

Usage of Positional Representations in Tasks of Revealing Logical Regularities

Yury Laptin¹, Alexander Vinogradov² and Yury Zhuravlev²

¹ Glushkov Institute of Cybernetics of the Ukrainian National Academy of Sciences
Academician Glushkov pr. 40, 03680 Kiev, Ukraine

² Dorodnicyn Computing Centre of the Russian Academy of Sciences
Vavilov str. 40, 119333 Moscow, Russian Federation

Abstract. A new approach to the problem of description of clusters in the form of sets of logical regularities is developed. Special attention is paid to detection of steady interrelations in data which are valid for considerable shares of the sample. The approach is based on usage of fast bit operations and positional method of data representation. New criterion of adjacency is proposed for high-level points of the representation, and it's used further in the process of assembling such points in maximal hyper-parallelepipeds corresponding to the best regularities. The method can be used as preliminary step in various tasks related to the search of essential logical regularities and substantial interpretation of data.

1 Introduction

Numerous approaches to the problem of extraction of knowledge from numerical data sample are used recently. Special attention is paid to methods of detection of steady interrelations in data which are valid for considerable shares of the sample. So, in the area of IP there are widely used methods where some hierarchical division of set of objects into blocks of the increasing size is set in advance. At various levels of hierarchy interrelations of specific types appear presented and can be processed by specialized procedures. It is supposed that at top levels the most important laws can be found out and used further. Among such approaches it is necessary to note various transformations of type coarse-to-fine, various scale-spaces, representations with multiple resolution, space filling curves, quad - and octo-trees [1] [2]. Positional representation is among them, too [3].

At all advantages of the mentioned methods, it is possible to note some obvious shortcomings. Use of concrete hierarchy is rather a severe constraint which not always corresponds to the nature of the data. In some cases shortcomings are connected also with spatial restrictions having origin from borders of an image and from the direction of basic axes. As it is known, transformation and conformity problems between different representations describing the same object in the form of quad- or octo-trees, are among the main obstacles to wider use of these representations. Many of noted lacks

are overcome in the approaches based on the use of language of *Logical Regularities* (LR) for the describing clusters in recognition tasks [4] [5]. In this work some possibilities of use of positional data representation for search of description of the classes and the decision rule in terms of LR are considered.

2 Logical Regularities and Positional Coordinates

Standard statement of a problem of recognition in R^N for K classes is considered. The exact decision of a problem results in a segmentation of all the space on K subareas. Let $X \subset R^N$ is one of segments corresponding to the class k . Logical regularity LR_j is represented by a conjunction of kind $\bigwedge R_i$ where each condition R_i is a couple of inequalities $A_i < X_i < B_i$. To conjunction R_i there corresponds a hyper-parallelepiped in X . Some disjunction $\bigvee LR_j$ can provide a covering of segment of the class k . The covering can possess those or other properties depending on the choice of A_i, B_i . In formal representation here the usual logic is used, and the form of condition R_i has simple intuitive sense. As a result of automatic search of representation $\bigvee LR_j$ there is a list of LR that is sufficient for a class. Logical regularities are presented in a form that is clear to the human acquisition, and they can contain important new information. The less number of indices i and j used in $\bigvee LR_j$, the bigger average share of sample covered by each hyper-parallelepiped, and the more important law it represents. Whatever the origin of exact solution of a problem of recognition, it is possible to try to transform it to appropriate representation in terms of LR and to use noted above advantage in substantial interpretation.

We will build this transformation by means of positional data representation. Positional representation implies setting some discrete grid $D^N \subset R^N$ that is common for all segments, where $|D| = 2^d$. For a grid points $x = (x_1, x_2, \dots, x_N)$ the conversion of the points in their positional representations corresponds to effectively carried out transformation on bit slices in D^N when each m -bit of binary representation of the number x_n results in $p(n)$ -bit of binary representation of the m -digit in 2^N -ary representation of value representing the vector as a whole. Here $m \leq d$, and function $p(n)$ defines some permutation on the set $\{1, 2, \dots, N\}$. As a result there is linearly ordered scale S of length 2^{dN} , representing one-to-one all points of the grid in the form of a curve filling the space D^N densely. We consider Z -scanning of a grid, when $p(n) = n$. At a choice of other types of scanning $p(n)$ it is possible to achieve more smooth filling of the grid, for example, in the form of Peano curve, etc.

3 Searching Maximal Hyper-parallelepipeds in Positional Representation of Segments

For a given digitization D^N , all segments $X_k \subset R^N$ corresponding to classes k are found as solutions of the task of recognition, and they set together a K -valued function f defined on the scale S . As it is known, m -bit in 2^N -ary positional representation corresponds to some n -dimensional cube of size $2^{(m-1)N}$. We will name such cube as a m -point. We will search for situations when there are homogeneous cubes as parts of a

segment of some class k which can be united in hyper-parallelepiped of bigger volume. We will see, how it can be made on the basis of research of the structure of function f .

Let's select at first all segments of scale S where value f does not vary. We will sequentially check points of S and build a list \hat{f} , containing records of sort $(b_1b_2\dots b_{d-m+1}, k)$. Here the binary sequence $(b_1b_2\dots b_{d-m+1})$ consists of $d - m + 1$ bits, and the last bit ends the positional representation of some m -point as whole. In such code we will represent in \hat{f} all detected cubes with homogeneous filling k .

Lemma 1. *At linear search of points S , the transition from one m -point to another occurs at zeroing of all low-order digits in 2^N -ary positional representation of a current point.*

The proof follows immediately from the way of construction of the list \hat{f} .

Obviously, the list \hat{f} is filled at one pass of the scale S . For such filling it is enough to watch the moments of simultaneous zeroing of all 2^N -ary digits, lower than m .

The received description of segments \hat{f} already has a structure of type $\bigwedge R_i$. However, this representation is not the best because it does not meet yet the requirements of the Section 1. In particular, it is redundant, because some n -dimensional cubes corresponding to different m -points of positional representation can be united in common hyper-parallelepiped, and they can receive therefore a representation in the form of a single logical regularity $\bigwedge R_i$. Two cubes of identical size we will name *adjacent* if they adjoin on the common $N - 1$ -dimensional edge. A $N - 1$ -dimensional edge, orthogonal to axes n , we will name as a *n-edge*. We will use the following criterion of association of pair of m -points belonging to the segment k .

Lemma 2. *Two m -points $C1, C2$ are adjacent on a n -edge iff: 1) there is a m' -point C such that $m' > m$; 2) record on $C1$ precedes record on $C2$ in \hat{f} ; 3) in binary record of all 2^N -ary digits $m, m + 1, \dots, m' - 1$ bits with number n have in record for $C1$ ($C2$) from the list \hat{f} values 1 (accordingly, 0); 4) all other bits in records for $C1, C2$ in \hat{f} coincide.*

The proof. If there is a digit $m' > m$ with the specified properties then volumes of cubes $C1, C2$ coincide, and in corresponding m' -cube C all m -points, in binary bits of 2^N -ary digits of which n -th bits can vary only, make a uniform hyper-parallelepiped of length $2^{(m'-1)}$. As Z-scanning is used, $C1$ takes a highest position in the $m' - 1$ -sub-cube, and $C2$ takes a lowest position in some other $m' - 1$ -sub-cube. But $C1$ precedes $C2$ in the list \hat{f} , and therefore, they are allocated in the middle of the hyper-parallelepiped and consequently are adjacent.

On the opposite. Any two sub-cubes of equal volume $C1, C2$ from D^N are allocated on S without intersections, therefore, one of them is placed in \hat{f} earlier. Let it be $C1$. If $C1, C2$ are adjacent on n -edge, there is minimum m' such that some m' -cube C contains $C1, C2$, but any $m' - 1$ -sub-cube does not contain them together. Therefore, each of them belongs to unique own $m' - 1$ -sub-cube. But $C1, C2$ are adjacent, and those should be own $m' - 1$ -sub-cubes, as any two sub-cubes of identical dimension in positional representation either are adjacent, or are not intersected. Repetition of this reasoning for all 2^N -ary digits up to m allows to conclude that all bits in binary representation of all 2^N -ary digits $m, m + 1, \dots, m' - 1$ in records for $C1, C2$ in \hat{f} coincide,

excepting bits with number n . It is obvious that change of value of n -th binary bit in any 2^N -ary digit from 0 to 1 results in moving corresponding sub-cube along the axis n (and backward, at change 1 on 0). For $C1$ within its own $m' - 1$ -sub-cube (and within any sub-cube of smaller size up to m) it is impossible to make such change. It means that in record for $C1, C2$ in the list \hat{f} all these bits have values 1 (accordingly, 0 for $C2$). The Lemma 2 is proved.

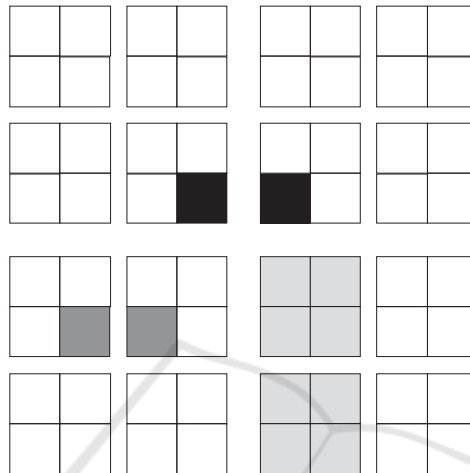


Fig. 1. An example in case $N = 2, K = 3, d = 3$. Adjacency of two gray 1-points ($k = 1$) is revealed for 1-edge by considering left-bottom 3-point; adjacency of two black 1-points ($k = 2$) is revealed for 1-edge by considering the whole grid (4-point); adjacency of two light-gray 2-points ($k = 3$) is revealed for 2-edge by considering right-bottom 3-point.

Using the given criterion, it is possible to build hyper-parallelepipeds of the increasing size. Gathering hyper-parallelepipeds of maximal volume results in exact description of segments of the sample in terms of LR and of sort $\sqrt{LR_j}$.

4 Discussion

It is necessary to remind that this construction depends essentially on the properties of the grid of digitization D^N . In particular, the value $D^N = 2^{dN} = 2^{32}$ can serve as natural estimate of length of the scale S that makes the approach applicable to data processing on PC. Though bit conversions are performed efficiently, the representation received on this way can differ considerably from the solutions found by more exact (and more slow) methods. The reason is that some features of boundaries of the segments cannot be represented adequately on a discrete grid. Nevertheless, if in the exact solution there is found some LR_j , the presented method reveals some logical regularity LR'_j as its correspondent on the grid, and vice-versa. On the basis of all told above and proved in Lemmas 1, 2 it is possible to assert that LR_j and LR'_j will differ from each other only in parameters A_i, B_i , and no more, than on the value of the step of digitization (and cuts of empty regions on boundaries of the grid D^N).

5 Conclusions

A new approach is developed for the task of mining new knowledge from the data sample in the form of logical regularities. The main attention is paid to the search of maximal logical regularities, that are valid for essential shares of the sample. The approach is based on usage of the fast bit operations providing conversion of data in positional representation and search of blocks in it. In such representation all sample elements receive hierarchically ordered positions on the common linear scale where the main relations of closeness are preserved. A new method is developed for description of the structure of classes in the form of restricted list where only important elements from various levels of hierarchy are present. The criterion of adjacency for pares of such elements is developed and used further in procedure of combining these elements into hyper-parallelepipeds of maximal size which correspond to the most essential logical regularities.

The approach can be used as preliminary step in solution of various tasks related to the search of logical regularities and substantial interpretation of data. The approach can find application in data processing and mining in 2D and 3D, and also in the search of logical regularities in abstract feature space of higher dimensions.

Acknowledgements

This work was done in the framework of Joint project of the National Academy of Sciences of Ukraine and the Russian Foundation for Basic Research No 08-01-90427 'Methods of automatic intellectual data analysis in tasks of recognition objects with complex relations'.

References

1. Mark de Berg, Marc van Kreveld, Mark Overmars, and Otfried Schwarzkopf: Computational Geometry (2nd revised ed.). Springer-Verlag (2000) 506 pages.
2. T. Lindeberg: Scale-Space Theory in Computer Vision. Kluwer Academic Publishers (1994) 440 pages.
3. Aleksandrov V.V., Gorskiy N.D.: Algoritmy i programmy strukturnogo metoda obrabotki dannykh. L. Nauka (1983) 208 str. (russian)
4. Yu.I.Zhuravlev, V.V.Ryazanov, O.V.Senko: RASPOZNAVANIE. Matematicheskie metody. Programmaya sistema. Prakticheskie primeneniya. Izdatelstvo "FAZIS", Moscow (2006) 168 str. (russian)
5. Ryazanov V.V.: Logicheskie zakonomernosti v zadachakh raspoznavaniya (parametricheskiy podkhod). Zhurnal vychislitelnoy matematiki i matematicheskoy fiziki, T. 47/10 (2007) str. 1793-1809. (russian)