# Using generalized precedents for big data sample compression at learning*

## *V. V. Ryazanov*[1], *A. P. Vinogradov*[1], *and Yu. P. Laptin*[2]

vngrccas@mail.ru

[1]Dorodnicyn Computing Centre of the Russian Academy of Sciences, 40 Vavilova st., Moscow, Russia
[2]Glushkov Institute of Cybernetics of the Ukrainian National Academy of Sciences,
40 Glushkova ave., Kiev, Ukraine

The role of intrinsic and introduced data structures at constructing efficient recognition algorithms is analyzed. The concept of generalized precedent as representation of stable local regularity in data and based on its use methods of reduction of the dimension of tasks has been investigated. Two new approaches to the problem based on positional data representation and on cluster means for elementary logical regularities are proposed. The results of computational experiment with data compression in parametric spaces for several practical tasks are presented.

**Keywords**: *generalized precedent; logical regularity; positional representation; bit slice; hypercube; correct decision rule*

## 1 Introduction

Methods of solving problems of recognition and data analysis use various structures in data. At a choice of appropriate structure, two main objectives are pursued:

a) identification natural clusters of density in the feature space in which the vectors of realizations are condensed; and

b) optimization of the computational expenses necessary for creation of the decision rule and subsequent calculations.

Both purposes are closely related with each other and at their realization compete for computational resources. For this reason, the overwhelming share of principles of structurization can be refereed to both directions simultaneously, and the choice of a concrete method in many respects is determined by assignment of priorities for (a) and (b). Now, the huge number of approaches, algorithms, and methods, more or less successful, are applied to achievement of both purposes. Note some survey publications on this subject [1, 2] where the most actual and perspective decisions are outlined. In them, both conceptual and technical aspects of the choice of the compromise are concerned.

In this paper, the close relationship which exists between the concepts 'precedent' and 'cluster' is investigated. The question how the mobility of the border between admissible realization of the concepts 'precedent' and 'cluster' in the computational environment can be used at the search of a compromise for (a) and (b) has been studied. Let the sum

$$F(x) = \sum_i \mu_i e^{-0.5(x_i - x)^\top \sigma (x_i - x)} \tag{1}$$

be parametrical approximation of empirical distribution by uniform normal mix with constant covariance matrix $\sigma^{-1}$. The component $\mathcal{N}(x_i, \sigma^{-1})$ represents compact spatial cluster $C_i$ with

the center $x_i$ which is unambiguously described by the couple $(x_i, \mu_i)$. The natural treatment (1) implies that each cluster of $C_i$ is filled by vectors corresponding to casual deviations from the parameters of the central object $x_i$. The recognized object $x_0$ can also be considered as a single realization of distribution of probable localizations of the true center which also form cluster $C_0$ with the center $x_0$ and with the same form of distribution $\mu_i e^{-0.5(\boldsymbol{x_i}-\boldsymbol{x})^\top \boldsymbol{\sigma}(\boldsymbol{x_i}-\boldsymbol{x})}$ where coordinates of the center $x_0$ and variable $x$ interchange positions according to the Bayes's law. Thereby, internally inherent structure of the sample gains simple representation; however, this simplicity is reached at the price of creation of representation (1) as a solution of hard multi-parametric inverse problem, and also with difficulties of reference of the cluster $C_0$ to one of the classes, each of which is represented by several clusters of type $C_i$. Certainly, the example is exaggerated, but it correctly reflects relationship between two concepts.

Opposite example in which injected structure of data appears, one can find in IP (Internet protocol) technologies where rigid hierarchy of clusters forcedly introduced into the $R^2$ plane in the form of quadtree provides high computational efficiency at training and recognition, but the hierarchy is thus invariable, and in orthodox approaches, it is not adjusted in any way to internal structure of the training sample [3]. The coordinates of clusters of quadtree are unambiguously fixed, and substantial information is coded by only the density of filling of clusters at different levels.

Further in this work, recognition problems will be considered in which the balance between the accuracy of representation and computational efficiency can be reached via structural reduction within the pair 'precedent–cluster.'

## 2 Generalized precedents: Feature space replacement

Application of models of type (1) assumes the use of Euclidean norm

$$\|x\| = \left(\sum_i x_i^2\right)^{1/2}$$

for estimation and comparison vectors in $R^N$. The norm binds together the values of different parameters, in particular, qualitatively incomparable ones. That is often convenient, but can cause questions at substantial interpretation of results.

On the contrary, for hierarchy of clusters in a quad- or oktree, the scales in different dimensions does not interact. In case of IP, it is the main drawback of quadtree-type models which limits their use [4]. Really, in case of images or scenes, it is usually assumed that spatial directions possess equal properties. At the same time, the models of this type are noninvariant to rotations in $R^2$ and $R^3$ and, therefore, the results achieved with their use are difficult to reproduce after rotation of the basis.

In abstract feature space, the assumption of 'equality axes in rights' takes place rarely. Moreover, invariance of a model to independent scaling of the main dimensions (in general, to independent nonlinear changes of scale on axes) becomes an important advantage. One of the successful approaches based on the use of this invariance is the approach with logical regularities [5–7]. In this approach, the clusters are hyperparallelepipeds in $R^N$, each cluster is described by conjunction of the following kind:

$$L^i = \&_n R_n^i, R_n^i = (A_n^i < x_n < B_n^i), \;\; n = 1, \ldots, N, \tag{2}$$

and substantively interpreted as a recurring joint manifestation of feature values $x_1, x_2, \ldots, x_N$ of the vector $x$ at intervals $A_n^i < x_n < B_n^i, n = 1, \ldots, N$. The principle of proximity to each other

precedents of the same phenomenon here is embodied in the requirement of filling the interior of a certain type of cluster by the objects of the same class. The shape of clusters becomes of particular importance, and multiple joint appearance of feature values at the selected intervals in this approach is seen as an independent phenomenon called *elementary logical regularity.*

In all approaches mentioned above, just limited number of parameters is used to describe the spatial arrangement of the cluster and its filling. In case of quadtree, each cluster is encoded by one integer and one real parameter $(q_i, \mu_i)$; for the normal mixture (1), it is a pair of kind $(x_i, \mu_i)$; in case of logical regularities, it is a set of $2N$ border markers on axes $A_n^i$ and $B_n^i$, $n = 1, \ldots, N$, and also, the weight of regularity $\mu_i$.

Recently, V. V. Ryazanov has proposed the idea of reduction of dimension of the problem through the use of substantial clusters such as hyperparallelepiped or component $\mathcal{N}(x_i, \sigma^{-1})$ with significant aprioristic weight as new training objects. Each combined object is regarded as geometric manifestation of some separate regularity in initial data and is called *generalized precedent*. Such generalized precedents are just proposed to use in training. Generalized precedents are described by geometric parameters of corresponding clusters and dimensions of the new feature spaces in the above examples are $2, N + 1$, and $2N + 1$, respectively. Thus, dimension of the space of generalized precedents may change as the upward and downward, but big training sample receives more compact representation as the result.

# 3 Examples of usage of generalized precedents for sample reduction

## 3.1 Positional representation

In case of positional data representation, structural elements belong also to the special family of logical regularities of the 1st type (2), when real numbers are truncated to real ones, and the intervals used $(A_n^i < x_n < B_n^i), n = 1, \ldots, N$, are equal in length. Thus, hyperparallelepipeds become hypercubes of restricted variety of kinds.

Positional notation is the development of quadtree model in dimensions greater than 2. The main advantage is that the structuring of positional hierarchy is already automatically injected into any numerical data when registering them, and it is immediately ready for use. It was also noted above that in models of this type, independent scaling of the main axes is naturally implemented, and this fact makes prospects of using the proposed approach in a variety of recognition problems, including the ones with incomparable numerical features.

Let finite sets $X_k$ are preset in $R^N$ and represent classes $k$, $k = 1, \ldots, K$, of the training sample $X$.

Positional representation [8] of data in $R^N$ is defined by a bit grid $D^N \subseteq R^N$ where $|D| = 2^d$ for some integer $d$.

The parameter $d$ is not fixed in advance. As it will be shown, its value is determined by the results of the analysis of the mutual arrangement of classes in the training sample.

Each grid point $x_1, x_2, \ldots, x_N$, $n = 1, \ldots, N$, corresponds to effectively performed transformation on bit slices in $D^N$, when the $m$th bit in binary representation $x_n \in D$ of the $n$th coordinate of $x$ becomes $p(n)$-bit of binary representation of the $m$th digit of $2^N$-ary number that represents vector $x$ as whole. Here, it is supposed $0 < m \leqslant d$, and function $p(n)$ defines a permutation on $1, \ldots, N, p \in S_N$. The result is a linearly ordered scale $S$ of length $2^{dN}$, representing one-to-one all the points of the grid in the form of a curve that fills the space $D^N$ densely. For chosen grid $D^N$, an exact solution of the problem of recognition with $K$ classes results in $K$-valued function $f$ defined on the scale $S$. As known, $m$-digit in $2^N$-ary positional representation corresponds to $n$-dimensional cube of volume $2^{N(m-1)}$. It is called $m$-point. For each $m$, the entire set of $m$-points is called $m$-slice. Thus, one has

**Lemma 1.** *There are just one d-point, $2^N$ ones of $(d-1)$-points, and $2^{dN}$ ones of 1-points on the scale $S$.*

Each of $m$-points, $0 < m \leqslant d$, can be regarded as separate cluster in $D^N$. If it is nonempty and filled with data of certain class only, one has got generalized precedent.

Further, for every $k, k = 1, \ldots, K$, and every $m, 0 < m \leqslant d$, let us look for the set of all of $m$-points, which are generalized precedents, i.e., elementary logical regularities of class $k$. The larger uniform regions in the domain of function $f$ (corresponding to generalized precedents as elder $m$-points), the better the decision rule. In the description of positional generalized precedent, the filled volume is represented latently by parameter $m$ (i.e., by the level in the hierarchy) and actual new feature space is formed of pair $(p_i, m_i)$.

Here, let describe the scheme of algorithm **A** that realizes this search on hierarchy of $m$-points of the grid $D^N$ from top to bottom.

The search is carried out for all classes $k, k = 1, \ldots, K$, simultaneously. Data of the training sample $X = \bigcup X_k \subseteq R^N$ are transformed into $2^N$-ary indices of the grid $D^N$.

All objects of the sample are processed in turn. Each next object $x$ marks all $m$-points, $m > 1$, of the own branch in hierarchy $D^N$ with the index $k$. Notice that for $m > 1$, there are no more than $\sum_{m=2}^{d} 2^{N(d-m)}$ different $m$-points. For dimensions $N > 3$, this number is negligible in comparison with the total number of 1-points of the grid $D^N$, and this fact provides the mechanism of compression of the sample.

Upon termination of search in each marked point of hierarchy $D^N$, the final attributing is carried out: if some $(m+1)$-point was marked with indexes of various classes (i.e., is not the generalized precedent), and all $m$-points subordinated to it are the generalized precedents, then the entire last are included in the decision rule. Further specification and attributing of subordinated $(m-1)$-points are not required.

As for all classes $k, k = 1, \ldots, K$, the analysis began with the same $d$-point as the top of hierarchy, one has

**Lemma 2.** *Algorithm **A** finds all generalized precedents of specified kind in the training sample $X = \bigcup X_k \subseteq R^N$.*

Since the number of $m$-points is final, any $m$-point that hashes classes at actual choice of the parameter $d$, will be further resolved by next iteration of algorithm **A** under this $m$-point regarded as new top and, thus, one has got

**Lemma 3.** *Iterative process on the basis of algorithm **A** provides creation of exact decision rule that is correct on the training sample $X = \bigcup X_k \subseteq R^N$.*

Thus, one has to decide what is better in this or that case: big $d$ or many iterations of **A**.

Since data of training sample $X = \bigcup X_k \subseteq R^N$ are analyzed consecutively, further retraining of any recognition algorithm constructed on this way will demand investigation of objects not more than inside one generalized precedent for each new object.

When sets of generalized precedents for all classes are built, one can combine within each class some collected $m$-points as hypercubes in larger hyperparallelepipeds according to criteria of contiguity [7]. So, one more way of building space of generalized precedents taking the form of elementary logical regularities of the 1st kind can be realized.

## 3.2 Cluster means as generalized precedents

Generally, large number of various ways of creation logical regularities of the 1st kind is developed now, and the choice of one of them is determined by properties of data and the character

of a problem of recognition or forecasting [5]. As the second example of use of generalized precedents, here, a new method of compression of data is considered which consists in transformation of feature space $R^N$ to the space $(c^i, \mu^i)$ with dimension $N+1$. Class means $c^i$ in clusters of regularities $L^i = \&_n R_n^i, R_n^i = (A_n^i < x_n < B_n^i), n = 1, \ldots, N$, are used as generalized precedents in this method. The space $(A_n^i, B_n^i, \mu^i), n = 1, \ldots, N$, itself is used thus at the intermediate stage.

Notice that in hierarchy of $m$-points of $D^N$, very rigid criterion of selection of generalized precedents was applied. Existence of the only object of alien class as a part of any $m$-point (hypercube of large volume, when $m$ is close to $d$) excludes the last from among the generalized precedents and strongly reduces thereby potential efficiency of compression of the sample. For this reason, in the majority of methods of creation of logical regularities, softer selection criteria are used when existence of certain share of objects of others classes as a part of this or that hyperparallelepiped (corresponding to elementary logical regularity) is allowed. Thus, flexibility of the model of logical regularities in general is reached and possibility of creation of simple decision rule with small set of elements of the sort $L^i = \&_n R_n^i, R_n^i = (A_n^i < x_n < B_n^i)$ is provided where each of them represents essential part of the training sample.

The proposed method of compression uses the specified opportunity fully, but realizes also a way of disposal of the difficulties related with the existence of alien objects in the cluster of regularity $L^i = \&_n R_n^i, R_n^i = (A_n^i < x_n < B_n^i)$. Let $x_t^L, t = 1, \ldots, T^L$, be a set of objects of the $k$th class as a part of cluster of elementary logical regularity $L$. Construct in the space $R^{N+1}$ a new sample that is made of vectors of averages $c^L = \sum_t^{T^L} x_t^L$, and their shares $T^L$ in each regularity $L$. Thereby, the space of the generalized precedents $(c^i, T^L)$ is set, each point of which corresponds to nonuniformly filled cluster of initial space $R^N$ where the objects of class $k$ dominate. The role of cluster geometry thus partially loses its value, important is only that the share $T^L$ of objects of the $k$th class within the cluster is big.

## 4 Reconstruction of the decision rule in initial feature space

Reconstruction of the decision rule in initial feature space consists in the return replacement of sets of the essential generalized precedents with clusters of the chosen for them geometric forms. Replacement is carried out directly and does not cause difficulties. In Fig. 1, it is shown how it takes place in case of hypercubes of positional data representation. For cluster means, this transition is even more direct since in this case, the space of generalized precedents differs from initial feature space only in additional equipment of weight coordinate $T^L$.
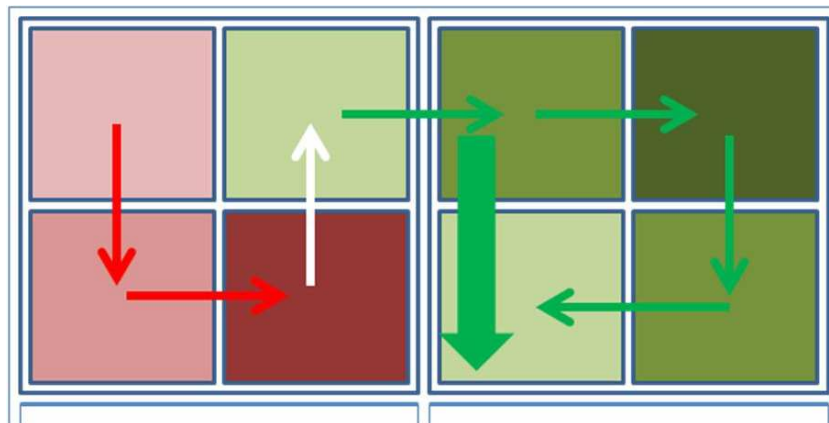
## 5 Computational experiment

Computational experiment in the framework of this training sample compression model was made for several types of generalized precedents on real tasks. The best accuracy was achieved by approach on the basis of cluster means. Here, the generalized precedents are used for representation of the training sample in the form of sets of new precedents that match as the source precedents and classes and the results of analysis of the initial training sample. As additional information for each class $K_\lambda, \lambda = 1, 2, \ldots, l$, there are used multiple logical regularities of classes $P_\lambda = \{P_t(\mathbf{x})\}$, i.e., predicates of the form
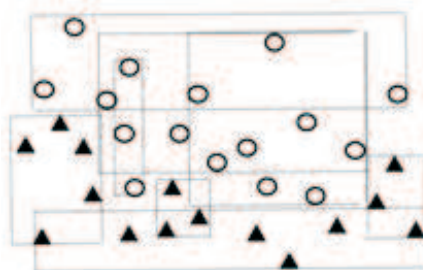
$$P^{\Omega_1, \mathbf{c}^1, \Omega_2, \mathbf{c}^2}(\mathbf{x}) = \bigwedge_{j \in \Omega_1} (c_j^1 \leqslant x_j) \bigwedge_{j \in \Omega_2} (c_j^1 \geqslant x_j),$$

$$\Omega_1, \Omega_2 \subseteq 1, 2, \ldots, n, \mathbf{c}^1, \mathbf{c}^2 \in R^n,$$

where

**Figure 1** Fragment of quadtree with the scale $S$ of Peano-type. Two classes (red and green) are separated in $S$ by white arrows. Color intensity depicts the density of filling. Left big square represents an $m$-point that hashes classes. Right square is filled with objects of the green class only detected in all subordinated $(m-1)$-points and so, this $m$-point represents generalized precedent as large uniform region included in the decision rule (big green arrow)



**Figure 2** Intervals of two-dimensional regularities of two classes are the marked boxes

1) $\exists \mathbf{x}_t \in K_2^0 | P^{\Omega_1, \mathbf{c}^1, \Omega_2, \mathbf{c}^2}(\mathbf{x}_t) = 1$;
2) $\forall \mathbf{x}_t \neg \in K_2^0 | P^{\Omega_1, \mathbf{c}^1, \Omega_2, \mathbf{c}^2}(\mathbf{x}_t) = 0$; and
3) $P^{\Omega_1, \mathbf{c}^1, \Omega_2, \mathbf{c}^2}(\mathbf{x}_t)$ represents a local optimum of the standard criterion of predicates' quality.

Here, through $\tilde{K}_\lambda$, the training sample objects from class $\lambda$ are designated. Two schemes of definition of generalized predicates are used.

In the first scheme, sets of objects that satisfy the predicates of $P_\lambda$ correspond to the set $\tilde{K}_\lambda$. Figure 2 shows a model example. An analog of the "nearest neighbor" algorithm was used.

Object $\mathbf{x}$ is assigned to the class, the regularity of which is considered the closest, the "distance" to the patterns is calculated by the formula:

$$d_\alpha(\mathbf{x}) = \frac{\sum_{\mathbf{x}_t : P^{\Omega_1, \mathbf{c}^1, \Omega_2, \mathbf{c}^2}(\mathbf{x}_t) = 1} \rho(\mathbf{x}, \mathbf{x}_t)}{|\{\mathbf{x}_t : P^{\Omega_1, \mathbf{c}^1, \Omega_2, \mathbf{c}^2}(\mathbf{x}_t) = 1\}|}$$

where $\rho$ is the Euclidean metric in $R^n$.

Comparison was carried out on the data of the credit scoring (2 classes, 15 features, 348 test objects) [8]. The accuracy of the standard and the modified method of "nearest neighbor" was on the test data, respectively, 75.6% and 77.5% of correct answers.

In the second scheme, generalized precedent is considered as the set of values of all logical regularities of the object, disjunction of their negations, the set of values of all logical regularities of another class, and disjunction of their negations (classification with 3 classes and more used

$(a)$



$(b)$

**Figure 3** Visualization of the original training sample ($a$) and of the sample of generalized precedents in the parametric space where classes become linearly separable ($b$)

scheme "one against all"). Thus, each object corresponds to a vector of numbers $\{0, 1\}$, and the generalized precedent is simply a description of the object in the new feature space. Figure 3 shows the visualization of the original training sample and the sample in new parametric space on the task of recognition of breast cancer [9]. The objects of different classes are presented in a plane gray and black circles. Generalized precedents of the training sample are linearly separable.

The version of support vector machine implemented in [5] was used as the main classification method. The results of the comparison of methods of recognition of test data on various tasks are presented in Table 1.

In general, the achieved positive results testify to prospects of the approach and to need of further development of this direction of researches.

## 6  Concluding remarks

The use of some inherent and injected structures in data has been considered. The opportunities arising from the use of generalized precedents for creation of detailed decision rule have been analyzed. It was shown that in case of positional data representation, the feature space $R^N$ can be reduced to two-dimensional space where training data become represented by compact clusters. Reduced representation realizes the one-dimensional scan of $R^N$, which is loaded with weights of generalized precedents. A scheme for an iterative process is proposed that yields to construct exact solutions which are correct on the training data. A new method of training

**Table 1** Results of comparison of recognition methods on various tasks

| Task | Classes | Dimension | Objects | Reference objects | Accuracy on reference objects | Accuracy on generalized precedents |
|------|---------|-----------|---------|-------------------|-------------------------------|-----------------------------------|
| "Breast" | 2 | 9 | 344 | 355 | 94.6 (0.8) | 96.1 |
| "Credit" | 2 | 15 | 342 | 348 | 80.5 (4.3) | 64.5 |
| "Image" | 7 | 16 | 210 | 2100 | 68.8 (27.7) | 92.0 (0.6) |

data compression has been developed and investigated based on the use of cluster means for elementary logical regularities and on its use as generalized precedents in transformed $(N+1)$-dimensional feature space. Computational experiment was made for several types of generalized precedents on real tasks. Good results approve the new opportunities and open prospects of the use of generalized precedents in recognition tasks with big data samples.

# References

[1] De Berg, M., M. van Kreveld, M. O. Overmars, and O. Schwarzkopf. 2000. *Computational geometry: Algorithms and applications*. 2nd ed. Springer. 291–306. doi: `http://dx.doi.org/10.1007/978-3-662-04245-8`

[2] Berman, J. 2013. *Principles of big data*. Elsevier. 1–14.

[3] Samet, H., and R. Webber. 1985. Storing a collection of polygons using quadtrees. *ACM Trans. Graph.* 4(3):182–222. doi: `http://dx.doi.org/10.1145/282957.282966`

[4] Eberhardt, H., V. Klumpp, and U. D. Hanebeck. 2010. Density trees for efficient nonlinear state estimation. *13th Conference (International) on Information Fusion Proceedings*. Edinburgh. 1–8. doi: `http://dx.doi.org/10.1109/ICIF.2010.5712086`

[5] Zhuravlev, Yu. I., V. V. Ryazanov, and O. V. Senko. 2006. *Raspoznavanie. Matematicheskie metody. Programmnaya sistema. Prakticheskie primeneniya*. Moscow: FAZIS. 168 p. (In Russian.)

[6] Ryazanov, V. V. 2007. Logicheskie zakonomernosti v zadachakh raspoznavaniya (parametricheskiy podkhod). Zhurnal vychislitelnoy matematiki i matematicheskoy fiziki 47(10):1793–1809. (In Russian.)

[7] Vinogradov, A., and Yu. Laptin. 2010. Usage of positional representation in tasks of revealing logical regularities. *VISIGRAPP-2010, Workshop IMTA-3 Proceedings*. Angers. 100–104.

[8] Aleksandrov, V. V., and N. D. Gorskiy. 1983. *Algoritmy i programmy strukturnogo metoda obrabotki dannykh*. Leningrad: Nauka. 208 p. (In Russian.)

[9] Mangasarian, O. L., and W. H. Wolberg. 1990. Cancer diagnosis via linear programming. *SIAM News* 23(5):1–18.

# Использование обобщенных прецедентов для сжатия больших выборок при обучении[*]

*В. В. Рязанов*[1], *А. П. Виноградов*[1], *Ю. П. Лаптин*[2]

vngrccas@mail.ru

[1]Вычислительный центр РАН им. А.А.Дородницына, Россия, г. Москва, ул. Вавилова, 40

[2]Институт кибернетики им. В. М. Глушкова Национальной академии наук Украины, Украина, г. Киев, пр. Ак. Глушкова, 40

Анализируется роль внутренне присущих и привнесенных структур данных при построении эффективных алгоритмов распознавания. Исследуется понятие обобщенного прецедента как способа представления устойчивой локальной закономерности в данных и методы снижения размерности задач на основе его использования. Предложены два новых подхода к проблеме, основанные на позиционном представлении и на средних по кластерам элементарных логических закономерностей. Представлены результаты вычислительного эксперимента по сжатию данных в параметрических пространствах для нескольких практических задач.

**Ключевые слова**: *обобщенный прецедент; логическая закономерность; позиционное представление; битовый слой; гиперкуб; корректное решающее правило*

## Литература

[1] *De Berg M., van Kreveld M., Overmars M. O. Schwarzkopf O.* Computational geometry: Algorithms and applications. — 2nd ed. — Springer, 2000. P. 291–306. doi: `http://dx.doi.org/10.1007/978-3-662-04245-8`

[2] *Berman J.* Principles of big data. — Elsevier, 2003. P. 1–14.

[3] *Samet H., Webber R.* Storing a collection of polygons using quadtrees // ACM Trans. Graph., 1985. Vol. 4. Iss. 3. P. 182–222. doi: `http://dx.doi.org/10.1145/282957.282966`

[4] *Eberhardt H., Klumpp V., Hanebeck U. D.* Density trees for efficient nonlinear state estimation // 13th Conference (International) on Information Fusion Proceedings. Edinburgh, 2010. doi: `http://dx.doi.org/10.1109/ICIF.2010.5712086`

[5] *Журавлев Ю. И., Рязанов В. В., Сенко О. В.* Распознавание. Математические методы. Программная система. Практические применения. — М.: Фазис, 2006. 168 с.

[6] *Рязанов В. В.* Логические закономерности в задачах распознавания (параметрический подход) // Ж. вычислительной математики и математической физики, 2007. Т. 47. № 10. С. 1793–1809.

[7] *Vinogradov A., Laptin Yu.* Usage of positional representation in tasks of revealing logical regularities // VISIGRAPP-2010, Workshop IMTA-3 Proceedings. Angers, 2010. P. 100–104.

[8] *Александров В. В., Горский Н. Д.* Алгоритмы и программы структурного метода обработки данных. — Л.: Наука, 1983. 208 с.

[9] *Mangasarian O. L., Wolberg W. H.* Cancer diagnosis via linear programming // SIAM News, 1990. Vol. 23. No. 5. P. 1–18.