

## СОГЛАСОВАНИЕ РАЗНОТИПНЫХ ШКАЛ

Н.Г.Загоруйко

### § I. Введение

Рассмотрению проблемы оценки "близости", "похожести", "расстояния" между объектами, свойства которых измерены в разных шкалах, посвящены многие работы [1,2]. При этом ищутся такие меры, которые удовлетворяли бы обычным аксиомам (непрерывности, симметричности и т.п.), были бы инвариантны к допустимым преобразованиям в данной шкале и не зависели бы от состава изучаемых объектов. Итоги этих рассмотрений сводятся в основном к тому, что меры, инвариантные к допустимым преобразованиям для многих шкал, можно указать, а мер, которые не зависели бы от состава выборки, не существует. Добавление к конечной выборке  $A$  или изъятие из нее какого-нибудь объекта может изменить прежние порядковые номера объектов  $i$  и  $j$  (для шкал порядка) или нормировку (для более сильных шкал), что приводит к изменению  $d_{ij}$ , расстояния между  $i$ -м и  $j$ -м объектами (см. [1]).

Какой же вывод нужно сделать из этих результатов? Не следует ли признать, что адекватных мер близости между объектами любой конечной выборки нет, а следовательно, нет и оснований верить результатам решения всех тех задач, в которых существенно используется меры близости или расстояния между объектами, т.е. задач таксономии, расположения образов, корреляционного, регрессионного анализа и т.п.?

Не будем спешить соглашаться с таким пессимистическим заключением. Вспомним, что  $d_{ij}$  меняется не при всяком изменении состава выборки  $A$ . Действительно, при нормировке сильных шкал по разности между самым большим и самым малым значением характеристи-

ки  $x$  в таблице, т.е. по  $(x_{\max} - x_{\min})$ , мера  $d_{ij}$  будет сохраняться всегда, пока изменения состава объектов не коснутся  $x_{\max}$  или  $x_{\min}$ . Для любых шкал нормированная мера  $d_{ij}$  останется неизменной, если в таблице продублировать все объекты любое число раз.

Очевидно, что справедливо и обратное утверждение: если конечная выборка  $A$  составлена из объектов, которые правильно отражают закономерности генеральной совокупности  $G$ , то меры  $d_{ij}$  будут одинаковыми для объектов  $i$  и  $j$ , независимо от того, рассматриваем ли мы их на фоне выборки  $A$  или на фоне генеральной совокупности  $G$ . Выводы (т.е. таксономия, решающие правила, региональные уравнения и т.д.), сделанные на основании мер  $d_{ij}$ , найденных по выборке, будут сохраняться и для генеральной совокупности  $G$ .

Отсюда мы приходим к существенному заключению: если выборка  $A$  представительна, то можно указать меры  $d_{ij}$ , которые будут адекватными и по отношению к соответствующим шкалам и по отношению к генеральной совокупности. А если выборка  $A$  непредставительна, то никакие гарантии инвариантности  $d_{ij}$  к допустимым преобразованиям шкал не имеют смысла – из-за непредставительности  $A$  индуктивные выводы для  $G$  все равно будут ложными.

Данная работа не ставит своей целью рассмотрение проблемы представительности выборки. Будем исходить из того, что обоснование гипотезы представительности выборки  $A$  каким-то путем сделано и нас теперь интересует только вопрос о выборе меры расстояния  $d_{ij}$ .

Нам хотелось бы иметь меры  $d_{ij}$ , обладающие следующими очевидными свойствами [1,2]:

а) непрерывности: мера  $d(x_i, x_j)$  должна быть непрерывной функцией своих аргументов;

б) симметричности: предполагая пространство значений аргументов изотропным [3], потребуем, чтобы выполнялось соотношение  $d(x_i, x_j) = d(x_j, x_i)$ ;

в) нормированности:  $d(x_i, x_j)$  должно изменяться в пределах от 0 до 1, причем  $d(x_i, x_j) = 0$ , если  $x_i = x_j$ ;

г) инвариантности [4]: для преобразования  $\varphi$ , допустимого для данной шкалы,  $d(x_i, x_j) = d[\varphi(x_i), \varphi(x_j)]$ ;

д) свойством треугольника: для любых трех объектов  $a, b, c$  мера  $d_{ab} \leq (d_{ac} + d_{cb})$ .

Не для всех задач анализа данных нужны меры, которые удовлетворяли бы всем этим требованиям. Часто достаточно, чтобы сохранилась информация о расстоянии между объектами лишь с точностью

до порядка, так что требование "г" можно было бы ослабить, а "д" снять совсем. Однако мы попытаемся найти более универсальную меру, удовлетворяющую всем требованиям "а"- "д".

## §2. Меры близости между объектами

Не будем подробно останавливаться на сильных шкалах (от шкалы интервалов и выше). Для них, как хорошо известно, свойствам "а"- "д" удовлетворяет мера

$$d_{ij}^c = \frac{|x_i - x_j|}{x_{\max} - x_{\min}}.$$

Оговоримся лишь, что для частного случая, когда  $x_{\max} = x_{\min}$ , неопределенное отношение  $\frac{0}{0}$  принимается за 0.

Рассмотрим шкалу порядка. Напомним, что при всех допустимых для этой шкалы преобразованиях  $\phi^n$  отношения из набора ( $<$ ,  $>$ ,  $=$ ) между двумя числами  $x_i$  и  $x_j$  должны сохраняться и для чисел  $\phi^n(x_i)$  и  $\phi^n(x_j)$ . Если мы построим матрицу  $M^n$  размером  $m \times m$  (где  $m$  - число объектов в выборке A), в которой для каждой пары объектов укажем их отношение в шкале порядка, то эта матрица не изменится при всех преобразованиях  $\phi^n$ . Этим свойством матрицы и воспользуемся для конструирования меры  $d_{ij}^c$ . Естественно считать, что те два объекта одинаковы, которые имеют одинаковые порядковые отношения со всеми другими объектами из выборки A. И, наоборот, наиболее непохожие объекты будут иметь наибольшие различия в своих отношениях с другими объектами. Различия (d) в отношениях i-го и j-го объектов к некоторому k-му объекту будем считать по такому правилу:

$$d_{ij}^k = \begin{cases} 0, & \text{если} \\ & \left\{ \begin{array}{l} x_i > x_k \text{ и } x_j > x_k; \\ x_i < x_k \text{ и } x_j < x_k; \\ x_i = x_k \text{ и } x_j = x_k; \end{array} \right. \\ 1, & \text{если} \\ & \left\{ \begin{array}{l} x_i > x_k, \text{ а } x_j < x_k; \\ x_i < x_k, \text{ а } x_j > x_k; \end{array} \right. \\ 0,5, & \text{если} \\ & \left\{ \begin{array}{l} x_i = x_k, \text{ а } x_j \neq x_k; \\ x_i \neq x_k, \text{ а } x_j = x_k. \end{array} \right. \end{cases}$$

Суммарное различие между  $i$ -м и  $j$ -м объектами

$$d_{ij}^{\Pi} = \sum_{k=1}^n d_{ij}^k.$$

Легко видеть, что если  $x_i = x_j$ , то  $d_{ij}^{\Pi} = 0$ , и что для объектов  $i$  и  $j$ , имеющих максимально разные порядковые позиции,  $d_{ij}^{\Pi} = I$ . Очевидно выполнение и других требований к  $d_{ij}^{\Pi}$ .

Заметим, что существует один канонический способ присыпывания чисел упорядоченному множеству объектов, при котором можно получать то же значение  $d_{ij}^{\Pi}$ , не строя матрицу отношений  $M^{\Pi}$ . Эти числа (назовем их "нормированными рангами") строятся по такому правилу: первому по порядку объекту присыпывается число 1, второму - 2 и так до конца. Если встречается 1 объектов с одинаковым порядковым номером (так называемые "серии"), то всем этим объектам присыпывается номер  $x^i = \frac{1}{1-p} \sum_{a=1}^p (a+p)$ , где  $p$  - количество объектов, предшествовавших серии. После такой канонизации  $d_{ij}^{\Pi}$  находится по правилу  $d_{ij}^{\Pi} = |x_i^i - x_j^i|$ .

В том, что эта мера  $d_{ij}^{\Pi}$  тождественно равна мере, найденной ранее по матрице  $M^{\Pi}$ , легко убедиться на примере, приведенном на рис. I.

$x_{\Pi}$	$x_{\Pi}^1$
1	12
...	...
...	6,3
1	8
...	...
...	...
j	12
...	5
...	...
m	123

a)

b)

в)

Рис. I. Правила определения различия в икалах порядка:

- а) претокой в иcale порядке,
- б) претокой в иcale нормированных рангов,
- в) инвариантная матрица отношений  $M^{\Pi}$ .

Перейдем теперь к иcale и а и м е н о в а н и й. Допустимо преобразования  $\phi^{\Pi}$  для икала этого типа всегда сохраняют отношения равенства, так что при всех возможных переключениях в мат-

рице  $M^H$  размера  $m \times m$  будут сохраняться значения отношений между всеми парами объектов из выборки A в виде символов ( $=, \neq$ ).

Расстояние  $d_{i,j}^H$  между i-м и j-м объектами можно найти по матрице  $M^H$ , используя следующее правило:

$$d_{i,j}^H = \frac{1}{m} \sum_{k=1}^m d_{i,j}^k,$$

где

$$d_{i,j}^k = 0, \text{ если } \begin{cases} x_i = x_k \text{ и } x_j = x_k; \\ x_i \neq x_k \text{ и } x_j \neq x_k; \end{cases}$$

$$d_{i,j}^k = 1, \text{ если } \begin{cases} x_i = x_k, \text{ а } x_j \neq x_k; \\ x_i \neq x_k, \text{ а } x_j = x_k. \end{cases}$$

Эта мера расстояния в шкале наименований удовлетворяет требованиям "а" - "д".

Как отмечено в [5], мера  $d_{i,j}^H$  может быть найдена и без построения матрицы  $M^H$ , а прямо через числа  $m_i$  и  $m_j$ , указывающие частоту встречаемости объектов, имеющих имена, одинаковые с i-м и j-м объектами, соответственно (см. рис.2):

$$d_{i,j}^H = \frac{m_i + m_j}{m}; \quad d_{i,j}^H = 0. \quad (i \neq j)$$

$x^H$	a	a	a	b	b	c
a	=	=	=	$\neq$	$\neq$	$\neq$
a	=	=	=	$\neq$	$\neq$	$\neq$
a	=	=	=	$\neq$	$\neq$	$\neq$
b	$\neq$	$\neq$	$\neq$	=	=	$\neq$
b	$\neq$	$\neq$	$\neq$	=	=	$\neq$
c	$\neq$	$\neq$	$\neq$	$\neq$	$\neq$	=

a) б)

$$m_a = 3, \quad m_b = 2, \quad m_c = 1; \quad d_{ab}^H = \frac{5}{6},$$

$$d_{ac}^H = \frac{4}{6}, \quad d_{bc}^H = \frac{3}{6}.$$

Рис.2. Пример определения меры близости в шкале наименований: а) протокол в шкале наименований, б) инвариантная матрица отношений  $M^H$ .

Рассмотрим некоторые свойства этой меры. Если признак использует только два имени, тогда  $m_i + m_j = m$  и  $d_{i,j}^H = 1$ . Так что для бинарных признаков эта мера однозначно соответствует мере близости, вычисляемой через хэммингово расстояние.

Если признак  $x^H$  содержит более двух имен, то две большие группы объектов с именами  $i$  и  $j$  будут считаться более "далекими", чем две группы с малым числом одноименных объектов. Это свойство хорошо согласуется со стремлением минимизировать ошибку, возникающую при перепутывании имен: при равной цене каждой ошибки суммарные потери будут пропорциональны числу перепутываемых объектов, так что две крупные группы одноименных объектов следует различать более уверенно, при таксономии объединять их в общий таксон в последнюю очередь и т.п. Перепутывание мелких групп, ошибочное объединение их в один таксон приведет к меньшим суммарным потерям.

Мерами  $d^C$ ,  $d^P$ ,  $d^H$  теперь можно пользоваться в многомерном случае, определяя расстояние  $d^P$  в пространстве разнотипных признаков по правилам следующего вида (для трех разнотипных признаков):

$$d^P = \frac{\alpha^C \cdot d^C + \alpha^P \cdot d^P + \alpha^H \cdot d^H}{\alpha^C + \alpha^P + \alpha^H};$$

$$d^P = \sqrt{\frac{(\alpha^C d^C)^B + (\alpha^P d^P)^B + (\alpha^H d^H)^B}{(\alpha^C)^B + (\alpha^P)^B + (\alpha^H)^B}};$$

и т.п. Меры такого типа также удовлетворяют требованиям "а"- "л".

Весовые положительные коэффициенты  $\alpha$  следовало бы задавать как функцию от априорно известной информативности признака, однако информативность признака зависит от решаемой задачи и становится известной обычно уже после того, как задача решена. В условиях такой неопределенности полагают все  $\alpha = 1$ . Но в случае разнотипных признаков можно, по-видимому, задавать  $\alpha$  пропорционально потенциальной информативности шкалы. Ведь известно, что один и тот же признак несет больше информации, если его измерять, например, в шкале отношений, чем в шкале порядка, и тем более в шкале наименований. Так что при отсутствии априорных данных об актуальной информативности каждого признака можно было бы считать, что  $\alpha^C > \alpha^P > \alpha^H$ .

Количественные соотношения потенциальной информативности шкал разного типа еще, к сожалению, не изучены. Так что никаких конкретных рекомендаций (кроме тривиальной  $\alpha = 1$ ) дать пока нельзя.

### §3. Мера расстояния между разнотипными признаками

При корреляционном, регрессионном анализе, при обработке групповых экспертных оценок и в других задачах анализа данных нужно измерять расстояния между признаками. В литературе известны разные меры, применяемые для пар однотипных признаков.

Здесь мы попытаемся сконструировать меру, которая была бы применима для пар как однотипных, так и разнотипных признаков и при этом уловлетворяла бы всем вышеприведенным требованиям "а" - "д".

Начнем с однотипных. Если признаки  $x_1$  и  $x_2$  измерены в шкалах, более сильных, чем шкала порядка, то всем требованиям "а"- "д" удовлетворяет мера  $d^{CC} = 1 - |r|$ , где  $r$  - коэффициент линейной корреляции.

Среди многочисленных мер расстояния двух признаков, измеренных в шкале порядка [6], своей простотой и естественностью отличается мера Кенделла-Кемени [7,8]. Для ее определения нужно перебрать все  $C_m^2$  парных сочетаний из  $m$  объектов и для каждой пары  $(a, b)$  сравнить порядковое отношение по признаку  $x_1$  и  $x_2$ . Если порядковые отношения одинаковы, т.е. если  $x_1(a) > x_1(b)$  и  $x_2(a) > x_2(b)$  или  $x_1(a) < x_1(b)$  и  $x_2(a) < x_2(b)$  или  $x_1(a) = x_1(b)$  и  $x_2(a) = x_2(b)$ , то  $d_{ab} = 0$ . Если отношения порядка на объектах  $(a, b)$  по признаку  $x_1$  прямо противоположны порядку по  $x_2$  (т.е. при  $x_1(a) > x_1(b)$  и  $x_2(a) < x_2(b)$  или  $x_1(a) < x_1(b)$  и  $x_2(a) > x_2(b)$ ), то  $d_{ab} = 1$ . В промежуточном случае, когда по одному признаку имеет место отношение ( $>$ ) или ( $<$ ), а по другому ( $=$ ), считается, что  $d_{ab} = 0,5$ .

Общее расстояние определяется как средняя мера "несогласия" двух признаков на всех парах объектов

$$d^{PP} = \frac{\sum_{a,b=1}^m d_{ab}}{C_m^2} .$$

Если упорядочивания на всех объектах совпадают, то  $d^{PP} = 0$ , если они на всех объектах противоположны, то  $d^{PP} = I$ .

Если один признак ("эксперт"  $x_1$ ) в порядковой шкале все объекты различает, а другой эти объекты считает одинаковыми (т.е. если  $x_2$  выдает "серию" линий  $m$ ), то мера  $d^{PP} = 0,5$ , что вполне естественно. Уместно отметить, что рекомендуемая во многих пособиях мера Спирмэна [5,7,8] в этом случае дает  $d^{PP} = \infty$ .

Мера расстояния между двумя признаками, измеренными в шкале наименований, определяется по правилу, аналогичному приведенному: перебираются все пары объектов ( $a, b$ ), и если отношение по признаку  $x_1$  совпадает с отношением по признаку  $x_2$ , т.е. если  $x_1(a) = x_1(b)$  и  $x_2(a) = x_2(b)$  или  $x_1(a) \neq x_1(b)$  и  $x_2(a) \neq x_2(b)$ , то  $d_{ab} = 0$ . Если же эти отношения по двум признакам не совпадают, т.е. при  $x_1(a) = x_1(b)$  и  $x_2(a) \neq x_2(b)$  или  $x_1(a) \neq x_1(b)$  и  $x_2(a) = x_2(b)$ , то  $d_{ab} = 1$ . В итоге

$$d^{HH} = \frac{\sum_{a,b=1}^n d_{ab}}{C^2}.$$

Величина  $d^{HH}$  в точности равна величине хэммингова расстояния между двоичными матрицами смежности, одна из которых построена по признаку  $x_1$ , а вторая - по признаку  $x_2$ .

Перейдем теперь к разнотипным парам признаков. Оба признака можно сделать однотипными, если один из них "обелить" до более слабого или второй "усилить" ("цифровать" [9]) до более сильного.

Сделаем по очереди то и другое, для каждого случая найдем свою меру расстояния. Общую меру расстояния будем определять как средневзвешенную величину этих двух частных расстояний. Рассмотрим пару  $x^c, x^\Pi$  (см. пример на рис. 3).

a)	$x^c$	$x^\Pi$	<u>ослабление</u>	$x^\Pi$	$x^\Pi$
	3, 5	12			5
	2, 3	6		2	2
	7, 0	8		4	3
	2, 1	1		1	1
	11, 0	50	$x^c$ до $x^\Pi$	6	6, 5
	15, 6	153		8	8
	13, 2	50		7	6, 5
	10	10		5	4

$$d^{\Pi\Pi} = \frac{2,5}{28} = 0,0893.$$

Рис.3. Пример ослабления сильной шкалы  $x^c$  до шкалы порядка  $x^\Pi$ : а) протокол в разнотипных шкалах, б) протокол в шкале порядка.

Ослабление сильной шкалы (признак  $x^c$ ) до шкалы порядка ( $x^\Pi$ ) состоит в том, что на числах  $x^c$  и  $x^\Pi$  учитываются только отношения порядка и находится  $d^{\Pi\Pi}$  по методу, изложенному выше.

Усиление ("оцифровка") шкалы порядка ( $x^{\Pi}$ ) до сильной ( $\hat{x}^c$ ) делается так, чтобы: а) значение порядка объектов по признаку  $\hat{x}^c$  совпадало с тем, который указал признак  $x^{\Pi}$ , и б) числовые значения признака  $\hat{x}^c$  были максимально коррелированы со значениями признака  $x^c$ . Достигается это способом, показанным на рис.4.Объект-

$x^c$	$x^{\Pi}$	$x^c$	$x^{\Pi}$	$x^c$	$\hat{x}^c$
3,3	12	2,1	1	2,1	2,1
2,3	6	7,0	8	7,0	4,2
2,0	8	3,5	12	3,5	4,3
2,1	1	2,4	16	2,4	4,4
1,0	50	10	10	10	10
15,6	153	13,2	50	11,0	12,0
13,4	50	11,0	50	8,8	12,2
10	10	15,6	153	15,6	15,6

a)

б)

в)

Рис.4. Пример усиления шкалы порядка  $x^{\Pi}$  до сильной шкалы  $\hat{x}^c$ : а) протокол в разнотипных шкалах, б) протокол, упорядоченный по  $x^{\Pi}$ , в) протокол в сильной шкале,  $d^{cc} = 0,074$ .

ты упорядочиваются по возрастанию значений признака  $x^{\Pi}$ . Если встречается "серия" (т.е.  $l > 1$  объектов с одинаковым порядковым номером), то они выписываются в порядке убывания их признака  $x^c$ . Затем, начиная с первого объекта, по значениям признака  $x^c$  выделяются "блоки инверсий", т.е. такие последовательности объектов, которые начинаются объектом  $i$  и заканчиваются самым далеким от него объектом  $j$ , такими, что  $x_i^c \geq x_j^c$ .

Если отношение порядка по  $x^{\Pi}$  и  $x^c$  совпадает, то в каждом блоке инверсии будет по одному объекту и усиление  $x^{\Pi}$  по  $\hat{x}^c$  делается присвоением  $i$ -му объекту значения  $\hat{x}_i^c = x_i^c$ .

Если порядки по  $x^{\Pi}$  и  $x^c$  не совпадают, то некоторые блоки инверсии ( $s$ ) будут сопретать по  $t > 1$  объектов. Требование "б" для этих объектов будет выполнено, если присвоить каждому из них их среднее значение по признаку  $x^c$ :  $\hat{x}_s^c = \frac{1}{t} \sum_{k=1}^t x_k^c$ . Для удовлетворения требования "а" объекты из  $S$ , различающиеся по  $x^{\Pi}$ , нужно "раздвинуть" на величину  $\Delta$ , так чтобы они равномерно заняли диапазон от

$x_s^c - \frac{t \cdot \Delta}{2}$  до  $x_s^c + \frac{t \cdot \Delta}{2}$ , где  $\Delta$  - разрешающая способность признака  $x^c$  (на рис. 4  $\Delta = 0,1$ ). Затем вычислим коэффициент линейной корреляции  $r$  и через его модуль  $|r|$  найдем величину  $d^{cc} = 1 - |r|$ . Средняя мера близости определяется так:

$$d^{cp} = \frac{\gamma^{\Pi} \cdot d^{\Pi\Pi} + \gamma^c \cdot d^{cc}}{\gamma^{\Pi} + \gamma^c}.$$

Здесь  $\gamma^{\Pi}$  и  $\gamma^c$  - весовые коэффициенты, показывающие величину нашего доверия к каждой из составляющих величин ( $d^{\Pi\Pi}$  и  $d^{cc}$ ). При одинаковом доверии к ним можно взять  $\gamma^{\Pi} = \gamma^c = 1$ , и тогда в нашем примере  $d^{cp} = \frac{0,0893 + 0,074}{2} = 0,0816$ .

Для пары  $x^c x^H$  (см. рис. 5) ослабление сильной шкалы  $x^c$  по шкале наименований ( $x^H$ ) состоит в том, что

a)	$x^c$	$x^H$	b)	$x^c$	$x^H$	v)	$x^H$	$x^H$
	3,5	b		k	k		3,5	3,5
	2,6	c		c	c		2,6	4,1
	3,5	a	ослабление	k	a	усиление	3,5	-0,1
	8,3	c		a	c		8,3	4,1
	1,4	c		p	c		1,4	4,1
	-3,7	a		p	a		-3,7	-0,1

Рис.5. Пример усиления и ослабления для шкал наименований  $x^H$  и сильных  $x^c$ : а) протокол в разных шкалах, б) протокол в шкале наименований, в) протокол в сильной шкале.

всем различным значениям  $x^c$  приписываются разные имена. Затем вычисляется мера  $d^{HH}$  по правилу, описанному для  $d^{HH}$ . При усилении ("оцифровке") признака  $x^H$  до  $x^c$  объекту  $a$  приписывается значение  $x^c(a) = x^c(a)$ . Если одинаковое имя  $b$  имеет "серия" из 1 объектов, то всем им приписывается

$$x^c(b) = \frac{1}{1} \sum_{i=1}^1 x_i^c(b).$$

По полученным числовым значениям вычисляется мера  $d^{cc}$  по правилу  $d^{cc}$  (через корреляцию). Средняя мера

$$d^{CH} = \frac{\gamma^H \cdot d^{HH} + \gamma^C \cdot d^{CC}}{\gamma^H + \gamma^C}.$$

Здесь  $\gamma^H$  и  $\gamma^C$  имеют тот же смысл, что и выше, и при  $\gamma^H = \gamma^C = 1$  в нашем примере  $d^{HH} = 0,33$ ,  $d^{CC} = 0,34$ ,  $d^{CH} = 0,335$ .

Наконец, для пары  $x^H$  и  $x^P$  (см. рис.6) обединение, как и в предыдущем случае, сводится к присваиванию разных имен объектам, имеющим разные порядковые номера. После чего находится  $d^{HP}$ . При усилении  $x^H$  до  $x^P$  признак  $x^P$  канонизируется по нормированных рангов, а затем вместо имен  $x^H$  объектам ставятся в соответствие порядковые номера так же, как и в предыдущем случае:  $x^P(a)$ ,  $x^P(b)$  для одиночных объектов и  $x^P(b) = \frac{1}{1} \sum_{i=1}^1 x_i^P(b)$  для серии из однаково поименованных объектов.

a)	$x^P$	$x^H$	б)	$x^P$	$x^H$	в)	$x^P$	$x^P$
1	a		a	a		1	1	
6	b		c	b		3	3	
6	d	ослабление	c	d	услечение	3	4	
6	b		c	b		3	3	
8	d		p	d		5	4	

Рис.6. Пример усиления и ослабления для шкал наименований  $x^H$  и порядка  $x^P$ : а) протокол в разных шкалах, б) протокол в шкале наименований, в) протокол в шкале порядка .

После нахождения  $d^{PP}$  (по правилу  $d^{PP}$ ), можно определить

$$d^{PH} = \frac{\gamma^P \cdot d^{PP} + \gamma^H \cdot d^{HH}}{\gamma^P + \gamma^H}.$$

При  $\gamma^P = \gamma^H = 1$  в примере на рис.6  $d^{HH} = 0,3$ ,  $d^{PP} = 0,15$ ,  $d^{PH} = 0,225$ .

Все промежуточные преобразования шкал входят в группу допустимых для своих шкал. Так что все приведенные меры также инвариантны к допустимым преобразованиям оцениваемых признаков. Усиление и ослабление шкал в некотором смысле вносят симметричную помеху (добавление и потерю информации), так что усреднение после этих процедур будет оправданным.

Применение указанных мер расстояния между объектами и между признаками позволяет использовать все то богатство математического обеспечения, которое имеется сейчас для анализа таблиц данных, измеренных в сильных шкалах. При этом нужно в программах добавить семантический блок, указывающий тип шкалы данного признака и заменить блок определения "расстояния" на блок вычисления соответствующей меры из набора, описанного выше.

### Л и т е р а т у р а

1. ВОРОНИН Ю.А. Введение мер сходства и связи для решения геолого-географических задач. -Докл. АН СССР, 1971, т. 199, № 5, с. 1011-1015.
2. ШУСТОРОВИЧ А.М. Об адекватных парных мерах сходства в задачах распознавания образов с разнородными признаками. -В кн.: Вопросы обработки информации при проектировании систем (Вычислительные системы, вып. 69). Новосибирск, 1977, с. 147-152.
3. ЗАГОРУЙКО Н.Г. Таксономия в анизотропном пространстве. -В кн.: Эмпирическое предсказание и распознавание образов (Вычислительные системы, вып. 76). Новосибирск, 1978, с. 26-34.
4. СУППЕС П., ЗИНЕС Дж. Основы теории измерений.-В кн.: Психологические измерения. -М., Мир, 1976, с. 120.
5. ЧЕРНЫЙ Л.Б. Порождение мер связи между объектами с помощью мер связи между признаками. -В кн.: Проблемы анализа дискретной информации. Новосибирск, 1975, с. 167-174.
6. ВОРОНИН Ю.А., ГРАДОВА Т.А. Регион-комплекс программ для постановки и решения задач районирования. Препринт № 258 ВЦ СО АН СССР, Новосибирск, 1980.
7. КЕНДЕЛ М. Ранговые корреляции.-М.: Статистика, 1976.-350 с.
8. КЕМЕНИ Дж., СНЕЛЛ Дж. Кибернетическое моделирование. -М.: Сов. радио, 1972. - 320 с.
9. ЕНОКОВ И.С., КУЛАКОВА Е.П. Числовые метки для качественных признаков в дискриминантном анализе. -В кн.: Прикладной многомерный статистический анализ. Ученые записки по статистике, т.33. М., 1978, с. 353-358.

Поступила в ред.-изд.отд.  
12 декабря 1982 года