

УДК 510.22

**ИЗВЛЕЧЕНИЕ ЗНАНИЙ
ДЛЯ ОЦЕНКИ КРЕДИТОСПОСОБНОСТИ:
ПОДХОД ТЕОРИИ МУЛЬТИМНОЖЕСТВ***

А.Б.Петровский¹

Мультимножество или множество с повторяющимися элементами служит удобной математической моделью для представления объектов, которые характеризуются многими разнородными (количественными и качественными), в том числе противоречивыми признаками. В работе предложен новый метод классификации таких многопризнаковых объектов как точек метрического пространства мультимножеств и использован для оценки кредитоспособности владельцев кредитных карт на основе анализа их финансового поведения.

Введение

Оплата товаров и услуг с помощью кредитных карт стало одной из обыденных и удобных возможностей современного мира. Вместе с тем банки и кредитные компании, выпускающие кредитные карты, ежегодно несут многомиллионные убытки из-за перерасхода денежных средств, допущенного владельцами кредитных карт. Поэтому оценка возможной кредитоспособности лица, желающего получить кредитную карту, является одной из актуальных и типичных задач в банковском деле. На практике для принятия решения о выдаче или отказе в выдаче кредитной карты желательно иметь одно или несколько достаточно простых решающих пра-

* Работа выполнена при поддержке Российского фонда фундаментальных исследований (проекты 02-01-01077, 04-01-00290), программ фундаментальных исследований РАН «Математическое моделирование и интеллектуальные системы», «Фундаментальные основы информационных технологий и систем», гранта НШ1964.2003.1 Президента Российской Федерации по поддержке ведущих научных школ.

¹ 117312, Москва, проспект 60-летия Октября, 9, ИСА РАН, pab@isa.ru

вил, которые позволяли бы отнести заявителя к соответствующей категории в зависимости от представленных им персональных данных.

К настоящему времени в базах данных кредитных организаций накоплены большие объемы информации о реальном финансовом поведении десятков и сотен тысяч владельцев кредитных карт. Каждый владелец карты описывается в базе данных набором многих разнородных признаков, включающих личные сведения (пол, возраст, образование, семейное положение, род занятий, доходы, место жительства и прочее) и финансовые показатели (баланс счета, получение и возврат кредитов, оплата покупок и услуг, выдача наличных и тому подобное). Анализ этой информации дает возможность выявить закономерности финансового поведения владельцев кредитных карт и построить их классификацию по платежеспособности, в простейшем случае разбить всех владельцев на две условные категории: «хорошие», не допускающие перерасхода средств или своевременно погашающие возникающие задолженности, и «плохие», совсем не возвращающие кредиты или возвращающие их с большой задержкой.

Большая размерность анализируемых массивов данных (десятки и сотни тысяч владельцев, несколько десятков признаков) и значительная разнородность признаков, включающих количественные и качественные (порядковые, номинальные) переменные, создает определенные трудности построения решающих правил для оценки кредитоспособности предполагаемых владельцев кредитных карт. Для решения задачи классификации таких многопризнаковых объектов разработан ряд методов, использующих линейную и логистическую регрессию, деревья решений, нейронные сети, линейное программирование, кластерный анализ [Shi *et al*, 2001]. Однако в основе этих подходов лежит тот или иной вид «обучения» или «настройки» алгоритмов классификации на некотором предварительно отобранном массиве данных. Вместе с тем, очевидно, что среди многих тысяч владельцев карт окажутся такие, которые имеют одинаковые наборы личных сведений, но разные финансовые показатели, то есть, относятся к разным категориям. Другими словами, эти объекты описываются противоречивыми признаками, что будет приводить к ошибкам при их классификации. Кроме того, такие объекты могут присутствовать и в выборочном массиве данных, что существенно снижает качество обучения классифицирующих алгоритмов.

Мультимножество или множество с повторяющимися элементами служит удобной математической моделью для описания объектов, которые характеризуются многими разнородными (количественными и качественными) признаками и могут существовать в нескольких экземплярах с отличающимися, в частности, противоречивыми значениями признаков, свертка которых или невозможна, или математически некорректна. Крат-

ность элементов – существенная особенность мультимножества, позволяющая отличать его от множества и рассматривать мультимножество как качественно новое математическое понятие. Множественность и повторяемость факторов, описывающих объекты, усложняет и затрудняет решение задачи их классификации.

В работе предложен новый метод классификации совокупности многопризнаковых объектов, основанный на их представлении в виде точек метрического пространства мультимножеств. Метод позволяет строить обобщенное решающее правило для сортировки объектов, выраженное через значения их признаков. Метод классификации может работать с реальными большими массивами объектов, имеющих любые наборы признаков, в том числе противоречивых, и не требует предварительной настройки на выборочном массиве.

1. Мультимножества и пространства мультимножеств

Мультимножеством A , порожденным обычным множеством $U = \{x_1, x_2, \dots\}$, все элементы которого различны, называется совокупность групп элементов вида $A = \{k_A(x) \bullet x \mid x \in U, k_A(x) \in \mathbb{Z}_+\}$. Здесь $k_A: U \rightarrow \mathbb{Z}_+$ называется функцией числа экземпляров мультимножества, определяющей кратность вхождения элемента $x_i \in U$ в мультимножество A , что обозначено символом \bullet . Если $k_A(x) = \chi_A(x)$, где $\chi_A(x) = 1$ при $x \in A$ и $\chi_A(x) = 0$ при $x \notin A$, то мультимножество A становится обычным множеством. Если все мультимножества семейства $A = \{A_1, A_2, \dots\}$ образуются из элементов множества G , то G называется доменом для семейства A , а множество $\text{Supp}A = \{x \mid x \in G, \chi_{\text{Supp}A}(x) = \chi_A(x)\}$ – опорным множеством или носителем мультимножества A . Мощность мультимножества $|A| = \sum_x k_A(x)$ определяется как общее число экземпляров всех его элементов; размерность мультимножества $/A/ = \sum_x \chi_A(x) = |\text{Supp}A|$ – как общее число различных элементов. Мультимножество называется пустым \emptyset , если $k_{\emptyset}(x) = 0$, и максимальным Z , если $k_Z(x) = \max_{A \in A} k_A(x), \forall x \in U$ [Петровский, 1995, 2003].

Введены следующие основные операции над мультимножествами:
объединение

$$A \cup B = \{k_{A \cup B}(x) \bullet x \mid k_{A \cup B}(x) = \max(k_A(x), k_B(x))\};$$

пересечение

$$A \cap B = \{k_{A \cap B}(x) \bullet x \mid k_{A \cap B}(x) = \min(k_A(x), k_B(x))\};$$

арифметическое сложение

$$A + B = \{k_{A+B}(x) \bullet x \mid k_{A+B}(x) = k_A(x) + k_B(x)\};$$

арифметическое вычитание

$$A - B = \{k_{A-B}(x) \bullet x \mid k_{A-B}(x) = k_A(x) - k_{A \cap B}(x)\};$$

симметрическая разность

$$A\Delta B = \{k_{A\Delta B}(x) \bullet x \mid k_{A\Delta B}(x) = |k_A(x) - k_B(x)|\};$$

дополнение

$$\bar{A} = Z - A = \{k_{\bar{A}}(x) \bullet x \mid k_{\bar{A}}(x) = k_Z(x) - k_A(x)\};$$

умножение на число

$$h \bullet A = \{k_{h \bullet A}(x) \bullet x \mid k_{h \bullet A}(x) = h \cdot k_A(x), h \in \mathbb{Z}_+\};$$

арифметическое умножение

$$A \bullet B = \{k_{A \bullet B}(x) \bullet x \mid k_{A \bullet B}(x) = k_A(x) \cdot k_B(x)\};$$

арифметическая n -ая степень

$$A^n = \{k_{A^n}(x) \bullet x \mid k_{A^n}(x) = (k_A(x))^n\};$$

прямое произведение

$$A \times B = \{k_{A \times B} \bullet \langle x_i, x_j \rangle \mid k_{A \times B} = k_A(x_i) \cdot k_B(x_j), x_i \in A, x_j \in B\};$$

прямая n -ая степень

$$(\times A)^n = \{k_{(\times A)^n} \bullet \langle x_1, \dots, x_n \rangle \mid k_{(\times A)^n} = \prod_{i=1}^n k_A(x_i), x_i \in A\}.$$

Носители операций над мультимножествами удовлетворяют следующим соотношениям:

$$\text{Supp}(A \cup B) = \text{Supp}(A + B) = (\text{Supp} A) \cup (\text{Supp} B);$$

$$\text{Supp}(A \cap B) = \text{Supp}(A \bullet B) = (\text{Supp} A) \cap (\text{Supp} B);$$

$$\text{Supp}(A \Delta B) = (\text{Supp}(A - B)) \cup (\text{Supp}(B - A));$$

$$\text{Supp}(h \bullet A) = \text{Supp}(A^n) = \text{Supp} A;$$

$$\text{Supp}(A \times B) = (\text{Supp} A) \times (\text{Supp} B).$$

Семейство мультимножеств, замкнутое относительно операций объединения, пересечения, сложения и дополнения, называется алгеброй мультимножеств $L(\mathbf{Z})$, где максимальное мультимножество \mathbf{Z} является единицей алгебры, а пустое мультимножество \emptyset – нулем. Действительная неотрицательная функция $m(A)$, определенная на алгебре $L(\mathbf{Z})$ и удовлетворяющая условию сильной аддитивности: $m(A+B) = m(A) + m(B)$, называется мерой мультимножества. Мера мультимножества $m(A)$ обладает следующими свойствами: $m(\emptyset) = 0$; слабая аддитивность $m(A \cup B) = m(A) + m(B)$, $A \cap B = \emptyset$; монотонность $m(A) \leq m(B) \Leftrightarrow A \subseteq B$; симметричность $m(A) + m(\bar{A}) = m(\mathbf{Z})$; непрерывность $\lim_{i \rightarrow \infty} m(A_i) = m(\lim_{i \rightarrow \infty} A_i)$; эластичность $m(h \bullet A) = hm(A)$. Мере мультимножества можно определить различными способами, например, как линейную комбинацию функций кратности: $m(A) = \sum_j w_j k_A(x_j)$, $w_j > 0$. Заметим, что мощность мультимножества $|A|$ также будет мерой мультимножества

Метрические пространства мультимножеств (A, d) были введены в [Петровский, 1995, 2003], где определены следующие виды расстояний между мультимножествами:

$d_1(A, B) = m(A \Delta B)$; $d_2(A, B) = m(A \Delta B)/m(Z)$; $d_3(A, B) = m(A \Delta B)/m(A \cup B)$. (1)
 Функции $d_2(A, B)$ и $d_3(A, B)$ удовлетворяют условию нормировки $0 \leq d(A, B) \leq 1$. По определению принимается $d_3(\emptyset, \emptyset) = 0$. Основное расстояние $d_1(A, B)$ является метрикой типа Хемминга, традиционно используемым во многих приложениях. Полностью усредненное расстояние $d_2(A, B)$ характеризует различие между двумя мультимножествами A и B , отнесенное к расстоянию, максимально возможному в исходном пространстве. Локально усредненное расстояние $d_3(A, B)$ задает различие, отнесенное к максимально возможной «общей части» только этих двух мультимножеств в исходном пространстве.

2. Решающее правило классификации

Пусть $A = \{A_1, \dots, A_k\}$ – совокупность объектов (владельцев кредитных карт), которые описываются m дискретными признаками Q_1, \dots, Q_m , имеющими количественные и качественные (номинальные, либо порядковые) значения $q_s^{e_s}$, $e_s = 1, \dots, h_s$, $s = 1, \dots, m$. Порядковые значения качественного признака предполагаются упорядоченными от лучшего значения к худшему $q_s^1 > q_s^2 > \dots > q_s^{h_s}$. Объекты A_i , $i = 1, \dots, k$ предварительно рассортированы по нескольким классам X_t , $t = 1, \dots, f$ путем прямой классификации. Принадлежность объекта A_i к некоторому классу X_t выражается правилом сортировки R , которое может считаться еще одним качественным признаком со значениями r_t . Признаки Q_1, \dots, Q_m выражают персональные сведения владельца карты, а R характеризует его кредитоспособность. Для описания каждого объекта используется только одно какое-то значение признака из каждой группы Q_1, \dots, Q_m, R . Других дополнительных предположений об особенностях классов, признаков объектов и их значений (важности, предпочтительности, характерности и прочее) не делается. Требуется построить одно или несколько решающих правил, составленных из небольшого числа значений признаков, которые относили бы объекты к заданным классам наилучшим (в смысле близости к предварительной сортировке) образом.

Каждый объект A_i , $i = 1, \dots, k$ из совокупности A обычно представляется m -мерным вектором $q_i = (q_{i1}^{e_1}, q_{i2}^{e_2}, \dots, q_{im}^{e_m})$ в пространстве $Q = Q_1 \times Q_2 \times \dots \times Q_m$, являющемся прямым произведением шкал значений признаков Q_s , и ищутся некоторые решающие правила для отнесения объекта к заданному классу X_t . Сопоставим теперь объект A_i с мультимножеством вида

$A_i = \{(k_{A_i}(q_s^{e_s}) \bullet q_s^{e_s}), (k_{A_i}(r_i) \bullet r_i)\}$ над доменом $G = \{Q_1, \dots, Q_m, R\}$. Для сокращения размерности мультимножества ряд признаков, особенно имеющих большое число значений, удобно агрегировать в небольшие группы. Например, вместо возраста владельца кредитной карты можно указывать возрастную группу (например, до 20 лет, 21-30 лет, 31-40 лет, и так далее), вместо точной величины дохода – некоторый интервал, вместо точного адреса – район или город, и тому подобное.

Представление объекта A_i мультимножеством A_i может трактоваться как способ выражения индивидуального правила сортировки вида:

$$\text{ЕСЛИ } \langle \text{условия} \rangle, \text{ ТО } \langle \text{решение} \rangle. \quad (2)$$

Терм $\langle \text{условия} \rangle$ ассоциируется с различными комбинациями значений признаков $q_s^{e_s}$, описывающими свойства объекта A_i , а терм $\langle \text{решение} \rangle$ – с принадлежностью объекта A_i к классу X_i .

Для простоты будем считать, что результатом классификации должно быть разложение совокупности объектов A только на два класса X_a и X_b (в нашем случае – на категории «хороших» или «плохих» владельцев карт). Требование бинарной декомпозиции совокупности $A = \{X_a, X_b\}$ не является принципиальным ограничением. Если необходимо рассортировать объекты на большее число классов, можно сначала разбить совокупность объектов на две группы, затем одну из них или обе группы – на подгруппы, и так далее.

Рассмотрим наиболее простой и типичный случай агрегирования многопризнаковых объектов, когда каждая группа объектов формируется как сумма соответствующих им мультимножеств. Тогда каждое из мультимножеств $X_t = \{(k_{X_t}(q_s^{e_s}) \bullet q_s^{e_s}), (k_{X_t}(r_i) \bullet r_i)\}$, $t=a, b$, представляющее свой класс объектов X_t , можно записать в виде следующего разложения на мультимножества по группам признаков:

$$X_t = \sum_{s=1}^m Q_{st} + R_t,$$

где каждое слагаемое есть в свою очередь разложение

$$Q_{st} = \sum_{e_s=1}^{h_s} Q_{st}^{e_s}, \quad Q_{st}^{e_s} = \sum_{i \in I_{st}^{e_s}} A_i, \quad R_t = \sum_{i \in I_{rt}} A_i$$

Здесь $I_{st}^{e_s}$ и I_{rt} – подмножества индексов i для объектов A_i , имеющих соответственно $k_{A_i}(q_s^{e_s}) \neq 0$ и $k_{A_i}(r_i) \neq 0$, $A_i \subset X_t$.

Очевидно, что объекты A_i , которые попали в разложение $\{R_a, R_b\}$, сделанное только по индивидуальным правилам сортировки, образуют наилучшую из всех возможных декомпозиций рассматриваемой совокупности объектов $A = \{A_1, \dots, A_k\}$ на два класса для имеющегося набора правил сортировки. Рассмотрим теперь метрическое пространство мультимножеств $(A,$

d) с метрикой d , определяемой одним из выражений (1), и обозначим через $d^*=d(\mathbf{R}_a, \mathbf{R}_b)$ расстояние между мультимножествами \mathbf{R}_a и \mathbf{R}_b . В каждом конкретном случае расстояние d^* является предельно возможным расстоянием между объектами, входящими в разные классы.

Задача поиска обобщенного решающего правила классификации многопризнаковых объектов сводится к m оптимизационным задачам для каждой группы признаков Q_s

$$d(Q_{sa}, Q_{sb}) \rightarrow \max d(Q_{sa}, Q_{sb}) = d(Q_{sa}^*, Q_{sb}^*), \quad (3)$$

где мультимножества Q_{sa}^* и Q_{sb}^* расположены на максимально возможном расстоянии в метрическом пространстве мультимножеств (A, d) и принадлежат к разным классам. Решение каждой из задач (3) является наилучшей бинарной декомпозицией $\{Q_{sa}^*, Q_{sb}^*\}$ имеющейся совокупности многопризнаковых объектов $A = \{A_1, \dots, A_k\}$ по s -ой группе признаков. Когда число h_s значений $q_s^{e_s}$ каждого из признаков невелико, решение задачи (3) не вызывает существенных трудностей и может быть получено даже путем перебора.

Каждое мультимножество Q_{st}^* ($t=a, b$), относящееся к одному и тому же классу, представляется как сумма двух подмультимножеств $Q_{st}^* = Q_{st}^{*1} + Q_{st}^{*2}$. Значение признака q_s^* , которое определяет границу между сгенерированными слагаемыми Q_{st}^{*1} и Q_{st}^{*2} внутри каждой пары, назовем граничным. Комбинация граничных значений признаков $\{q_s^*\}$ из разных групп признаков Q_s , задает условия отнесения объекта A_i к соответствующему классу X_i и образует искомое обобщенное правило классификации совокупности многопризнаковых объектов вида (2).

Граничные признаки q_s^* можно упорядочить по величине расстояния $d(Q_{sa}^*, Q_{sb}^*)$. Для построения обобщенных правил классификации следует использовать признаки q_s^* , занимающие первые места в такой ранжировке. Чем ближе значение $d(Q_{sa}^*, Q_{sb}^*)$ к расстоянию $d^*=d(\mathbf{R}_a, \mathbf{R}_b)$, тем более точной будет аппроксимация первоначальной индивидуальной сортировки объектов. Оценить точность аппроксимации по s -ой группе признаков можно, например, выражением

$$\rho_s = d(Q_{sa}^*, Q_{sb}^*)/d(\mathbf{R}_a, \mathbf{R}_b).$$

В обобщенное решающее правило следует тогда включать значения признаков q_s^* , имеющие показатель точности ρ_s , превышающей некоторый желаемый пороговый уровень ρ_0 . Заметим, что величина ρ_s показателя точности аппроксимации характеризует своего рода относительную важность s -ой группы признаков Q_s в обобщенном решающем правиле.

После построения обобщенных решающих правил, относящих объекты к каждому из заданных классов, необходимо провести содержательный анализ правил с целью выявления в них одинаковых комбинаций гранич-

ных значений признаков, которые (в зависимости от контекста решаемой задачи) нужно либо оставить, либо исключить из правил.

Заключение

Классификация объектов, которые описываются многими количественными и качественными признаками, причем каждый из объектов может существовать в нескольких различающихся «экземплярах», является достаточно трудной проблемой. Эти трудности имеют и содержательные основания (например, некорректность применения процедур «усреднения» качественных признаков), и формальные причины (например, большая размерность задачи). В работе предложен новый метод классификации многопризнаковых объектов, основанный на их представлении с помощью мультимножеств, который не содержит необоснованных преобразований исходной информации и не приводит к потере или искажению данных.

Аналогичный предложенному подход к построению обобщенного решающего правила для классификации объектов, описываемых качественными признаками, был проверен на результатах конкурсного отбора проектов по государственной научно-технической программе по высокотемпературной сверхпроводимости [Ларичев и др., 1989]. Было сформулировано несколько решающих правил, одно из которых полностью совпало с примененным на практике [Петровский, 2003]. Обобщенное решающее правило классификации объектов позволило также выделить наиболее важные для отбора проектов критерии.

Список литературы

- [Ларичев и др., 1989] Ларичев О.И., Прохоров А.С., Петровский А.Б., Стернин М.Ю., Шепелев Г.И. Опыт планирования фундаментальных исследований на конкурсной основе. // Вестник АН СССР, 1989, №7, 51-61.
- [Петровский, 1995] Петровский А.Б. Метрические пространства мультимножеств. // Доклады Академии наук. 1995, Т.344, №2, 175-177.
- [Петровский, 2003] Петровский А.Б. Упорядочение и классификация объектов с противоречивыми признаками. // Новости искусственного интеллекта. 2003, №4, 34-43.
- [Петровский, 2003] Петровский А.Б. Пространства множеств и мультимножеств. – М.: Едиториал УРСС, 2003.
- [Shi *et al*, 2001] Shi Y., Wise M., Luo M., Lin Y. Data mining in credit card portfolio management: a multiple criteria decision making approach. // M.Koksalan, S.Zionts (Eds.). Multiple Criteria Decision Making in New Millennium. – Springer-Verlag, Berlin, 2001, 427-436.

Петровский А. Б. Извлечение знаний для оценки кредитоспособности: подход теории мультимножеств // Труды Девятой национальной конференции по искусственному интеллекту с международным участием (КИИ-2004).— Т. 2. — М.: Физматлит, 2004.— С. 853–860.

```
@InProceedings{Petrovsky_2004c,  
  author = "Петровский, А. Б.",  
  title = "Извлечение знаний для оценки кредитоспособности: подход  
теории мультимножеств",  
  booktitle = "Труды Девятой национальной конференции по искусственному  
интеллекту с международным участием (КИИ-2004)",  
  volume = "2",  
  address = "М.",  
  publisher = "Физматлит",  
  year = "2004",  
  pages = "853--860",  
  language = "russian",  
}
```