

УДК 681.3.06:621.391

ЧАСТЬ I. ПРОГРАММЫ ТАКСОНОМИИ

В настоящей главе опубликован ряд программ для решения задач таксономии (группировки объектов, выделения групп близких точек). Под задачей таксономии обычно понимают задачу поиска наилучшего в некотором смысле разбиения S множества Q объектов (реализаций) q_i на таксоны (формальные элементы) S_m в заданном пространстве признаков X с помощью решающих функций \mathcal{D} определенного типа. При решении задачи таксономии, как правило, предполагается, что пространство признаков X и способ измерения расстояния ("близости") между точками выбраны (на основании знаний, опыта и интуиции исследователя) так, что точки, соответствующие "ближним", "похожим" объектам — объектам одного класса — образуют изолированные области, достаточно удаленные от образов объектов других классов.

Данные алгоритмы не используют сведения об априорном распределении точек по таксонам. Все построения для данного множества точек делаются только на основании предположения об изолированности областей, соответствующих различным образом.

По алгоритму "Фораль-І" ("Формальный элемент - І") объединяются в один таксон точки, близкие в евклидовом пространстве признаков. Границная область отдельного таксона представляет собой гиперсферу. В данном сборнике публикуются программы "Фораль-І" для ЭВМ БЭСМ-6, М-20 (М-220) и Минск-22. Описание алгоритма и его применений можно найти в работах [1, 2, 3]. В программах предусмотрена возможность нормализации исходных данных по дисперсиям и, как правило, программы должны использоваться в режиме "с нормализацией".

Эти программы могут быть использованы, в основном, для ре-

мения задач с количественными или ранжированными качественными признаками. Применять эти программы для чисто двоичных кодов нецелесообразно ввиду нерационального использования машинного времени и памяти. Для двоичных кодов разработаны специальные алгоритмы (группа "Форэль-5") и программы. Логика алгоритмов группы "Форэль-5" та же, что и алгоритмов "Форэль-1".

Программа "Форэль-5 а" [1, 3] позволяет осуществлять разбиение множества на таксоны при условии равнозначности (или почти равнозначности) разрядов, как по полному набору признаков, так и по любому их подмножеству без искажения массива исходных данных. С этой целью в информации к программе задается маска, размерность которой совпадает с размерностью исходного двоичного кода. В маске на местах тех признаков, по которым должна быть проведена классификация, стоят единицы. Если нужна классификация по всему набору признаков, то во всех разрядах должны стоять единицы. Так как при счете исходный массив данных не искается, то на одном и том же массиве данных можно без дополнительного ввода просчитать несколько вариантов задачи.

В алгоритме "Форэль-5 б", помимо задания признакового подпространства с помощью маски, предусмотрена возможность введения "весов" для признаков и отдельных разрядов. Это бывает необходимо в том случае, когда признаки неравнозначны или закодированы существенно различным количеством двоичных разрядов. Для возможности учета весовых коэффициентов в алгоритме "Форэль-5б" производится сравнение двоичных кодов Y_k ($\xi_{k1}, \dots, \xi_{kn}$) и Y_e ($\xi_{e1}, \dots, \xi_{en}$) по "звешенному" хэммингову расстоянию:

$$d(Y_k, Y_e) = \sum_{m=1}^n (\xi_{km} + \xi_{em}) \bmod 2 \cdot W_m,$$

где W_m - "вес", приписываемый каждому разряду двоичного кода. Разряды, относящиеся к одному и тому же признаку x_i , имеют равные "веса" W_i . Эти веса W_i должны быть заданы перед началом работы программы. В общем случае каждому двоичному разряду ξ_{ip} ($p = 1, \dots, d_i$) приписывается "вес" $W_i = \beta_i/d_i$, где β_i - "вес" признака x_i ; d_i - количество двоичных разрядов, отведенных под признак x_i . В случае, когда нет оснований выделять какие-либо признаки, т.е. $\beta_i = 1$, тогда $W_i = 1/d_i$. Если и все разряды равнозначны, то матрица весов состоит из единиц. В этом случае обе программы да-

ют один и тот же результат, но "Форэль-5 а" требует меньше машинного времени.

Как уже говорилось ранее, алгоритмы "Форэль-5 а" и "Форэль-5 б" предназначены для работы с данными, представленными двоичными кодами. Рассмотрим, какого типа информацию удобно представлять в таком виде. Прежде всего следует упомянуть о кодировании информации анкетного типа, например, такой как ответы испытуемого на те или иные вопросы, когда он должен сказать "да" или "нет". Признаки такого рода кодируются: 1 - "да", 0 - "нет". Двоичным кодом могут быть закодированы и количественные признаки, такие как вес, рост, размер и т.п. Представить число в виде двоичного кода можно различными способами. При подготовке данных для решения задач таксономии с помощью программ группы "Форэль-5" удобно пользоваться следующим методом. Диапазон изменения ($x_{imax} - x_{imin}$) конкретного числового признака x_i разбивается на некоторое число $d_i + 1$ градаций и соответственно под этот признак отводится d_i двоичных разрядов. Кодировать следует так, чтобы x_{imax} был представлен в виде набора из d_i единиц, а x_{imin} - из d_i нулей. Каждое значение признака x_i ($x_{imin} \leq x_i \leq x_{imax}$) кодируется соответствующим ему набором из d_i подряд идущих единиц и $d_i - d_i$ нулей. Например, для некоторого множества признак x_i принимает следующие значения: 5, 7, 17, 10, 8. Здесь $x_{imax} = 17$, $x_{imin} = 5$, $x_{imax} - x_{imin} = 12$. Пусть под этот признак отведено 5 разрядов. Тогда в закодированном виде приведенная выше последовательность значений признака имеет вид: 5 ~ 00000; 7 ~ 10000; 17 ~ 11111; 10 ~ 11100; 8 ~ 11000.

Выделение таксонов сложной формы осуществляется алгоритмом "Краб" [4, 5]. Этим алгоритмом производится поиск лучшего по критерию качества F (см. [4, 5]) разбиения множества на заданное число таксонов.

Функция F не зависит от числа точек L и количества таксонов K , что дает возможность не только оценивать различные разбиения одного и того же множества Q на равное число таксонов, но и определять наиболее предпочтительное для данного множества количество K таксонов S_m , а также сравнивать качество разбиений различных множеств Q_j на K_j таксонов.

Предварительно перед работой программы "Краб" все точки

исходного множества должны быть соединены в связный неориентированный граф без петель с минимальной суммарной длиной ребер [6, 7] (рис. I) по программе "Построение кратчайшего пути" (данний сборник).

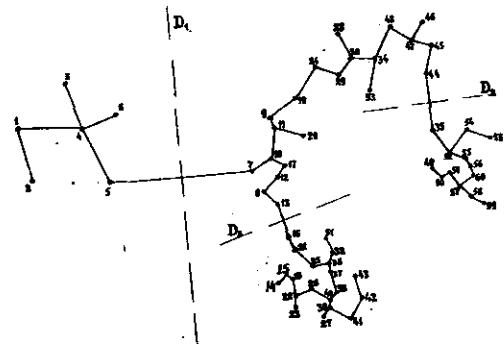


Рис. I. Соединение точек кратчайшим
незамкнутым путем.

Результаты работы программы ПКП являются исходными данными для программы "Краб".

- 6. Boruvka O., On a minimal problem. Prace moraske, Pridovedecky Spolecnosti, N 3, 1926.
- 7. Kruskal J.B., On a shortest spanning subtree of graph and the travelling salesman problem. Proc. Amer. Math.Soc., N 7, 1956.

В.Н. Ёлкина

Л и т е р а т у р а

1. Ёлкина В.Н., Загоруйко Н.Г. Об алфавите объектов распознавания. Сб. Вычислительные системы, вып. 22, Новосибирск, Изд-во "Наука", 1966.
2. Ёлкина В.Н., Ёлкин Е.А., Загоруйко Н.Г. О возможности применения методов распознавания образов в палеонтологии. Геология и геофизика, № 9, 1967.
3. Загоруйко Н.Г., Заславская Т.И., Ёлкина В.Н., Лбов Г.С. и др. Распознавание образов в социальных исследованиях. Новосибирск, Изд-во "Наука", 1968.
4. Ёлкина В.Н., Загоруйко Н.Г. Количественные критерии качества таксономии и их использование в процессе принятия решений. Сб. Вычислительные системы, вып. 36, Изд-во "Наука", Новосибирск, 1969.
5. Ёлкина В.Н. Выбор формальных элементов алгоритма (алгоритм таксономии). Автореферат диссертации. Новосибирск, 1969.