

ПРИМЕРНЫЕ ДОКЛАДЫ

# **КОГНИТИВНЫЕ ШТУДИИ**

**Выпуск 2**

**Когнитивная психология  
и  
искусственный интеллект**

Минск 2002  
НЕССИ

УДК 159.99  
ББК 88.4  
К57

**Редакционная коллегия**

Голенков В.В. (БГУИР)  
Лобанов А.П. (БГПУ им. М. Танка)  
Лосик Г.В. (ИТК НАН Беларуси)  
Радчикова Н.П. (БГПУ)  
Репеко А.П. (БГУ)  
Шатон Г.И. (ЖНИ "ЭНВИЛА")

**К 57**            **Когнитивные штудии.** Выпуск 2 "Когнитивная психология и искусственный интеллект": сборник статей. - Мн.: НЕССИ, 2002. - 156 с.

**ISBN 985-6188-57-1**

В настоящий сборник вошли научные работы, представленные на II междисциплинарном теоретическом семинаре по когнитивной науке, прошедшем в г. Минске в 2001 г. под эгидой БГПУ им. М. Танка, ИТК НАН Беларуси и ЖНИ "Энвила" и посвященном различным проблемам когнитивной психологии и искусственного интеллекта.

УДК159.99  
ББК88.4

ISBN 985-6188-57-1

© КОЛЛЕКТИВ АВТОРОВ

1. Римский Г.В., Бочкарёва Л.В., Таборовец О.В., Мазуренко В.А., Римский А.Г. Автоматно временной анализ процессов на семантическом полигоне. Весті Нацыянальнай акадэміі навук Беларусі, сер фіз.-тэхн. навук, № 1, 1999. – с. 76-83.
2. Римский Г.В., Бочкарёва Л.В. Преобразование информации в интеллектуальных фразеологических информационно-поисковых системах. Тезисы докладов 1-го Международного конгресса по информатике, безопасности, экологии. Минск: БелЭКСПО, 9-10 декабря 1999.
3. Римский Г.В., Бочкарёва Л.В., Бабинец О.О., Таборовец О.В. Построение интеллектуальных фразеологических информационно-поисковых систем // Компьютерная лингвистика и обучение языкам / Сборник научных статей. - Минск, МГЛУ, 2000. - с.131-142.
4. Дубов Ю.А., Травкин С.И., Якимец В.Н. Многокритериальные модели формирования и выбора вариантов систем. -М:Наука. 1986г.-296 с.
5. Римский Г.В. Чебаков С.В. Метод автоматического конструирования программ при многокритериальном бинарном отношении предпочтения между альтернативами // Программирование.-1985г.-N5.-с. 38-41.
6. Гафт М.Г. Озерной В.М. Выделение множества неподчиненных решений и их оценок в задачах принятия решений при векторном критерии // Автоматика и телемеханика. 1973г. N11. -с 85-95.
7. Римский Г.В. Теория систем автоматического проектирования. - Мн.: Навука і тэхніка, 1994. - 631 с.

#### **КАТЕГОРИИ И ПОНЯТИЯ: ВВЕДЕНИЕ В АНАЛИЗ ДАННЫХ**

Н.П. Радчикова  
Белорусский государственный педагогический университет,  
Институт технической кибернетики НАН Беларуси

## 1. ПРОБЛЕМА КАТЕГОРИЗАЦИИ

Для людей, сталкивающихся на каждом шагу с огромным потоком информации, важно уметь делить эту информацию на значимые части. Это и есть основная проблема категоризации. Категоризация начинается на ранних этапах восприятия, когда стимулы соединяются с более абстрактными понятиями, о чем свидетельствует факт, что одинаковая физическая разница между стимулами воспринимается большей или меньшей в зависимости от того, принадлежат они разным категориям или одной (Harnad, 1987). Более того, образование понятий и категориальной структуры - основная часть любого процесса обучения (А.И. Розов, 1985; Medin&Ross, 1997). Категории и их мыслительные представления (понятия) составляют основу для процессов мышления.

Категоризация не только имеет большое значение в повседневной жизни, но и является объектом изучения многих отраслей науки. Хорошо известны, например, усилия биологов в разработке таксономических систем флоры и фауны. В медицине и психиатрии категоризация определенных симптомов как некоторого заболевания также является вопросом чрезвычайной важности. Самой же проблемой категоризации (т.е. определению законов, правил и алгоритмов, по которым происходит отнесение объекта к определенному классу или категории) занимается когнитивная психология.

При изучении процесса категоризации обычно возникают два основных типа вопросов:

1. Какие объекты группируются вместе, образуя категорию?
2. Каким образом можно охарактеризовать (описать) категорию (и ее членов) с помощью их характеристик (свойств)?

Первая из этих проблем называется иногда агрегацией (aggregation), а вторая - характеристикой (characterization) (Fisher&Langley, 1986).

Для деления информации на части, т.е. для категоризации, применяется довольно много вычислительных методов и процедур: это и кластерный анализ, и

дискриминантный анализ, и анализ соответствий, и анализ иерархических классов, а также линейная и логистическая регрессия. Методы анализа данных могут решать как первую проблему, так и вторую, некоторые - даже обе сразу. Кластерный анализ по определению - метод, занимающийся решением первой проблемы, а регрессионный анализ и дискриминантный анализ - популярные техники решения второй проблемы.

В данной статье мы остановимся на рассмотрении только некоторых из этих методов, и критерием выбора вычислительной процедуры будет служить возможность найти ее в стандартном пакете для статистической обработки данных (например, STATISTICA, SPSS, SyStat).

Следует отметить, что до сих пор методы анализа данных и когнитивная психология развивались не в тесном сотрудничестве, а скорее независимо друг от друга (Mechelin&Michalski, 1993). С одной стороны, когнитивные исследования проблемы категоризации и образования понятий основывались на прямом опросе испытуемых об интересующих исследователя аспектах, а не на использовании информации, полученной с помощью методов анализа данных. С другой стороны, методы анализа данных редко проверялись на когнитивную валидность их основных предположений, моделей и, тем более, алгоритмов. В лучшем случае можно предположить, что результаты анализа данных соответствуют потенциальной когнитивной структуре, и требуется дополнительная проверка (теоретические или эмпирические подтверждения) эффективности этого соответствия.

Несмотря на такое положение дел, сотрудничество в области когнитивной психологической традиции и методов анализа данных может быть очень продуктивным. Подход с точки зрения методов анализа данных может требовать от психологов более четкой формализации главных понятий и классификационных принципов. С другой стороны, рассмотрение методов анализа данных с точки зрения когнитивной психологии может способствовать лучшему пониманию психологической значимости моделей, лежащих в основе анализа данных. Принципы категоризации, открытые психологами, могут быть использованы для разработки новых

процедур и алгоритмов анализа данных, которые, в свою очередь, могут быть использованы для изучения категорий и понятий. Именно поэтому психологи должны иметь представления об алгоритмах и принципах, лежащих в основе статистических процедур, которыми они пользуются при анализе данных. В свою очередь математики и специалисты, занимающиеся искусственным интеллектом, должны быть в курсе последних разработок в области психологии, чтобы реализованные в статистических процедурах принципы были по возможности когнитивно валидными.

### 1.1. Формальные принципы категоризации

Модели статистического классификационного анализа данных можно разделить на основании множества формальных принципов категоризации. Эти принципы, в сущности, составляют ядро процедуры анализа данных.

Наиболее известным является принцип сходства (similarity principle), который утверждает, что категории состоят из похожих объектов. Со времен Аристотеля сходство является фундаментальным понятием, которое помогает вносить порядок и определять структуру. К измерению сходства есть два подхода (Figer&De Воеск, 1993):

- сходство может быть оценено испытуемым, который сообщает о своих субъективных ощущениях;
- сходство может быть вычислено из перечислений характеристик объектов.

Обобщенный термин для сходства, который часто употребляется в когнитивной психологии и искусственном интеллекте, - близость (proximity). Близость включает также различие, ассоциацию, корреляцию и т.д. (Arabie, Carroll, & DeSarbo, 1987), и часто является синонимом слова *сходство*.

Принцип сходства лежит в основе всех непрямых кластерных процедур, которые используют данные о близости. С точки зрения когнитивной психологии этот принцип, который связан с принципом семейного сходства Людвига Витгенштейна (Wittgenstein, 1953) и Э. Рош (Rosch, 1978), являлся весьма перспективным. В последние годы, однако, некоторые исследователи обнаружили недостатки этого

принципа как основания для категоризации (см., например, Murphy, 1993; Rips, 1989; Michalski&Stepp, 1983).

Следующий принцип может быть назван монотетическим (monothetic principle). Он утверждает, что категории состоят из объектов, которые принимают одни и те же значения на некотором множестве характеристик (свойств). Этот принцип лежит в основе модели иерархических классов (De Voeck, Rosenbergg, &Van Mechelin, 1993). Монотетический принцип - это формальный перевод классического (формально-логического) подхода к категоризации - наличия существенных свойств, каждое из которых в отдельности необходимо, а все вместе - достаточны для членства в категории (Гетманова, 1997).

Третий важный принцип - корреляционный. Этот принцип утверждает, что категории формируются таким образом, что корреляция между свойствами членов категорий полностью объясняется категориальной структурой. Этот принцип лежит в основе анализа латентных классов (latent class analysis). Его отзвуки можно найти в выдвинутом Элеонор Рош когнитивном принципе о "воспринимаемой структуре мира", который утверждает, что категории отражают корреляционную структуру характеристик в воспринимаемом объективном мире (Rosch, 1978).

Четвертый формальный принцип, принцип краткости (parsimony principle), является принципом другого порядка, так как он всегда работает вместе с каким-либо из первых трех принципов. Одна его версия утверждает, что категории всегда образуют простую таксономическую структуру. Эта версия используется, например, в иерархическом кластерном анализе для разбиения кластерной иерархии. Он также используется в иерархическом анализе классов (hierarchical class analysis). Другая версия принципа краткости используется в методе выведения классификационного правила для категорий, заданных заранее. Она утверждает, что классификационное правило должно быть самым простым из всех возможных. Эта версия используется, например, в пошаговом применении дискриминантного анализа. Когнитивный эквивалент принципа краткости - это, безусловно, принцип когнитивной экономии.

## 1.2. Типы данных

Данные для проведения статистического классификационного анализа могут быть одного из следующих типов:

1. Данные о сходстве *объект*  $\times$  *объект*. Для всех пар данного множества из  $n$  элементов либо собираются данные о сходстве, либо из описаний объектов вычисляется мера сходства. В результате получается матрица размерности  $n \times n$ .
2. Описательные данные *объект*  $\times$  *характеристика*. Когда  $n$  объектов описывается с помощью  $m$  характеристик, получается матрица, которая может быть проанализирована без вычислений значений сходства.
3. Описательные данные *объект*  $\times$  *степень\_свойства*. Для некоторого множества объектов измеряется степень членства в категории, типичность, прототипичность или другая похожая мера.

Классификационный анализ данных позволяет ответить нам на следующие вопросы:

1. Какие объекты являются членами одной и той же категории (из множества категорий, которые заранее не известны)?
2. Какие характеристики (свойства) характеризуют категорию (которая известна заранее)?
3. Какие объекты являются членами одной и той же (неизвестной) категории и какими характеристиками (свойствами) эта (неизвестная) категория обладает?

Каждому из вопросов соответствует свой тип данных. Первому вопросу - данные о сходстве *объект*  $\times$  *объект*, второму вопросу - информация *объект*  $\times$  *характеристика* или *объект*  $\times$  *степень\_свойства*, третьему - *объект*  $\times$  *характеристика*.

Интересно, что во всех трех случаях информация, которая получается в результате анализа, является информацией совсем другого типа, чем исходные данные. Более того, во всех трех случаях полученные данные могут быть рассмотрены как информация более высокого уровня обобщенности, чем исходная. Для первого вопроса данные



содержат информацию об отношениях объект-объект, а в результате анализа получают данные об отношениях категория-категория и категория-объект. Для второго вопроса собираются данные о членстве объектов в категории, а выводы делаются о свойствах (характеристиках) категории. Для третьего вопроса из данных о свойствах (характеристиках) объектов делаются выводы об отношениях категория-категория, категория-объект и категория-характеристика.

## 2. МЕТОДЫ КЛАССИФИКАЦИОННОГО АНАЛИЗА ДАННЫХ

### 2.1. Кластерный анализ

Термин *кластерный анализ* (впервые ввел Трюон, 1939) в действительности включает в себя набор различных алгоритмов классификации. Общий вопрос, задаваемый исследователями во многих областях, состоит в том, как организовать наблюдаемые данные в наглядные структуры, т.е. развернуть таксономии. Например, биологи ставят цель разбить животных на различные виды, чтобы содержательно описать различия между ними. В соответствии с современной системой, принятой в биологии, человек принадлежит к приматам, млекопитающим, позвоночным и животным. Заметим, что в этой классификации, чем выше уровень обобщенности, тем меньше сходства между членами в соответствующем классе. Человек имеет больше сходства с другими приматами (т.е. с обезьянами), чем с "отдаленными" членами семейства млекопитающих (например, собаками) и т.д.

Техника кластеризации применяется в самых разнообразных областях. Дж. Хартиган (Hartigan, 1975) дал прекрасный обзор многих опубликованных исследований, содержащих результаты, полученные методами кластерного анализа. Например, в области медицины кластеризация заболеваний, лечения заболеваний или симптомов заболеваний приводит к широко используемым таксономиям. В области психиатрии правильная диагностика кластеров симптомов, таких как паранойя, шизофрения и т.д., является решающей для успешной терапии. В археологии с помощью кластерного анализа исследователи пытаются установить таксономии каменных орудий, похоронных объектов и т.д.

Известны широкие применения кластерного анализа в маркетинговых исследованиях. В общем, всякий раз, когда необходимо классифицировать "горы" информации к пригодным для дальнейшего использования группам, кластерный анализ оказывается весьма полезным и эффективным.

Под кластером обычно понимают группу объектов, обладающих свойством плотности (плотность объектов внутри кластера выше, чем вне его), дисперсией, отделимостью от других кластеров, формой, размером. Интуитивно можно представить себе кластер как горы на географической карте или созвездия на небе (Боровиков, 2001b).

Методы кластеризации "делятся на агломеративные (от слова агломерат - скопление) и итеративные дивизивные (от слова division - деление, разделение)" (Боровиков, 2001 b: 183).

В агломеративных методах происходит последовательное объединение наиболее близких объектов в один кластер. Процесс такого последовательного объединения можно показать на графике в виде дендрограммы, или дерева объединения. Это удобное представление позволяет наглядно представить кластеризацию агломеративными методами.

Исходными данными для анализа могут быть собственно объекты и их характеристики или параметры, а могут быть и матрицы расстояний между объектами. Если расстояния не даны сразу, то агломеративные алгоритмы начинаются с вычисления расстояний между объектами.

Расстояние между объектами - одна из мер сходства. Понятно, что чем меньше расстояние между объектами, тем более они схожи. Однако расстояние между объектами можно определить по-разному. Часто используют обычную евклидову метрику. Например, если объект описывается двумя параметрами, то он может быть изображен точкой на плоскости, а расстояние между объектами - это расстояние, вычисленное по теореме Пифагора. Вы просто возводите в квадрат расстояния по каждой координате, суммируете их и из полученной суммы извлекаете квадратный корень. Если вы не будете возводить в квадрат по координатные расстояния, а просто возьмете их абсолютные значения и просуммируете, то получите так называемое манхэттенское расстояние, или

"расстояние городских кварталов". Такое расстояние связано с перемещением человека по улицам города, а не с движением по ровной местности.

Однако в реальных задачах могут встретиться объекты, которые по самой своей природе трудно представить точками на плоскости или в пространстве. В этих задачах используются меры сходства, которые измеряются непосредственно. Например, проводится массовый опрос: сходно ли Ваше отношение к определенным событиям? К двум цветам?

Важной мерой сходства, которая традиционно используется в социальных науках, являются коэффициенты корреляции. Для бинарных данных часто просто вычисляют количество параметров, которые совпадают у объектов. Далее это число делят на общее число параметров и получают меру сходства. Таким образом построенные меры называют коэффициентами ассоциативности.

## 1.2. Дискриминантный анализ

Дискриминантный анализ используется для принятия решения о том, какие переменные различают (дискриминируют) две или более возникающие совокупности (группы). Например, некий исследователь в области образования может захотеть исследовать, какие переменные относят выпускника средней школы к одной из трех категорий: (1) поступающий в вуз, (2) поступающий в техникум или (3) отказывающийся от дальнейшего образования или профессиональной подготовки. Для этой цели исследователь может собрать данные о различных переменных, связанных с учащимися школы. После выпуска большинство учащихся естественно должно попасть в одну из названных категорий. Затем можно использовать *Дискриминантный анализ* для определения того, какие переменные дают наилучшее предсказание выбора учащимися дальнейшего пути.

Медик может регистрировать различные переменные, относящиеся к состоянию больного, чтобы выяснить, какие переменные лучше предсказывают, что пациент, вероятно, выздоровел полностью (группа 1), частично (группа 2) или совсем не выздоровел (группа 3). Биолог может записать различные характеристики сходных типов (групп) цветов,

чтобы затем провести анализ дискриминантной функции, наилучшим образом разделяющей типы или группы.

С вычислительной точки зрения дискриминантный анализ очень похож на дисперсионный анализ. Рассмотрим следующий простой пример. Предположим, что вы измеряете рост в случайной выборке из 50 мужчин и 50 женщин. Женщины в среднем не так высоки, как мужчины, и эта разница должна найти отражение для каждой группы средних (для переменной *Рост*). Поэтому переменная *Рост* позволяет вам провести дискриминацию между мужчинами и женщинами лучше, чем, например, вероятность, выраженная следующими словами: "Если человек большой, то это, скорее всего, мужчина, а если маленький, то это вероятно женщина".

Вы можете обобщить все эти доводы на менее "тривиальные" группы и переменные. Например, предположим, что вы имеете две совокупности выпускников средней школы - тех, кто выбрал поступление в университет, и тех, кто не собирается это делать. Вы можете собрать данные о намерениях учащихся продолжить образование в университете за год до выпуска. Если средние для двух совокупностей (тех, кто в настоящее время собирается продолжить образование, и тех, кто отказывается) различны, то вы можете сказать, что намерение поступить в университет, как это установлено за год до выпуска, позволяет разделить учащихся на тех, кто собирается и кто не собирается поступать в университет.

В завершение заметим, что основная идея дискриминантного анализа заключается в том, чтобы определить, отличаются ли совокупности по среднему какой-либо переменной (или линейной комбинации переменных), и затем использовать эту переменную, чтобы предсказать для новых членов их принадлежность к той или иной группе.

### **Пошаговый дискриминантный анализ**

Вероятно, наиболее общим применением дискриминантного анализа является включение в исследование многих переменных с целью определения тех из них, которые наилучшим образом разделяют совокупности между собой. Например, исследователь в области образования, интересующийся предсказанием выбора, который сделают выпускники средней школы относительно

своего дальнейшего образования, произведет с целью получения наиболее точных прогнозов регистрацию возможно большего количества параметров обучающихся, например, мотивацию, академическую успеваемость и т.д. Другими словами, вы хотите построить "модель", позволяющую лучше всего предсказать, к какой совокупности будет принадлежать тот или иной образец.

**Пошаговый анализ с включением.** В пошаговом анализе дискриминантных функций модель дискриминации строится по шагам. Точнее, на каждом шаге просматриваются все переменные и находится та из них, которая вносит наибольший вклад в различие между совокупностями. Эта переменная должна быть включена в модель на данном шаге, и происходит переход к следующему шагу.

**Пошаговый анализ с исключением.** Можно также двигаться в обратном направлении; в этом случае все переменные будут сначала включены в модель, а затем на каждом шаге будут устраняться переменные, вносящие малый вклад в предсказания. Тогда в качестве результата успешного анализа можно сохранить только "важные" переменные в модели, то есть те переменные, чей вклад в дискриминацию больше остальных.

Пошаговый дискриминантный анализ основан на использовании статистического уровня значимости. Поэтому по своей природе пошаговые процедуры рассчитывают на случай, так как они "тщательно перебирают" переменные, которые должны быть включены в модель для получения максимальной дискриминации. При использовании пошагового метода исследователь должен осознавать, что используемый при этом уровень значимости не отражает истинного значения *альфа*, то есть вероятности ошибочного отклонения гипотезы  $H_0$  (нулевой гипотезы, заключающейся в том, что между совокупностями нет различия).

### **Классификация**

Другой главной целью применения дискриминантного анализа является проведение классификации. Как только модель установлена и получены дискриминирующие функции, возникает вопрос о том, как хорошо они могут *предсказывать*, к какой совокупности принадлежит конкретный образец?

Обычно, если вы оцениваете на основании некоторого множества данных дискриминирующую функцию, наилучшим образом разделяющую совокупности, и затем используете *те же самые* данные для оценивания того, какова точность вашей процедуры, то вы во многом полагаетесь на волю случая. В общем случае, получают, конечно, худшую классификацию для образцов, не использованных для оценки дискриминантной функции. Другими словами, классификация действует лучшим образом для выборки, по которой была проведена оценка дискриминирующей функции (*апостериорная* классификация), чем для свежей выборки (*априорная* классификация). Поэтому оценивание качества процедуры классификации никогда не производят по той же самой выборке, по которой была оценена дискриминирующая функция. Если желают использовать процедуру для классификации будущих образцов, то ее следует "испытать" (произвести кросс-проверку) на новых объектах.

**Функции классификации.** Функции классификации не следует путать с дискриминирующими функциями. Функции классификации предназначены для определения того, к какой группе наиболее вероятно может быть отнесен каждый объект. Имеется столько же функций классификации, сколько групп. Как только вы вычислили показатели классификации для наблюдений, легко решить, как производить классификацию наблюдений. В общем случае наблюдение считается принадлежащим той совокупности, для которой получен наивысший показатель классификации (кроме случая, когда вероятности *априорной* классификации становятся слишком малыми). Поэтому, если вы изучаете выбор карьеры или образования учащимися средней школы после выпуска (поступление в университет или получение работы) на основе нескольких переменных, полученных за год до выпуска, то можете использовать функции классификации, чтобы предсказать, что наиболее вероятно будет делать каждый учащийся после выпуска. Однако вы хотели бы определить *вероятность*, с которой учащийся сделает предсказанный выбор. Эти вероятности называются *апостериорными*, и их также можно вычислить. Для каждой совокупности в выборке вы можете определить положение точки, представляющей средние для всех переменных в многомерном пространстве,

определенном переменными рассматриваемой модели. Эти точки называются *центроидами* группы. Для каждого наблюдения вы можете затем вычислить его расстояние Махаланобиса от каждого центроида группы. Вы признаете наблюдение принадлежащим к той группе, к которой он ближе, т.е. когда расстояние Махаланобиса до нее минимально.

**Апостериорные вероятности классификации.** Используя для классификации расстояние Махаланобиса, вы можете теперь получить вероятность того, что образец принадлежит к конкретной совокупности. Это значение будет не вполне точным, так как распределение вокруг среднего для каждой совокупности будет не в точности нормальным. Так как принадлежность каждого образца вычисляется по априорному знанию модельных переменных, эти вероятности называются *апостериорными* вероятностями. Короче, *апостериорные* вероятности - это вероятности, вычисленные с использованием знания значений других переменных для образцов из частной совокупности. Некоторые пакеты автоматически вычисляют эти вероятности для всех наблюдений (или для выбранных наблюдений при проведении кросс-проверки).

**Априорные вероятности классификации.** Имеется одно дополнительное обстоятельство, которое следует рассмотреть при классификации образцов. Иногда вы знаете заранее, что в одной из групп имеется больше наблюдений, чем в другой. Поэтому *априорные* вероятности того, что образец принадлежит такой группе, выше. Например, если вы знаете заранее, что 60% выпускников вашей средней школы обычно идут в университет, (20% идут в армию и остальные 20% идут работать), то вы можете уточнить предсказание таким образом: при всех других равных условиях более вероятно, что учащийся поступит в университет, чем сделает два других выбора. Вы можете установить различные *априорные* вероятности, которые будут затем использоваться для уточнения результатов классификации наблюдений (и для вычисления *апостериорных* вероятностей).

На практике исследователю необходимо задать себе вопрос, является ли неодинаковое число наблюдений в различных совокупностях в первоначальной выборке

отражением истинного распределения в популяции, или это только (случайный) результат процедуры выбора. В первом случае вы должны положить *априорные* вероятности пропорциональными объемам совокупностей в выборке; во втором - положить *априорные* вероятности одинаковыми для каждой совокупности. Спецификация различных *априорных* вероятностей может сильно влиять на точность классификации.

В общем, дискриминантный анализ - это очень полезный инструмент (1) - для поиска переменных, позволяющих отнести наблюдаемые объекты в одну или несколько реально наблюдаемых групп, (2) - для классификации наблюдений в различные группы.

#### 2.4. Анализ соответствий

Анализ соответствий (correspondence analysis) - это разведочный метод анализа, позволяющий визуально и численно исследовать структуру таблиц сопряженности большой размерности. В настоящее время он интенсивно применяется в разных областях - в социологии, экономике, медицине, при анализе текстов и т.д. Известно применение анализа соответствий для исследования голосования в ООН по принципиальным вопросам (например, в 1967 году исследовалось 127 стран по 13 важным вопросам). Анализ показал, что по первому фактору страны отчетливо разделяются на две группы: одна с центром США, другая - с центром СССР (двухполюсная модель мира). Другие факторы могут интерпретироваться как изоляционизм, неучастие в голосовании и т.д. (Боровиков, 2001а: 551).

Анализ соответствий опирается на статистику хи-квадрат и во многом похож на факторный анализ. Строки или столбцы исходной таблицы представляются точками пространства, между которыми вычисляется расстояние хи-квадрат (аналогично тому, как вычисляется статистика хи-квадрат для сравнения эмпирических и теоретических частот). Далее требуется найти пространство небольшой размерности, как правило, двумерное, в котором вычисленные расстояния минимально искажаются, и в этом смысле максимально точно воспроизводят структуру исходной таблицы с сохранением связей между признаками.



### 3. ЗАКЛЮЧЕНИЕ

Подход к изучению категорий и проблемы категоризации, использующий методы анализа данных, может быть очень продуктивным, так как он позволяет начинать исследование с типов данных, которые нельзя получить простым опросом о явлении, которое нас интересует. Таким образом можно избежать многих проблем, связанных с методом прямого опроса. В любом случае могут быть получены свидетельства, согласующиеся с результатом прямого опроса, т.е. использование прямого опроса и непрямых методов анализа данных не исключают, а дополняют друг друга. Напротив, взаимодействие двух методов может быть весьма продуктивным, чтобы, например, получить представление об отношениях между данными, относящимися к категориям на разных уровнях обобщенности.

Ядро методов анализа данных состоит из нескольких формальных принципов категоризации. Понятно, что эти принципы должны быть когнитивно валидными, чтобы и анализ данных имел какие-то шансы на валидность. Однако в настоящее время когнитивные психологи скорее заняты тем, что показывают невалидность и неприемлемость любых принципов классификации. Даже основной принцип сходства подвергся суровой экспериментальной критике. Поэтому основной задачей психологов будет, по-видимому, разработка новых принципов классификации и оценка старых с точки зрения когнитивной валидности.

Хороший пример — исследование Медина, Ваттенмейкера и Михальского (Medin, Wattenmaker, & Michalski, 1987), изучавших валидность конъюнктивного типа классификационного правила у людей и в машинном обучении. Это как раз тема, где специалисты, работающие в области искусственного интеллекта, и психологи могут помочь друг другу.

Поиск новых, когнитивно валидных принципов категоризации — многообещающее направление сотрудничества. Найденные принципы могут использоваться для разработки новых процедур анализа данных, а когнитивная валидность результатов анализа данных, в свою

очередь, может быть оценена, что требует четких статистических и психологических критериев.

Пожалуй, последний вопрос, который следует поднять в связи с возможным сотрудничеством психологии и искусственного интеллекта, — могут ли психологические принципы и результаты статистического анализа объяснить достаточную часть данных о процессе категоризации или же они скорее могут быть фальсифицированы с помощью необъясненной части данных? До сих пор исследователи, работающие в области категоризации, пытались объяснить все имеющиеся факты. Этот подход, при котором один контр-пример может опровергнуть теорию или процедуру анализа данных, привел к последовательности все более сложных моделей и техник. Классический (Аристотелевский) подход был заменен на теорию прототипов и базисного уровня, та в свою очередь, уступила место теории, "основанной на теории". Разработать теорию, полностью объясняющую все стороны исследуемого процесса, безусловно, является заманчивой задачей, однако, может быть, и не вполне выполнимой. Поэтому не менее важным является построение моделей, которые с достаточной точностью отражают реальность.

## ЛИТЕРАТУРА

- Боровиков В. (2001a) *STATISTICA: искусство анализа данных на компьютере. Для профессионалов*. - СПб.: Питер.
- Боровиков В. (2001b) *Программа STATISTICA для студентов и инженеров*. - Компьютер Пресс: Москва.
- Гетманова А.Д. (1997) *Учебник по логике*. М.
- Розов А.И. (1985) Проблемы категоризации: теория и практика. *Вопросы психологии* № 3, с. 90-97.
- Arabie, P., Carroll, J.D., and DeSarbo, W.S. (1987) *Three-way Scaling and Clustering*, Newbury Park, CA: Sage.
- De Boeck, P., Rosenbergg, S., and Van Mechelin, I. (1993) The hierarchical classes approach: a review. In Mechelin, I.V., Hampton, J., Michalski, R.S., Theuns, P. (eds.) *Categories and Concepts: Theoretical Views and Inductive Data Analysis*, Academic Press, pp. 287-308.
- Harnad, S (ed.) (1987) *Categorical Perception*, Cambridge: Cambridge University Press.

- Hartigan, J. A. (1975). *Clustering algorithms*. New York: Wiley.
- Figer, H., and De Boeck, P. (1993) Categories and Concepts: Introduction to Data Analysis. In Mechelin, I.V., Hampton, J., Michalski, R.S., Theuns, P. (eds.) *Categories and Concepts: Theoretical Views and Inductive Data Analysis*, Academic Press, pp. 203-224.
- Fisher, D., and Langley, P. (1986) Conceptual clustering and its relation to numerical taxonomy, in W.A. Gale (ed.), *Artificial Intelligence and Statistics*. Reading, MA: Addison-Wesley, pp. 77-116.
- Mechelin, I. Van, Michalski, R.S. (1993) General Introduction. In Mechelin, I.V., Hampton, J., Michalski, R.S., Theuns, P. (eds.) *Categories and Concepts: Theoretical Views and Inductive Data Analysis*, Academic Press, pp. 1-8.
- Medin, D.L., Ross, B.H. (1997) *Cognitive Psychology*. Harcourt Brace College Publishes.
- Medin, D.L., Wattenmaker, W.D., and Michalski, R.S. (1987) Constraints and preferences in inductive learning: an experimental study of human and machine performance. *Cognitive Science*, 11: pp. 299-339.
- Michalski, R.S., and Stepp, R.E. (1983) Automatic construction of classifications: conceptual clustering versus numerical taxonomy. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 5: pp. 396-410.
- Murphy, G.L. (1993) Theories and Concept Formation. In Mechelin, I.V., Hampton, J., Michalski, R.S., Theuns, P. (eds.) *Categories and Concepts: Theoretical Views and Inductive Data Analysis*, Academic Press, pp. 173-201.
- Rips, L.J. (1989) Similarity, typicality and categorization. In S. Vosniadou & A. Ortony (Eds.) *Similarity and analogical reasoning*. Cambridge: Cambridge University Press, pp. 21-59.
- Rosch, E.H. (1978) Principles of categorization. *Cognition and categorization* (Rosch E.H. and Lloyd B.B. (eds.)). Hillside, N.J.: Lawrence Erlbaum Associates: 27-48.
- Tryon, R. C. (1939). *Cluster Analysis*. Ann Arbor, MI: Edwards Brothers.
- Wittgenstein, L. (1953) *Philosophical investigations*. New York: Macmillan.