

**МОДЕЛЬ ОЦЕНКИ КРЕДИТОСПОСОБНОСТИ
ВЛАДЕЛЬЦЕВ КРЕДИТНЫХ КАРТ ПО ПРОТИВОРЕЧИВЫМ ДАННЫМ *****ABSTRACT**

The new model of credit card portfolio management decisions is suggested in this paper. Cardholders are represented in the model as multi-attribute objects with possibly contradictory values of attributes. Decision rules for an assessment of potential cardholder credibility are based on classifying multi-attribute objects in the multiset metric space.

ВВЕДЕНИЕ

Банки и кредитные компании, выпускающие кредитные карты, ежегодно несут многомиллионные убытки из-за перерасхода денежных средств, допущенного владельцами кредитных карт. Чтобы уменьшить объемы потерь, банки, принимая решение о выдаче или отказе в выдаче кредитной карты, стремятся оценить возможную кредитоспособность заявителя, основываясь на его 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 [1]. Однако в основе многих из этих подходов лежит тот или иной вид «обучения» или «настройки» алгоритмов классификации на

* Работа частично поддержана программами фундаментальных исследований РАН «Математическое моделирование и интеллектуальные системы», «Фундаментальные основы информационных технологий и систем», Российским фондом фундаментальных исследований (проекты 02-01-01077, 04-01-00290), грантом НШ1964.2003.1 Президента Российской Федерации по поддержке ведущих научных школ.

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 находиться и в выборочном массиве данных, что существенно снижает качество обучения классифицирующих алгоритмов.

В работе предложен новый метод классификации совокупности многопризнаковых объектов, основанный на их представлении в виде точек метрического пространства мультимножеств. 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. Основная идея метода классификации заключается в следующем. Все исходные мультимножества суммируются в два мультимножества, представляющие классы «хороших» и «плохих» объектов. Мультимножества-суммы в свою очередь разбиваются на несколько мультимножеств-слагаемых по числу признаков, характеризующих объекты. По каждой группе признаков для каждой пары слагаемых мультимножеств генерируется пара новых мультимножеств, максимально удаленных друг от друга в метрическом пространстве. Граница между новыми слагаемыми в каждой паре определяется некоторым значением соответствующего признака. Различные комбинации таких «граничных» значений признаков дают обобщенные решающие правила для классификации объектов. Величина взаимной удаленности пар сгенерированных мультимножеств в метрическом пространстве различна, а значит, будет различаться и точность сортировки объектов с помощью обобщенных решающих правил, составленных из разных «граничных» значений признаков. В итоговое решающее правило следует включить небольшое число признаков, обеспечивающих желаемый уровень точности классификации объектов. Предложенный метод классификации владельцев кредитных карт может работать с реальными массивами объектов, имеющих любые наборы признаков, в том числе противоречивых, и не требует предварительной настройки на выборочном массиве.

МУЛЬТИМНОЖЕСТВА И ОПЕРАЦИИ С НИМИ

Мультимножеством A , порожденным обычным множеством $U = \{x_1, x_2, \dots\}$, все элементы которого различны, называется совокупность групп элементов вида $A = \{k_A(x) \cdot x \mid x \in U, k_A(x) \in \mathbb{Z}_+\}$. Здесь $k_A: U \rightarrow \mathbb{Z}_+ = \{0, 1, 2, \dots\}$ называется функцией числа экземпляров мультимножества, определяющей кратность вхождения элемента $x_i \in U$ в мультимножество A , что обозначено символом \cdot . 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 $|A| = \sum_x k_A(x)$ опре-

деляется 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 $k_Z(x) = \max_{A \in \mathcal{A}} k_A(x)$, $\forall x \in U$ [2-5].

Введем следующие основные операции с мультимножествами:

| | |
|--------------------------------|--|
| объединение | $A \cup B = \{k_{A \cup B}(x) \cdot x \mid k_{A \cup B}(x) = \max(k_A(x), k_B(x))\};$ |
| пересечение | $A \cap B = \{k_{A \cap B}(x) \cdot x \mid k_{A \cap B}(x) = \min(k_A(x), k_B(x))\};$ |
| арифметическое сложение | $A + B = \{k_{A+B}(x) \cdot x \mid k_{A+B}(x) = k_A(x) + k_B(x)\};$ |
| арифметическое вычитание | $A - B = \{k_{A-B}(x) \cdot x \mid k_{A-B}(x) = k_A(x) - k_B(x)\};$ |
| симметрическая разность | $A \Delta B = \{k_{A \Delta B}(x) \cdot x \mid k_{A \Delta B}(x) = k_A(x) - k_B(x) \};$ |
| дополнение | $\bar{A} = Z - A = \{k_{\bar{A}}(x) \cdot x \mid k_{\bar{A}}(x) = k_Z(x) - k_A(x)\};$ |
| умножение на число | $h \cdot A = \{k_{h \cdot A}(x) \cdot x \mid k_{h \cdot A}(x) = h \cdot k_A(x), h \in \mathbb{Z}_+\};$ |
| арифметическое умножение | $A \cdot B = \{k_{A \cdot B}(x) \cdot x \mid k_{A \cdot B}(x) = k_A(x) \cdot k_B(x)\};$ |
| арифметическая n -ая степень | $A^n = \{k_{A^n}(x) \cdot x \mid k_{A^n}(x) = (k_A(x))^n\};$ |
| прямое произведение | $A \times B = \{k_{A \times B} \langle x_i, x_j \rangle \mid k_{A \times B} = k_A(x_i) \cdot k_B(x_j), x_i \in A, x_j \in B\};$ |
| прямая n -ая степень | $(\times A)^n = \{k_{(\times A)^n} \langle x_1, \dots, x_n \rangle \mid k_{(\times A)^n} = \prod_{i=1}^n k_A(x_i), x_i \in A\}.$ |

Носители операций с мультимножествами определяются следующими выражениями:

$$\begin{aligned} \text{Supp}(A \cup B) &= \text{Supp}(A + B) = (\text{Supp}A) \cup (\text{Supp}B); \\ \text{Supp}(A \cap B) &= \text{Supp}(A \cdot B) = (\text{Supp}A) \cap (\text{Supp}B); \\ \text{Supp}(A \Delta B) &= (\text{Supp}(A - B)) \cup (\text{Supp}(B - A)); \\ \text{Supp}(h \cdot A) &= \text{Supp}(A^n) = \text{Supp}A; \quad \text{Supp}(A \times B) = (\text{Supp}A) \times (\text{Supp}B). \end{aligned}$$

Семейство мультимножеств, замкнутое относительно операций объединения, пересечения, сложения и дополнения, называется алгеброй мультимножеств $L(\mathbf{Z})$, где максимальное мультимножество 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: $m(A) + m(B) = m(A + B)$, называется мерой мультимножества. Мера мультимножества $m(A)$ обладает следующими свойствами: $m(\emptyset) = 0$; монотонность $m(A) \leq m(B) \Leftrightarrow A \subseteq B$; симметричность $m(A) + m(\bar{A}) = m(Z)$; непрерывность $\lim_{i \rightarrow \infty} m(A_i) = m(\lim_{i \rightarrow \infty} A_i)$; эластичность $m(h \cdot A) = hm(A)$. Мету 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) введены в [2,5], где определены следующие виды расстояний между мультимножествами:

$$d_{1p}(A, B) = [m(A \Delta B)]^{1/p}; \quad d_{2p}(A, B) = [m(A \Delta B) / m(Z)]^{1/p}; \quad d_{3p}(A, B) = [m(A \Delta B) / m(A \cup B)]^{1/p}, \quad (1)$$

p – целое. Функции $d_{2p}(A, B)$ и $d_{3p}(A, B)$ удовлетворяют условию нормировки $0 \leq d(A, B) \leq 1$. По определению принимается $d_{3p}(\emptyset, \emptyset) = 0$. Основное расстояние $d_{1p}(A, B)$ является мет-

рикой типа Хемминга, традиционно используемым во многих приложениях. Полностью усредненное расстояние $d_{2p}(A, B)$ характеризует различие между двумя мультимножествами A и B , отнесенное к расстоянию, максимально возможному в исходном пространстве. Локально усредненное расстояние $d_{3p}(A, B)$ задает различие, отнесенное к максимально возможной «общей части» только этих двух мультимножеств в исходном пространстве.

МЕТОД КЛАССИФИКАЦИИ МНОГОПРИЗНАКОВЫХ ОБЪЕКТОВ

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 $Q_s = \{q_s^{e_s}\}$ отражают персональные данные владельца, а $R = \{r_t\}$ характеризует степень его кредитоспособности. 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 $\{q_s^{e_s}\}$ значений признаков, и ищутся некоторые решающие правила для отнесения объекта к заданному классу X_t . Сопоставим теперь объект A_i с мультимножеством вида

$$A_i = \{k_{Ai}(q_1^1) \bullet q_1^1, \dots, k_{Ai}(q_1^{h_1}) \bullet q_1^{h_1}, \dots, k_{Ai}(q_m^1) \bullet q_m^1, \dots, k_{Ai}(q_m^{h_m}) \bullet q_m^{h_m}, k_{Ai}(r_1) \bullet r_1, \dots, k_{Ai}(r_f) \bullet r_f\}$$

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 доходов, 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 ⟨условия⟩ ассоциируется с различными комбинациями значений признаков $q_s^{e_s}$, описывающими свойства объекта A_i , а терм ⟨решение⟩ – с принадлежностью объекта A_i к классу X_t . 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.

Рассмотрим наиболее простой и типичный случай агрегирования многопризнаковых объектов, когда каждая группа объектов формируется путем сложения соответствующих им мультимножеств. Тогда мультимножество

$X_t = \{k_{X_t}(q_1^1) \cdot q_1^1, \dots, k_{X_t}(q_1^{h_1}) \cdot q_1^{h_1}, \dots, k_{X_t}(q_m^1) \cdot q_m^1, \dots, k_{X_t}(q_m^{h_m}) \cdot q_m^{h_m}, k_{X_t}(r_a) \cdot r_a, k_{X_t}(r_b) \cdot r_b\}$, представляющее свой класс объектов, можно записать как сумму мультимножеств

$$X_t = \sum_{i \in I_t} A_i.$$

Здесь I_t – подмножество индексов i для объектов класса X_t , $t=a,b$, $I_a \cup I_b = \{1, \dots, k\}$, $I_a \cap I_b = \emptyset$. Перепишем мультимножество X_t в виде следующего разложения на новые мультимножества:

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

где отдельные слагаемые суть «однопризнаковые» мультимножества

$$Q_{st} = \sum_{i \in I_t} A_{iqs}, \quad A_{iqs} = \{k_{A_i}(q_s^1) \cdot q_s^1, \dots, k_{A_i}(q_s^{h_s}) \cdot q_s^{h_s}\};$$

$$R_t = \sum_{i \in I_t} A_{ir}, \quad A_{ir} = \{k_{A_i}(r_a) \cdot r_a, k_{A_i}(r_b) \cdot r_b\}.$$

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 $d^* = d(R_a, R_b)$ расстояние между мультимножествами R_a и 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

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

Требуется найти в каждой группе признаков Q_s новые мультимножества Q_{sa}^* и Q_{sb}^* , которые 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 $d^*=d(R_a, R_b)$, тем более точной будет аппроксимация первоначальной индивидуальной сортировки объектов. Оценить точность аппроксимации по s -ой группе признаков можно, например, выражением

$$\rho_s = d(Q_{sa}^*, Q_{sb}^*)/d(R_a, R_b).$$

В обобщенное решающее правило следует включать граничные значения признаков q_s^* , имеющие показатель точности ρ_s , превышающий 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 классификации.

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 подход к построению обобщенного решающего правила был протестирован на результатах конкурсного отбора проектов для государственной научно-технической программы по высокотемпературной сверхпроводимости [3]. Было сформулировано несколько решающих правил, одно из которых полностью совпало с примененным на практике. Обобщенное решающее правило классификации объектов позволило также выделить наиболее важные для отбора проектов критерии.

ЛИТЕРАТУРА

1. 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.
2. Петровский А.Б. Метрические пространства мультимножеств.//Доклады Академии наук. 1995, Т.344, №2, 175-177.

3. Petrovsky A. Method for approximation of diverse individual sorting rules. // Informatica. 2001, V.12, №1, 109-118.
 4. Петровский А.Б. Мультимножества как модель представления многопризнаковых объектов в принятии решений и распознавании образов. // Искусственный интеллект. 2002, №2, 236-243.
 5. Петровский А.Б. Пространства множеств и мультимножеств. – М.: Едиториал УРСС, 2003.
-

Петровский А. Б. Модель оценки кредитоспособности владельцев кредитных карт по противоречивыми данным // Искусственный интеллект. — Т. 2.— Донецк, Украина: Наука і освіта, 2004.— С. 155–161.

```
@InProceedings{Petrovsky_2004b,  
  author =      "Петровский, А. Б.",  
  title =      "Модель оценки кредитоспособности владельцев кредитных  
                карт по противоречивыми данным ",  
  booktitle =  "Искусственный интеллект",  
  volume =    "2",  
  address =    "Донецк, Украина",  
  publisher =  "Наука і освіта",  
  year =      "2004",  
  pages =     "155--161",  
  language =   "russian",  
}
```