18) удовлетворяет условиям (1) и (2), определения 2. Множество всех допустимых разбиений обозначим через \mathcal{K} . Пусть, далее, U — пространство допустимых значений вектора \bar{u} пороговых оценок, определяемое неравенствами (6), (8б) и (10а).

Определим на \mathcal{K} некоторый функционал \mathcal{J} , сопоставляющий любой допустимой классификации $K \in \mathcal{K}$ некоторое число $\mathcal{J}(K)$ и рассмотрим следующую задачу.

Для заданных начального „состояния” $K(0) \in \mathcal{K}$ и начального управления $\bar{u}(0) \in U$ найти такое допустимое управление $\bar{u}(t)$, $t = 0, 1, \dots$, которое является решением задачи

$$(19) \quad \mathcal{J} = \min_{\bar{u} \in U} \mathcal{J}(K(\bar{u})).$$

Отметим две особенности (и вытекающие из них трудности) указанной задачи. Во-первых, \mathcal{K} — дискретное неупорядоченное множество и поэтому понятие траектории в \mathcal{K} не имеет „физического” смысла и, в лучшем случае, сводится к апостериорному упорядочению некоторых элементов из \mathcal{K} в виде последовательности

$$K(\bar{u}(0)), K(\bar{u}(1)), \dots, K(\bar{u}(i)), \dots$$

Во-вторых, очевидно, что функционал \mathcal{J} является ступенчатой функцией от $\bar{u}(t)$, что исключает возможность использования классических процедур минимизации. В связи с этим были разработан метод минимизации ступенчатых функций, в основу которого был положен метод нелинейного программирования, впервые описанный в [5] (и независимо разработанный в [2]).

На основе программного комплекса АСУМ МС [2], [3] была разработана и реализована на ЭВМ БЭСМ-6 подсистема для решения задачи (19). В частности, задача (19) была решена для следующих типов функционалов.

1. Во многих задачах априори известно количество L^* классов. В этом случае в качестве функционала задачи (19) берется выражение

$$(20) \quad \mathcal{J}(\bar{u}) = |L^*(K(\bar{u})) - L^*|,$$

где $L^*(K(\bar{u}))$ — „рассчитанное” количество классов для классификации $K(\bar{u})$. Минимизируя функционал (20), мы найдём классификацию K^{opt} , содержащую L^* классов. Очевидно, что при столь общем критерии трудно гарантировать единственность решения задачи (19), тем не менее, если варьировать значение начального управления $\bar{u}(0)$, можно построить ряд классификаций

$$K_1, K_2, \dots,$$

содержащих одинаковое число L^* классов, и использовать полученную информацию для принятия решения о выборе направления дальнейших исследований.

2. Из „физических” соображений очевидно, что в общем случае качество классификации тем лучше, чем меньше расстояния между объектами внутри классов и чем большие расстояния между самими классами. Учитывая связь

$$(21) \quad \sum_{K_i \in \mathcal{K}} \sum_{\substack{s_p \in K_i \\ s_q \in K_i \\ p \neq q}} R(p, q) + \sum_{\substack{K_i \in \mathcal{K} \\ K_j \in \mathcal{K} \\ i \neq j}} \sum_{\substack{s_p \in K_i \\ s_q \in K_j \\ p \neq q}} R(p, q) = \sum_{\substack{s_p \in S \\ s_q \in S \\ p \neq q}} R(p, q) = \text{const.}$$

для поиска классификации, оптимальной в указанном смысле достаточно использовать следующий функционал

$$(22) \quad \mathcal{J}(\bar{u}) = \sum_{K_i \in \mathcal{K}(\bar{u})} \sum_{\substack{s_p \in K_i \\ s_q \in K_i}} R(p, q).$$

Решая задачу (19) с функционалом (22), мы получаем классификацию K^{opt} , при которой сумма расстояний между объектами внутри классов является минимальной (при этом, в силу (21), суммарное расстояние между классами будет, очевидно, максимальным).

3. Очевидно, что наибольший вклад с точки зрения формирования некоторого класса дают те признаки объектов, которые являются в данном классе наиболее близкими по своим значениям.

Обозначим через w_{ij} — вес j -го признака в i -ом классе и положим, по определению,

$$(23) \quad w_{ij} = 1 - \sigma_{ij}/\Delta_j,$$

где σ_{ij} — средне-квадратическое отклонение значения j -го признака в i -ом классе, а Δ_j — величина абсолютного колебания значений j -го признака (значений столбца s_j^+) в таблице T .

Рассмотрим произвольную классификацию K . Пусть класс $K_i \in K$. Упорядочим все N признаков по убыванию весов (23) в классе K_i

$$(24) \quad w_{i,j_1} \geq w_{i,j_2} \geq \dots \geq w_{i,j_d},$$

Введём в качестве порога целое число $\delta_D < N$ и следующее

Определение 3. Назовём понятием о классе K_i по порогу δ_D совокупность признаков с номерами

$$(25) \quad j_1, j_2, \dots, j_{\delta_D}.$$

Другими словами, понятием о классе является совокупность δ_D признаков с наибольшими весами.

Мы построили классификацию K и, в частности, получили класс $K_i \in K$, используя все N признаков объектов. Если теперь для объектов класса K_i мы будем решать задачу распознавания (15) не по всем N признакам (что, в силу свойства (17) классификации K было бы тавтологией), а по понятию о классе K_i , т.е. только по совокупности признаков (25), то возможно, что некоторые объекты класса K_i будут отнесены к другим классам. Обозначим

количество таких „переотнесенных” объектов, точнее, их долю от общего числа объектов класса K_i — через ν_i и положим

$$(26) \quad \nu(K(\bar{u})) = \max_{K_i \in K(\bar{u})} \nu_i.$$

Решая задачу (19) с функционалом (26) мы получаем классификацию K^{opt} , оптимальную с точки зрения информативности совокупности понятий о классах этой классификации K^{opt} . Это позволяет, в общем случае, заменить в каждом классе описание объектов с помощью N признаков их описанием через понятие о классе, т.е. как правило, существенно меньшим числом признаков, что в конечном счете упрощает задачу распознавания новых объектов с помощью полученной классификации K^{opt} .

Отметим важность изучаемой задачи разбиения для проблемы прогнозирования, содержанием которой является установление связи между прогнозируемым параметром или, иначе говоря, состоянием объектов и переменными $\{x_i, i = 1, 2, \dots, n\}$, определяющими эти состояния. При этом возникает необходимость в решении задачи оптимальной классификации множества реализаций контрольных выборок [4].

Литература

- [1] Ю. И. Журавлев, В. В. Никифоров, Кибернетика 3 (1971), Киев.
- [2] А. А. Зенкин, Канд. дисс., МГУ, Москва 1974.
- [3] —, Доклады АН СССР 230 (1976), 1051.
- [4] —, О математических методах прогнозирования, Знание, Москва 1976.
- [5] Р. Хуук, Т. А. Джинс, R. Hoole, T. A. Jeeves, J. Assoc. Comput. Mach. 8 (1961), 212.

*Presented to the Semester
Discrete Mathematics
(February 15–June 16, 1977)*



DISCRETE MATHEMATICS
BANACH CENTER PUBLICATIONS, VOLUME 7
PWN—POLISH SCIENTIFIC PUBLISHERS
WARSAW 1982

О ТОЧНОСТИ МЕТОДА „УСРЕДНЕНИЯ”

В. К. ЛЕОНТЬЕВ

Вычислительный центр АН СССР, Москва, СССР

Многие задачи теории кодирования, распознавания образов, экспертных оценок можно рассматривать как оптимизационные метрические задачи на единичном n -мерном кубе E^n . При этом стандартной метрикой обычно является метрика Евклида или метрика Хэминга, а оптимизируемый функционал является линейным от некоторой функции этой метрики. Задача состоит в нахождении множества фиксированной мощности на котором этот функционал достигает своего максимума или минимума.

Хорошо известным методом нахождения верхних (или нижних) оценок функционалов такого рода является „метод случайного кодирования” или метод оценки экстремума функционала средним значением этого функционала, вычисленным по „подходящему” ансамблю подмножества фиксированной мощности из E^n . Этот метод играет существенную роль не только в таком рода оптимизационных задачах, но и во многих других проблемах комбинаторного анализа.

Настоящая работа посвящена изучению точности метода „усреднения” для специального класса метрических оптимизационных задач, о котором шла речь выше. Основной результат работы состоит в следующем.

Введен класс функций R_0 для функционалов от которого „метод усреднения” даёт асимптотически точный результат. Показано, что класс R_0 является выпуклым конусом в пространстве всех вещественнонезначимых функций, определенных на $[0, +\infty]$. Описаны некоторые свойства этого конуса, а для подкласса выпуклых вниз и убывающих на $[0, +\infty]$ функций дано абстрактное описание подконуса R_0 в виде некоторого аналитического соотношения. Показано, что для функций из класса R_0 градиентный алгоритм, применённый к функционалам от этих функций приводит к асимптотически точному результату.

Пусть E^n — множество вершин единичного n -мерного куба или множество двоичных наборов длины n . Всюду в дальнейшем множество E^n мы будем рассматривать как n -мерное векторное пространство над полем из