

MULTI-ATTRIBUTE CLASSIFICATION OF CREDIT CARDHOLDERS: MULTISET APPROACH

Alexey B. Petrovsky¹

Institute for System Analysis
Russian Academy of Sciences

Prospect 60 Let Otyabrya 9, Moscow 117312, Russia

[mailto: pab@isa.ru](mailto:pab@isa.ru)

Keywords: credit cardholders, multi-attribute objects, multiset metric space, classification, decision rules

Summary: *The new model of constructing decisions for a credit card portfolio management is suggested. Cardholders are represented in the model as multi-attribute objects with possibly contradictory values of quantitative and qualitative attributes. Decision rules for an evaluation of potential cardholder credibility are based on classifying such multi-attribute objects in the multiset metric space.*

1. Introduction

Every year banks and credit card companies lose millions of dollars because of excess expenditures of credit cardholders. In order to diminish such losses banks try to predict a financial behavior of potential cardholder. The good cardholders reimburse credits in time, the bad cardholders forget to return money and contribute to bankruptcy.

Decision rules for a multiple criteria classification of credit cardholders' credibility can be produced by data mining of large-size files that are collected today in banks and describe real-life financial histories of credit cardholders. In these databases each cardholder is represented with a set of manifold attributes, which include personal data (sex, age, residence, occupation, affiliation, and so on), and financial data (income, balance, payments, purchase, cash, credit rating, and others). A large dimension of data files and variety of attributes (numerical and verbal, ordinal and nominal) cause difficulties of constructing decision rules for a classification of cardholders' credibility.

A number of classification techniques has been developed in order to construct the credit card portfolio management decisions. These approaches use the linear and logistic regression, linear programming, decision trees, neural networks, cluster analysis, and other techniques (ref. see Shi *et al.*, 2001). The most of them are based on learning classification algorithms. Obviously, there exist persons among thousands of real cardholders who have the same or very closed personal attributes but the opposite financial behavior. It means, these persons belong to the different categories, thus they have contradictory descriptions. A presence of such inconsistencies in data samples may deteriorate a quality of learning procedures, and leads to classification errors and poor robustness of results.

¹ This work is supported by the Scientific Programs "Mathematical Modeling and Intellectual Systems", "Fundamentals of Information Technology and Systems" of the Russian Academy of Sciences; the projects 02-01-01077, 04-01-00290 of the Russian Foundation for Basic Research; the grant 1964.2003.1 of the President of the Russian Federation for a support of the prominent scientific schools.

A new approach to classifying credit cardholders in terms of their attributes is suggested in this paper. The main idea of a classification algorithm is the following. All real-life cardholders represented as multisets are aggregated in the given classes according to their credit ratings, which are determined beforehand by one or several experts. Multiset-sums that correspond to the given classes of cardholders are decomposed on the several one-attribute multiset-terms. The pairs of new submultisets are generated for every one-attribute multiset-terms. These submultisets are to be placed at the maximal distance in the metric space of multisets. Combination of attributes, which defines boundaries between the generated submultisets within one-attribute categories, produces the general decision rule. The method proposed for classifying multi-attribute objects operates, generally, with an arbitrary-size data file, including contradictory data, without a previous adjustment on a sampling collection.

2. Presentation of Credit Cardholders

Let $\mathbb{A} = \{A_1, \dots, A_k\}$ be a collection of k objects (credit cardholders), which are described with many discrete attributes (personal and financial data) Q_1, Q_2, \dots, Q_m . Each attribute Q_s has quantitative or qualitative (ordinal or nominal) values $\{q_s^{e_s}\}$, $e_s = 1, \dots, h_s$, $s = 1, \dots, m$. Ordinal values of attribute are supposed to be ordered from the best to the worst $q_s^1 \succ q_s^2 \succ \dots \succ q_s^{h_s}$. All objects (cardholders) A_1, \dots, A_k are related previously to some classes (credit ratings) X_1, \dots, X_f by one or several experts. The sorting rule is an additional qualitative attribute R with gradations $\{r_t\}$, $t = 1, \dots, f$, and shows that an object A_i belongs to the class X_t . When many experts estimate credit ratings of cardholders, individual expert decision rules r_t may be similar, diverse, or contradictory.

Multi-attribute objects A_i , $i = 1, \dots, k$ are presented usually as vectors $q_i = (q_{i1}^{e_1}, \dots, q_{im}^{e_m})$ in the space $Q = Q_1 \times \dots \times Q_m$, which is a direct product of attribute values $q_s^{e_s}$. The classification rules coincide with discriminant surfaces in this space.

Let us present now a multi-attribute objects A_i as a multiset (A1)

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

drawn from the domain of attributes $G = \{Q_1, \dots, Q_m, R\}$. In order to shorten a multiset dimensionality, some attributes, especially which has a large number of values, may be aggregated in small groups. For instance, the real cardholder age may be replaced by an age group (less than 20 years, 21-30 years, 31-40 years, and so on), the income of cardholder – by some interval of incomes, the cardholder address – by a residence region, and the like.

For simplicity, assume that the collection of objects $\mathbb{A} = \{A_1, \dots, A_k\}$ is to be sorted only in two classes X_a and X_b (good and bad credit cardholders). Consider the most simple and typical case of a multi-attribute objects' aggregation, when every class of objects is formed as an addition of corresponding multisets (Petrovsky, 2003a). In this case, all properties of all members of the group X_t are aggregated. Then the multiset

$$X_t = \{k_{X_t}(q_1^1) \bullet q_1^1, \dots, k_{X_t}(q_1^{h_1}) \bullet q_1^{h_1}, \dots, k_{X_t}(q_m^1) \bullet q_m^1, \dots, k_{X_t}(q_m^{h_m}) \bullet q_m^{h_m}, k_{X_t}(r_a) \bullet r_a, k_{X_t}(r_b) \bullet r_b\}$$

that corresponds to the object class X_t is a sum of multisets

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

Here I_t is a subset of indexes i for the object of the category X_t , $t = a, b$, $I_a \cup I_b = \{1, \dots, k\}$, $I_a \cap I_b = \emptyset$. Let us rewrite the multiset X_t as a decomposition of new multisets:

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

where terms are the following one-attribute multisets

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

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

The demand to decompose the collection of objects only in two categories is not the principle limitation.

The collection of multi-attribute objects $A=\{A_1,\dots,A_k\}$ and a set of their attributes $G=\{Q_1,\dots,Q_m,R\}=\{g_j\}$ may be expressed also by matrix $C=||k_{A_i}(g_j)||$. Here k is a number of objects, $m+1$ is a number of attributes groups, $h=h_1+\dots+h_m+f$ is a number of object attributes, $j=1,\dots,h$. The matrix C is often used in data analysis, pattern recognition and called the “object-attribute” table, information table or decision table. In our case, the decision table C presents information on properties of the multi-attribute objects A_i and their membership to a certain decision class. This table C has a dimension $k\times h$ and consists of $m+1$ boxes, which correspond to attributes Q_1,\dots,Q_m,R . The reduced decision table $C^*=||k_{X_i}(g_j)||$ has a dimension $2\times h$ and consists of two rows $k_{X_a}(g_j), k_{X_b}(g_j)$, which correspond to the classes X_a and X_b .

3. Algorithm of Classification

The object description in the form of multiset A_i (1) is another presentation of the sorting rule:

$$\text{IF } \langle \text{conditions} \rangle, \text{ THEN } \langle \text{decision} \rangle. \quad (2)$$

The term $\langle \text{conditions} \rangle$ corresponds to the various combinations of attribute values $k_{A_i}(q_s^e)$; the term $\langle \text{decision} \rangle$ is associated with the individual expert decisions. Need to generate the simple general decision rule, which would coincide maximally with a large set of individual (and may be contradictory) sorting rules, include a small number of attributes and assign objects to the given classes with the accuracy admitted.

The scheme of an algorithm, which approximates of a large family of contradictory sorting rules with a compact general decision rule, looks as follows. Consider a multiset metric space (A, d) with one of the distances (A2)-(A4). Obviously, objects A_i corresponding to the multiset decomposition $R=\{R_a, R_b\}$ is the best partition of the object collection $A=\{A_i\}$ into the categories X_a and X_b for the given set of individual sorting rules. Thus the distance $d^*=d(R_a, R_b)$ between multisets R_a and R_b is maximal in the space (A, d) . In the case of ideal classification (without inconsistencies of individual expert sorting rules), the maximal distance is equal correspondingly to $d_{1p}^*=[kn]^{1/p}$, $d_{2p}^*=[1/h]^{1/p}$, or $d_{3p}^*=1$. Here n is a number of experts.

The problem of an approximation of the diverse rules for sorting a collection of multi-attribute objects is transformed into the following m optimization problems:

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

So, it is necessary to find in every attribute group Q_s the new multisets Q_{sa}^* and Q_{sb}^* , which are placed at the maximal distance from each other in the metric space (A, d) and belong to the different classes. The solution of every optimization problem (3) is the best binary decomposition $Q_s^*=\{Q_{sa}^*, Q_{sb}^*\}$ of the object collection referred to the corresponding s -th attribute group.

Every multisets Q_{st}^* , $t=a,b$, is represented as the sum of two submultisets $Q_{st}^*=Q_{st}^{*1}+Q_{st}^{*2}$. A value q_s^* of the attribute Q_s that defines a boundary between the submultisets Q_{st}^{*1} and Q_{st}^{*2} is said to be a boundary value. A combination of attribute boundary values $\{q_s^*\}$ for various attribute groups forms a decision rule of type (2) for classifying the multi-attribute object collection A_1,\dots,A_k .

The attribute boundary values q_s^* may be ordered according to values of distances $d(Q_{sa}^*, Q_{sb}^*)$. The attributes q_s^* , which occupy first places in this ranking, are to be included in the general decision rule. The nearer the distances $d(Q_{sa}^*, Q_{sb}^*)$ to the maximal distance $d^*=d(R_a, R_b)$, the more accurate is the approximation of individual sorting rules. The accuracy of sorting rules approximation can be estimated by the rate

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

The attribute boundary values q_s^* that provide the required accuracy of approximation $ac_s \geq ac_0$ are to be included in the term $\langle \text{conditions} \rangle$ of the general decision rule (2) for an object classification. Note that the value of accuracy rate ac_s characterizes also a relative importance of the s -th attribute group Q_s within the general decision rule.

4. Conclusion

A classification of objects, that are described by many diverse (quantitative and qualitative) attributes, and may exist in several copies with different, in particular, contradictory values of attributes, is a quite difficult problem. These difficulties have substantial grounds (for example, an incorrectness of “averaging” qualitative attributes), and some formal reasons (for instance, a large dimension of problem). In this paper, the new method is proposed for classifying a collection of objects, which are represented as multisets in metric spaces. This technique does not contain ungrounded transformations of data (like “averaging”, “mixing”, “weighting” attributes, and so on), and may be especially fruitful in the case of inconsistencies of objects’ properties and individual sorting rules.

The analogous approach to constructing a general classification rule had been tested on the individual expert decisions related to a competitive selection of research projects (Petrovsky, 2001). One of the found general rules coincided completely with the real-life decision rule (Larichev *et al.*, 1989).

5. Appendix (Multisets and Multiset Metric Spaces)

A multiset (also called a bag) is a known notion that is used in combinatorial mathematics and other fields. A multiset, or a set with repeating elements is the very convenient mathematical model to present and analyze a collection of objects that are described with many quantitative and qualitative attributes, and can exist in several copies with various values of attributes. Review briefly the theory of multisets and metric spaces of multisets (Knuth, 1969, Yager, 1986, Petrovsky, 1994,2001,2003). A multiset A drawn from an ordinary (crisp) set $U=\{x_1, x_2, \dots, x_j, \dots\}$ with different elements is defined as the following collection of elements’ groups

$$A = \{k_A(x_1) \bullet x_1, \dots, k_A(x_j) \bullet x_j, \dots\} = \{(k_A(x) \bullet x) | x \in U, k_A(x) \in \mathbb{Z}_+\}. \quad (A1)$$

Here $k_A: U \rightarrow \mathbb{Z}_+ = \{0, 1, 2, \dots\}$ is called a counting function of multiset, which defines the number of times that the element $x_i \in U$ occurs in the multiset A , and this is indicated with the symbol \bullet . A multiset A becomes an ordinary set when $k_A(x) = \chi_A(x)$, where $\chi_A(x) = 1$, if $x \in A$, and $\chi_A(x) = 0$, if $x \notin A$.

If all multisets of the family $A = \{A_1, A_2, \dots\}$ are composed from the elements of set G then G is said to be a generic domain for a family A . A crisp set $\text{Supp}A = \{x | x \in G, \chi_{\text{Supp}A}(x) = \chi_A(x)\}$ is named a support set or carrier of the multiset A . The multiset cardinality $|A| = \sum_x k_A(x)$ is defined as a total number of all copies of its elements, and the multiset dimensionality $/A/ = \sum_x \chi_A(x) = |\text{Supp}A|$ is defined as a total number of different elements. The maximal value of the counting function $\text{hgt}A = \max_{x \in G} k_A(x)$ is called a height of the multiset A , and an element $x_{A*} = \arg \max_{x \in G} k_A(x)$ is called a peak of A . The multiset is named the empty multiset \emptyset , if $n_{\emptyset}(x) = 0$, the maximal multiset Z , if $n_Z(x) = \max_{A \in A} k_A(x)$, and the constant multiset $C_{[h]}$, if $k_{C_{[h]}}(x) = h = \text{const}, \forall x \in G$. So the empty multiset \emptyset is the constant multiset of the height $\text{hgt}\emptyset = 0$, and an ordinary set B is the constant multiset of the height $\text{hgt}B = 1$.

Define the following operations with multisets:

union of multisets

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

intersection of multisets

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

arithmetic addition of multisets

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

arithmetic subtraction of multisets

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

complement of multiset

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

symmetric difference of multisets

$A\Delta B = \{k_{A\Delta B}(x) \bullet x \mid k_{A\Delta B}(x) = |k_A(x) - k_B(x)|\};$
multiplication of multiset by a scalar (a reproduction of multiset)

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

$A \bullet B = \{k_{A \bullet B}(x) \bullet x \mid k_{A \bullet B}(x) = k_A(x) \cdot k_B(x)\};$
raising to an arithmetic power

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

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

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

where $\langle x_1, \dots, x_n \rangle$ is a cortege of n elements.

The support sets of operations with multisets are defined as follows:

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

In general, the operations of arithmetic addition, multiplication by a scalar, arithmetic multiplication, and raising to an arithmetic power are not defined in the theory of sets. These operations may be analogous respectively to the operations with vectors $\mathbf{a} + \mathbf{b} = (a_1 + b_1, \dots, a_n + b_n)$, $h \bullet \mathbf{a} = (ha_1, \dots, ha_n)$, and matrixes $A + B = \|a_{ij} + b_{ij}\|_{m \times n}$, $h \bullet A = \|h \cdot a_{ij}\|_{m \times n}$, $A \bullet B = \|a_{ij} \cdot b_{ij}\|_{m \times n}$. The last operation is different from the traditional matrix multiplication. The operation of multiset selection suggested by Yager (1986) is a special case of multiset arithmetic multiplication, where one of the factors is an ordinary set. Note that a product of multiset A and scalar h may be also presented as a sum of h multisets A : $h \bullet A = A + \dots + A$, or as a product of constant multiset $C_{[h]}$ and multiset A : $h \bullet A = C_{[h]} \bullet A$. When multisets are reduced to sets, the operations of arithmetic multiplication and raising to an arithmetic power degenerate into a set intersection, but the operations of set arithmetic addition and set multiplication by a scalar will be impracticable.

A family of multisets, which is closed under the operations of union, intersection, addition and complement, is said to be an algebra $L(\mathbf{Z})$ of multisets, where the maximal multiset \mathbf{Z} is the unit and the empty multiset \emptyset is the zero of the algebra. A measure m of multiset A is a real-valued function, which is defined on the algebra $L(\mathbf{Z})$, and has the following properties: $m(A) \geq 0$, $m(\emptyset) = 0$; strong additivity $m(\sum_i A_i) = \sum_i m(A_i)$; weak additivity $m(\cup_i A_i) = \sum_i m(A_i)$ for $A_i \cap A_j = \emptyset$; weak monotony $m(A) \leq m(B) \Leftrightarrow A \subseteq B$; symmetry $m(A) + m(\bar{A}) = m(\mathbf{Z})$; continuity $\lim_{i \rightarrow \infty} m(A_i) = m(\lim_{i \rightarrow \infty} A_i)$; elasticity $m(h \bullet A) = hm(A)$. The multiset measure may be written in the various ways, for instance, as a linear combination of counting functions $m(A) = \sum_j w_j k_A(x_j)$, $w_j > 0$. Remark that the multiset cardinality $|A|$ is also a measure of multiset.

Different metric spaces of multisets (A, d) introduced by Petrovsky (1994, 2001, 2003) have the following distances between multisets:

$$d_{1p}(A, B) = [m(A \Delta B)]^{1/p}, \quad (\text{A2})$$

$$d_{2p}(A, B) = [m(A \Delta B)/m(\mathbf{Z})]^{1/p}, \quad (\text{A3})$$

$$d_{3p}(A, B) = [m(A \Delta B)/m(A \cup B)]^{1/p}, \quad (\text{A4})$$

where $p > 0$ is integer. Functions $d_{2p}(A, B)$ and $d_{3p}(A, B)$ satisfy the normalization condition $0 \leq d(A, B) \leq 1$. Note that due to the continuity of the multiset measure, the distance $d_{3p}(A, B)$ is not defined for $A = B = \emptyset$. So $d_{3p}(\emptyset, \emptyset) = 0$ by the definition. The metric d_{1p} is a Hamming-type distance that is traditional for many applications. The metric d_{2p} characterizes a difference between two multisets related to the common properties of all others, and the metric d_{3p} reflects a difference related to properties of only both multisets.

Various peculiarities of the multiset metric spaces are considered and discussed in (Petrovsky, 2003b).

References

Knuth, D.E. (1969) *The Art of Computer Programming. Vol.2. Seminumerical Algorithms*, Reading: Addison-Wesley.

Larichev, O.I., Prokhorov, A.S., Petrovsky, A.B., Sternin, M.Yu., Shepelev, G.I. (1989) The Experience of Planning of the Basic Research on the Competitive Base. *Vestnik of the USSR Academy of Sciences*, 7, 51-61 (in Russian).

Petrovsky, A.B. (1994) An Axiomatic Approach to Metrization of Multiset Space. In: Tzeng, G.H., Wang, H.F., Wen, U.P., Yu, P.L., editors, *Multiple Criteria Decision Making*, New York:Springer-Verlag, 129-140.

Petrovsky, A.B. (2001) Multiattribute Sorting of Qualitative Objects in Multiset Spaces. In: Koksalan M., Zionts S., editors, *Multiple Criteria Decision Making in the New Millennium. Lecture Notes in Economics and Mathematical Systems*, N507, Berlin: Springer-Verlag, 124-131.

Petrovsky, A.B. (2003a) Cluster Analysis in Multiset Spaces. In: Goldevsky M., Mayr H., editors, *Information Systems Technology and its Applications*, Bonn: Gesellschaft fur Informatik, 199-206.

Petrovsky, A.B. (2003b) *Spaces of Sets and Multisets*, Moscow: Editorial URSS (in Russian).

Shi Y., Wise M., Luo M., Lin Y. (2001) Data Mining in Credit Card Portfolio Management: a Multiple Criteria Decision Making Approach. In: Koksalan M., Zionts S., editors, *Multiple Criteria Decision Making in the New Millennium. Lecture Notes in Economics and Mathematical Systems*, N507, Berlin: Springer-Verlag, 427-436.

Yager, R.R. (1986) On the Theory of Bags. *International Journal of General Systems*, 13, 23-37.

Petrovsky Alexey B. Multi-Attribute Classification Of Credit Cardholders: Multiset Approach // Electronic Proceedings of the 17-th International Conference Multiple Criteria Decision Making (August 6-11).— Whistler, Canada: 2004.

```
@InProceedings{Petrovsky_2004d,  
  author = "Petrovsky, {\relax Alexey B}.",  
  title = "Multi-Attribute Classification Of Credit  
  Cardholders: Multiset Approach",  
  booktitle = "Electronic Proceedings of the 17-th International  
  Conference Multiple Criteria Decision Making (August 6-11)",  
  volume = "",  
  address = "Whistler, Canada",  
  publisher = "",  
  year = "2004",  
  pages = "",  
}
```