

УДК 519.95.

КИБЕРНЕТИКА И ТЕОРИЯ РЕГУЛИРОВАНИЯ

Ю. Л. ВАСИЛЬЕВ, А. Н. ДМИТРИЕВ

**СПЕКТРАЛЬНЫЙ ПОДХОД К СРАВНЕНИЮ ОБЪЕКТОВ,
ОХАРАКТЕРИЗОВАННЫХ НАБОРОМ ПРИЗНАКОВ**

(Представлено академиком С. Л. Соболевым 23 II 1972)

1. Здесь рассматривается вопрос о заключительном этапе процедур, называемых диагностическими, распознаванием образов и т. п., а именно: вопрос об уточнении диагностики на основе той или иной числовой меры. Предлагаемый подход вводит на объектах и характеризующих их признаках числовую меру и дает способ ее вычисления, основанный на итерациях. Следует отметить, что эта мера в некотором смысле отражает естественное соотношение между сравниваемыми объектами (п. 3) и оправдывается при практических применениях (п. 6). Способ вычисления отличается малой трудоемкостью и делает доступными задачи с весьма большим числом признаков (пп. 3, 4, 6). На обрабатываемые таблицы налагается некоторое естественное структурное ограничение (связность, п. 4).

2. Пусть имеется h объектов W_1, \dots, W_h , охарактеризованных двузначными признаками P_1, \dots, P_l . Пусть $T = (t_{ij})$ — таблица размера $h \times l$, которая составлена из единиц и нулей и в которой i -я строка (t_{i1}, \dots, t_{il}) отвечает объекту W_i , $i = 1, \dots, h$, j -й столбец (t_{1j}, \dots, t_{hj}) отвечает признаку P_j , $j = 1, \dots, l$; таблица T отражает выраженность признаков у объектов: если для i -го объекта и j -го признака она превышает некоторый уровень, то $t_{ij} = 1$, а если меньше, то $t_{ij} = 0$. Условимся, что $h \geq 2$, $l \geq 2$ и что в таблице T нет строк и нет столбцов, составленных сплошь из нулей.

Если признаки k -значны, $k > 2$, то t_{ij} может принимать соответствующее число значений из отрезка $[0, 1]$. Ограничимся рассмотрением случая $k = 2$, так как переход на случай $k > 2$ будет очевиден.

Таблица описывает объекты в связи с тем или иным фактором X , и задача состоит в оценке последнего по картине, даваемой таблицей T . Предполагается, что признаки P_1, \dots, P_l существенны для X , что набор их достаточно полон и что объекты W_1, \dots, W_h родственны по отношению к X . Пусть T — геологическая таблица, X — запасы месторождений, объекты — месторождения одного и того же полезного ископаемого, которые родственны по признакам и сравнимы по запасам, но оценка последних требует уточнения.

Нас интересуют случаи, когда число объектов сравнительно невелико (10—20), число признаков велико (десятки, сотни), влияния признаков находятся в сложном переплетении. Поскольку статистический подход встречает здесь серьезные препятствия, возникли другие подходы. Сначала они проверяются на всесторонне изученных примерах, затем распространяются на менее изученные ситуации вплоть до решения задач прогноза. Хорошо согласующиеся с практикой оценки дает основанный на аппарате тупиковых тестов⁽¹⁾ известный тестовый подход⁽²⁻⁴⁾, на который мы будем ориентироваться при сопоставлении результатов*.

* Напомним, что для таблицы T тестовый вес j -го столбца (j -го признака) $c_j = (N_j/N)$, где N_j — число тупиковых тестов, содержащих j -й столбец, $j = 1, \dots, l$, N — число всех тупиковых тестов; тестовый вес i -й строки (i -го объекта) $r_i = t_{i1}c_1 + \dots + t_{il}c_l$, $i = 1, \dots, h$.

3. Вектор $\mathbf{a} = (a_1, \dots, a_m)$ называется положительным, если $a_1, \dots, a_m > 0$. Норма вектора $\mathbf{a} = (a_1, \dots, a_m)$ — число $|\mathbf{a}|$, равное $\max_{1 \leq i \leq m} |a_i|$; вектор \mathbf{a} называется нормальным, если $|\mathbf{a}| = 1$.

Числовую меру для объектов и признаков естественно задавать в виде положительных нормированных векторов $\omega = (\omega_1, \dots, \omega_h)$ и $\pi = (\pi_1, \dots, \pi_l)$. Ниже мы определим их как нагрузку строк и нагрузку столбцов таблицы T . По ним предлагается судить о проявленности фактора X соответственно в объектах W_1, \dots, W_h , а также о степени влияния на нее признаков P_1, \dots, P_l .

Определение нагрузки подсказывается анализом грубой оценки объектов и признаков, при которой им приписывается вес, равный количеству единиц в соответствующих строках и столбцах таблицы T . При этом

	P_1	P_2	P_3	P_4	Вес
W_1	1	1	0	0	2
W_2	0	1	1	0	2
W_3	0	0	1	1	2
Вес	1	2	2	1	

Рис. 1. Грубая оценка

1	1	0	0	2	0,75	0,75	0,714
0	1	1	0	2	1	1	1
0	0	1	1	2	0,75	0,75	0,714
1	2	2	1				
0,5	1	1	0,5				
0,429	1	1	0,429				
0,429	1	1	0,429				

Рис. 2. Итерационный процесс («качели»)

строки и столбцы выступают как бы порознь. В примере на рис. 1 строки получают равные веса, однако они различны по отношению к столбцам: во второй строке обе единицы отвечают признакам, представленным единицами и в других строках, а в первой и третьей строках имеется лишь по одной такого рода единице. Содержательно это может означать, что вторая строка отвечает более сильному проявлению фактора X .

Определение нагрузки отражает итог описанных ниже пересчетов весов строк с учетом весов столбцов, и наоборот. При каждом пересчете вес строки W_i , $i = 1, \dots, h$, определяется через веса столбцов, найденные на предшествующем шаге как сумма весов тех столбцов, по которым в строке W_i стоят единицы; вес столбца P_j , $j = 1, \dots, l$, аналогично определяется через найденные на предшествующем шаге веса строк; затем полученные два набора чисел нормируются, начинается новый пересчет и т. д. Процесс сходится к некоторым предельным векторам ω и π , которые примем в качестве нагрузок. Для таблицы на рис. 1 $\omega \approx (0,707; 1; 0,707)$, $\pi \approx (0,414; 1; 1; 0,414)$. По сравнению с грубой оценкой нагрузки ω выделяет вторую строку. На рис. 2 представлены результаты трех первых пересчетов.

4. Исходя из таблицы T и пары векторов $\mathbf{w} = (w_1, \dots, w_h)$, $\mathbf{p} = (p_1, \dots, p_l)$, определим пару векторов $\mathbf{w}' = (w'_1, \dots, w'_h)$, $\mathbf{p}' = (p'_1, \dots, p'_l)$, а также нормирующие множители α , β и пару векторов \mathbf{w}'' и \mathbf{p}'' :

$$\mathbf{w}' = T(\mathbf{p}), \text{ т. е. } w'_i = t_{i1}p_1 + \dots + t_{il}p_l, \mathbf{w}'' = \alpha \cdot \mathbf{w}', \quad (1)$$

$$i = 1, \dots, h, \quad |\mathbf{w}''| = 1$$

$$\mathbf{p}' = T^*(\mathbf{w}), \text{ т. е. } p'_j = t_{1j}w_1 + \dots + t_{hj}w_h, \mathbf{p}'' = \beta \cdot \mathbf{p}', \quad (2)$$

$$j = 1, \dots, l, \quad |\mathbf{p}''| = 1$$

Формулы (1), (2) дают формальное определение упомянутых пересчетов.

Определения. Пару нормированных векторов $\varphi = (\varphi_1, \dots, \varphi_h)$, $\psi = (\psi_1, \dots, \psi_l)$ назовем фракционирующей парой таблицы, а сами векторы — фракционирующими векторами, если $\varphi' = \varphi$ и $\psi' = \psi$.

Фракционирующую пару векторов ω , π назовем нагрузкой таблицы T , а сами векторы — нагрузкой строк и нагрузкой столбцов, если эти векторы положительны.

Таблицу T , нагрузка которой существует и единственна, назовем измеримой.

Таблицу T назовем разрозненной, если среди всех ее h строк найдутся h_1 таких строк, каждая из которых ортогональна* каждой из прочих $h - h_1$ строк ($h_1 < h$). Иными словами, совокупность объектов распадается на две такие группы по h_1 и $h - h_1$ объектов, что каждый из признаков обязательно представлен нулями во всех объектах какой-то одной из этих двух групп (грубо говоря, эти две группы «не имеют ничего общего»). Перестановками строк и столбцов любую разрозненную таблицу

можно привести к виду $\begin{pmatrix} T_1 & 0 \\ 0 & T_2 \end{pmatrix}$. Любая разрозненная таблица неизмерима**.

Таблицу T , не являющуюся разрозненной, назовем связной.

Теорема 1. Любая связная таблица измерима.

Теорема 2 (о качелях). Нагрузка строк ω и нагрузка столбцов π измеримой таблицы T являются пределами соответственно последовательностей векторов

$$\omega^0, \omega^1, \dots, \omega^n, \dots, \quad \text{где } \omega^0 = \overbrace{(1, \dots, 1)}^{h \text{ единиц}}, \quad \omega^n = (\omega^{n-1})', \quad (3)$$

$$\pi^0, \pi^1, \dots, \pi^n, \dots, \quad \text{где } \pi^0 = \underbrace{(1, \dots, 1)}_{l \text{ единиц}}, \quad \pi^n = (\pi^{n-1})'. \quad (4)$$

5. Доказательство теоремы 1 сводит ее к известной теореме Фробениуса⁽⁵⁾ о спектре неотрицательной неразложимой матрицы Φ . Роль матрицы Φ играют матрицы TT^* и T^*T , объединяемые матрицей $K^2 = \begin{pmatrix} TT^* & 0 \\ 0 & T^*T \end{pmatrix}$, где $K = \begin{pmatrix} 0 & T \\ T^* & 0 \end{pmatrix}$.

Последняя связана с (3), (4) и служит посредником между Φ и T . Сокращения с.в. и с.з. означают собственный вектор и собственное значение. Для векторов $\mathbf{a} = (a_1, \dots, a_h)$ и $\mathbf{b} = (b_1, \dots, b_l)$ обозначим через $\mathbf{a} \vee \mathbf{b}$ вектор $(a_1, \dots, a_h, b_1, \dots, b_l)$. Аккомодацией назовем преобразование вектора $\mathbf{a} \vee \mathbf{b}$ в пару нормированных векторов $v_a \cdot \mathbf{a}$, $v_b \cdot \mathbf{b}$ ($|v_a \cdot \mathbf{a}| = |v_b \cdot \mathbf{b}| = 1$). Пусть X — множество нормированных с.в. матрицы K и пусть Y — множество фракционирующих пар таблицы T .

(А) Аккомодация взаимно однозначно отображает X на Y ; при этом положительные собственные векторы переходят в нагрузки, и обратно. С.з. матриц TT^* и T^*T совпадают (различиями в кратности с.з. 0 пренебрегаем).

(Б) Если $K(\mathbf{a} \vee \mathbf{b}) = \lambda(\mathbf{a} \vee \mathbf{b})$, то $\|\mathbf{a}\| = \|\mathbf{b}\|$ (норма эвклидова) и $TT^*(\mathbf{a}) = \lambda^2 \mathbf{a}$, $T^*T(\mathbf{b}) = \lambda^2 \mathbf{b}$; имеет место и обратное.

Пусть Φ — любая из матриц TT^* и T^*T . Матрица Φ неотрицательна; так как таблица T связная, то Φ неразложима. Согласно теореме Фробениуса, матрица Φ имеет такое с.з. λ_{\max} и отвечающий ему с.в. \mathbf{u} , что а) $\lambda_{\max} > 0$, λ_{\max} больше всех других с.з., λ_{\max} является простым с.з.; б) с.в. \mathbf{u} положителен. Из (б), (Б) и (А) следует существование нагрузки таблицы T ; из простоты с.з. λ_{\max} , (Б), ортогональности с.в. матрицы K и из (А) следует единственность нагрузки.

* Два вектора (x_1, \dots, x_l) и (y_1, \dots, y_l) ортогональны, если $x_1 y_1 + \dots + x_l y_l = 0$.

** Можно доказать, что доля разрозненных таблиц среди всех $(h \times l)$ - таблиц T стремится к нулю при $h \leq l \leq 2^{h-1}$ и $h \rightarrow \infty$.

Доказательство теоремы 2 сводит ее к обычному итерационному построению ⁽⁶⁾ с.в. u матрицы Φ , отвечающего с.з. λ_{\max} .

6. Примеры решения прогнозно-поисковых задач из геологии.

Пример 1. При сравнительном изучении гигантских месторождений нефти ⁽⁷⁾ требуется из общей совокупности (около 350) характеристических признаков выделить группу признаков, имеющих поисковое значение. Для этого каждая из выделенных групп испытывается на спо-



Рис. 3. Группа месторождений Аравийской платформы

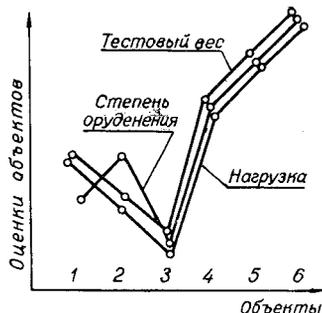


Рис. 4. Дифференцированные трапповые интрузии

собность фракционировать месторождения по масштабу запасов. Одна из таких групп, именуемая «структурная ловушка» ($l=28$), нами взята для сопоставления фракционирующих возможностей тестового и спектрального подходов (рис. 3).

Пример 2. При изучении дифференцированных трапповых интрузий, перспективных на медно-никелевое оруденение, рассматривались интрузии норильского типа ⁽⁸⁾. Упорядоченность последних по степени оруденения, а также фракционирование их по результатам обработки таблиц, содержащих данные о петрохимических признаках ($l=32$), тестовым и спектральным подходами представлены на рис. 4.

Сравнительное изучение тестового и спектрального подходов было проведено для большого набора геологических таблиц. Как и в приведенных примерах, имеет место близость оценок объектов по обоим подходам. Время обработки таблиц: для одних и тех же таблиц требовалось 30—90 мин. на БЭСМ-6 при тестовом подходе и не более 1 мин. на М-20, М-220 при спектральном подходе.

Наряду с близостью нагрузки и тестового веса для объектов наблюдаются существенные различия между нагрузкой и тестовым весом для признаков.

Институт математики

Поступило

Институт геологии и геофизики
Сибирского отделения Академии наук СССР
Новосибирск

17 II 1972

ЦИТИРОВАННАЯ ЛИТЕРАТУРА

- ¹ И. А. Чегнис, С. В. Яблонский, Тр. Матем. инст. им. В. А. Стеклова АН СССР, 51, 270 (1958). ² А. Н. Дмитриев, Ю. И. Журавлев, Ф. П. Кренделев, Сборн. Дискретный анализ, в. 7, 3 (1966). ³ А. Н. Дмитриев, Ю. И. Журавлев, Ф. П. Кренделев, Геология и геофизика, № 5, 50 (1968). ⁴ С. В. Яблонский, Н. Г. Демидова и др., Геология рудных месторожд., № 2, 3 (1971). ⁵ Ф. Р. Гантмахер, Теория матриц, М., 1967. ⁶ Д. К. Фаддеев, В. Н. Фаддеева, Вычислительные методы линейной алгебры, М., 1960. ⁷ В. С. Вышемирский, А. Н. Дмитриев, А. А. Трофимук, Мировой нефтяной конгресс, М., 1971. ⁸ А. Н. Дмитриев, В. В. Золотухин, Ю. Р. Васильев, Сов. геол., № 12, 64 (1968).