

Ю. А. ВОРОНИН

**ВВЕДЕНИЕ МЕР СХОДСТВА И СВЯЗИ ДЛЯ РЕШЕНИЯ
ГЕОЛОГО-ГЕОФИЗИЧЕСКИХ ЗАДАЧ**

(Преображенено академиком Г. И. Марчуком 8 II 1971)

При решении задач выделения, описания, группирования и распознавания геолого-геофизических объектов прямое использование теоретико-вероятностной меры встречается с рядом трудностей, обусловленных малыми объемами выборок, разнородностью и многомерностью описания объектов. Большие трудности возникают и при интерпретации результатов, полученных таким путем. По этим причинам многие предпочитают использовать при решении упомянутых задач другие меры, которые обычно называют мерами сходства и связи. Во многих случаях использование таких мер, традиционных для геологии, как представляется, дает очевидные преимущества. Однако эти меры сходства и связи, в отличие от теоретико-вероятностной меры, не имеют аксиоматического обоснования: в одних и тех же конкретных ситуациях разными авторами они вводятся различно, часто недостаточно четко. Представляется желательным разработать аксиоматический подход к построению таких мер.

Пусть задана таблица — объекты (a_1, a_2, \dots, a_n) , свойства $(f_a^1, f_a^2, \dots, f_a^m)$:

$$\{f_a^i\}, \quad f_a^i \in (f_a), \quad i = 1, 2, \dots, n, \quad a = 0, 1, \dots, m, \quad (1)_1$$

где через f_a^i обозначено значение свойства f_a на объекте a_i , через (f_a) — множество всех возможных значений свойства f_a . Будем говорить, что свойство f_a является: арифметическим, если на (f_a) определены арифметические операции; логическим 1-го рода, если на (f_a) определены только отношения порядка и эквивалентности; логическим 2-го рода, если на (f_a) определено только отношение эквивалентности.

Принято под мерой сходства между объектами a_i и a_j по свойствам $F = (f_1, f_2, \dots, f_m)$ понимать функцию $\Lambda(F^i, F^j)$, отвечающую аксиомам (1): 1) $0 \leq \Lambda(F^i, F^j) \leq 1$; 2) $\Lambda(F^i, F^j) = \Lambda(F^j, F^i)$; 3) $F^i = F^j \Rightarrow \Lambda(F^i, F^j) = 1$. Эта система аксиом неполна, она не позволяет однозначно определить вид $\Lambda(F^i, F^j)$. Дополним ее так, чтобы можно было однозначно определить вид $\Lambda(F^i, F^j)$ и указать способ определения ее параметров с учетом конкретной ситуации.

Пусть f_a — арифметическое свойство. Потребуем, чтобы мера сходства между объектами a_i и a_j по свойству f_a удовлетворяла аксиомам: 1) $0 \leq \Lambda(f_a^i, f_a^j) \leq 1$; 2) $\Lambda(f_a^i, f_a^j) = \Lambda(f_a^j, f_a^i)$; 3) $f_a^i = f_a^j \Rightarrow \Lambda(f_a^i, f_a^j) = 1$; 4) $\Lambda(f_a^i, f_a^j) = 1 \Rightarrow f_a^i = f_a^j$; 5) $\Lambda(f_a^i, f_a^j)$ имела производную по f_a^j , линейно зависящую от f_a^j .

Тогда получим

$$\Lambda(f_a^i, f_a^j) = [1 - ((f_a^i - f_a^j)/\Delta f_a)^2], \quad (2)$$

$$\Delta f_a = f_a^{**} - f_a^*, \quad f_a^{**} = \max f_a, \quad f_a^* = \min f_a.$$

Если f_a — свойство логическое 1-го рода, то $\Lambda(f_a^i, f_a^j)$ будем определять тоже формулой (2), если же f_a — свойство логическое 2-го рода, $\Lambda(f_a^i,$

f_a^i) определим так:

$$\Lambda(f_a^i, f_a^j) = \begin{cases} 1, & f_a^i = f_a^j, \\ 0, & f_a^i \neq f_a^j. \end{cases} \quad (2)_2$$

Если принять еще одну аксиому: 6) мера сходства $\Lambda(F^i, F^j)$ должна определяться формулой

$$\Lambda(F^i, F^j) = \sum_{a=1}^m \delta_a \Lambda(f_a^i, f_a^j), \quad \sum_{a=1}^m \delta_a = 1, \quad (3)$$

то получим ряд интересных возможностей. Параметры δ_a можно определять стандартно из условия

$$\sum_{ij} [\Lambda(f_a^i, f_a^j) - \Lambda(F^i, F^j)]^2 = \min, \quad (4)$$

$$i, j = 1, 2, \dots, n, \quad i < j,$$

где f_0 — распознаваемое свойство.

Таблицу (1) можно описать так, чтобы ее описание не зависело от нумерации строк и столбцов, линейного преобразования столбцов $f_a' = e_a^0 + e_a^1 f_a$, числа строк и столбцов. Такое описание можно построить, например, следующим образом. Приведем таблицу (1) в соответствие таблицу:

$$\{\Lambda(F^i, F^j)\}, \quad i, j = 1, 2, \dots, n. \quad (5)$$

Найдем такие a_k , для которых обращается в максимум выражение

$$T(i) = \Lambda(i) \left\{ 1 - \frac{1}{n} \sum_{j=1}^n [\Lambda(i) - \Lambda(F^i, F^j)]^2 \right\}, \quad (6)$$

$$\Lambda(i) = \frac{1}{n} \sum_{j=1}^n \Lambda(F^i, F^j).$$

Возьмем любой a_k и назовем его голотипом. Зафиксируем μ — константу для разбиения объектов a_1, a_2, \dots, a_n на компоненты связности, $0 \leq \mu \leq 1$ *. Пусть $A_1, A_2, \dots, A_{n(\mu)}$ — компоненты связности, а $a_{k1}, a_{k2}, \dots, a_{kn(\mu)}$ — их голотипы. Определим

$$\Lambda_s(\mu) = \min \Lambda(F^{ks}, F^j), \quad a_j \in A_s, \quad s = 1, 2, \dots, n(\mu); \quad (7)$$

$$\tilde{\Lambda}_s(\mu) = \Lambda(F^{ks}, F^s), \quad s = 1, 2, \dots, n(\mu), \quad (8)$$

и, кроме того, найдем

$$\Lambda(\mu) = \frac{1}{n(\mu)} \sum_{s=1}^{n(\mu)} \Lambda_s(\mu), \quad (9)$$

$$\tilde{\Lambda}(\mu) = \frac{1}{n(\mu)} \sum_{s=1}^{n(\mu)} \tilde{\Lambda}_s(\mu). \quad (10)$$

Тогда функции

$$\theta(\mu) = (n(\mu) - 1) / (n - 1), \quad (11)_x$$

$$\varphi(\mu) = (\Lambda(\mu) - \Lambda(0)) / (1 - \Lambda(0)), \quad (11)_z$$

$$\psi(\mu) = \tilde{\Lambda}(\mu) \quad (11)_z$$

* Совокупность объектов $A = (a_1, a_2, \dots, a_n)$ является компонентой связности, если для любых a_i и $a_j \in A$ можно найти такие $a_{kp} \in A$, что $\Lambda(F^i, F^k) \geq \mu$, $\Lambda(F^k, F^l) \geq \mu, \dots, \Lambda(F^{ka}, F^j) \geq \mu$.

будут давать нужное описание таблицы (1)₁. Описание (11) является более грубым, чем описание через закон распределения. Закон распределения имеет смысл только тогда, когда все свойства f_a являются арифметическими, и он зависит от числа столбцов таблицы (1)₁. Функции (11) могут быть описаны через некоторые параметры $p_1(z), p_2(z), \dots, p_i(z)$, $z = 0, \varphi, \psi$. Это позволяет использовать формулу (3) для определения мер сходства между таблицами (1)₁, произвольными множествами объектов $a_{i_1}, a_{i_2}, \dots, a_{i_h}$.

Обратимся к мерам связи между свойствами f_a и f_b . Потребуем, чтобы эта мера $\sigma(f_a, f_b)$ являлась функцией от мер (2) и удовлетворяла аксиомам: 1) $0 \leq \sigma(f_a, f_b) \leq 1$; 2) $\sigma(f_a, f_b) = \sigma(f_b, f_a)$; 3) $f_a = e_b^0 + e_b^1 f_b \Rightarrow \sigma(f_a, f_b) = 1$; 4) $\sigma(f_a, f_b) = 1 \Rightarrow f_a = e_b^0 + e_b^1 f_b$.

Тогда получим

$$\sigma(f_a, f_b) = \left(\sum_{ij} \Lambda(f_a^i, f_a^j) \Lambda(f_b^i, f_b^j) \right) / \left(\sqrt{\sum_{ij} \Lambda^2(f_a^i, f_a^j)} \sqrt{\sum_{ij} \Lambda^2(f_b^i, f_b^j)} \right), \quad (12)$$

$$i, j = 1, 2, \dots, n, \quad i < j.$$

Если положить

$$\Lambda(f_\gamma^i, f_\gamma^j) = \begin{cases} f_\gamma^i - \bar{f}_\gamma, & f_\gamma^i = f_\gamma^j, \\ 0, & f_\gamma^i \neq f_\gamma^j, \end{cases} \quad \gamma = a, b, \quad (2)$$

то (12) перейдет в выборочный коэффициент корреляции. Аналогично можно получить коэффициенты Спирмэна, Кендалла, Чупрова и др. (2).

Использование (12) тоже открывает ряд интересных возможностей. Можно найти меру связи между двумя любыми совокупностями свойств $(f_{a1}, f_{a2}, \dots, f_{ak})$ и $(f_{b1}, f_{b2}, \dots, f_{bl})$, фиксируя распознаваемое свойство f_0 . Если f_0 — распознаваемое свойство, «идеально хорошее», то легко построить таблицу

$$\{\Lambda(\tilde{f}_0^i, \tilde{f}_0^j)\}, \quad i, j = 1, 2, \dots, n, \quad i < j, \quad (13)$$

для «идеально плохого» свойства \tilde{f}_0 из условия $\sigma(f_0, \tilde{f}_0) = 0$. Тогда каждому свойству f_a можно приписать коэффициент информативности

$$\kappa(f_a / f_0) = \sigma(f_a, f_0) - \sigma(f_a, \tilde{f}_0). \quad (14)$$

Аналогично можно получить $\kappa(f_{a1}, f_{a2}, \dots, f_{ak} / f_0)$.

Из предыдущего следует, что предлагаемые меры позволяют ставить и решать, минуя введение теоретико-вероятностной меры, целый круг задач, важных для обработки любых данных, например, таких, как задача об оптимальном разбиении промежутка изменения свойства на интервалы, задача выбора оптимальной совокупности свойств, задача представительности выборки. Если (f_a^{**}, f_a^*) разбит на интервалы $f_{a1}, f_{a2}, \dots, f_{al}$, то меру (2) можно записать так:

$$\Lambda(f_a^i, f_a^j) = [1 - ((p - q)/(l - 1))^2], \quad (2)$$

$$f_a^i \in f_{a_p}, \quad f_a^j \in f_{a_q}, \quad p, q = 1, 2, \dots, l.$$

Можно построить алгоритм, позволяющий найти такое разбиение (f_a^{**}, f_a^*) на $f_{a1}, f_{a2}, \dots, f_{al}$, которое при минимальном l обращает в максимум (14). Используя (12), можно получить таблицу

$$\{\sigma(f_a, f_b)\}, \quad a, b = 1, 2, \dots, m', \quad a < b. \quad (15)$$

Опираясь на (14) и (15), можно построить алгоритм, позволяющий при минимальном k найти такую $f_{a1}, f_{a2}, \dots, f_{ak}$, которой отвечает максимальное значение $\kappa(f_{a1}, f_{a2}, \dots, f_{ak} / f_0)$. Таблица

$$\{f_a^i\}, \quad i = 1, 2, \dots, n', \quad n' \leq n, \quad a = 1, 2, \dots, m, \quad m' \leq m, \quad (1_2)$$

будет отвечать «представительной» выборке из таблицы (1)₁, если функции (11) для таблиц (1)₁ и (1)₂ оказываются, в некотором смысле, одинаковыми.

Имеется возможность получить ряд полезных теорем, учитывая, что меры (2) могут быть обобщены на случай, когда множество объектов несчетно.

Таблицу (5) можно интерпретировать как граф со «взвешенными» ребрами. Однако таблицы (1), по-видимому, можно рассматривать как новые математические объекты, поскольку для них иначе формулируются проблемы изоморфизма, изоморфного вхождения и пр. ('). По-видимому, предлагаемые меры следует рассматривать как основу для грубого, детерминированного, решения задач, упомянутых в самом начале. Опираясь на такие решения, в некоторых случаях можно будет ставить вопрос о более тонком, теоретико-вероятностном, решении этих задач.

Вычислительный центр
Сибирского отделения Академии наук СССР
Новосибирск

Поступило
25 I 1971

ЦИТИРОВАННАЯ ЛИТЕРАТУРА

- ¹ В. И. Васильев, Распознающие системы (справочник), Киев, 1969. ² В. Ю. Урбах, Математическая статистика для медиков и биологов, Изд. АН СССР, М., 1963. ³ А. А. Зыков, Теория конечных графов, Новосибирск, 1969.