

Algorithm for Predicting the Quality of the Product of Metallurgical Production

Damir N. Gainanov

Dmitriy A. Berenov

Ural Federal University
Gor'kogo St. 63, 620002 Ekaterinburg, Russia
damir.gainanov@gmail.com, berenov@dc.ru

Abstract

In this paper, the problem of the quality of the product of metallurgical production is investigated in the conditions when the reassignment can be organized in the process of realization of a specific technological route. The information on the completed technological routes forms a training sample for the pattern recognition problem with the teacher and the choice of the technological route for the continuation of the production process is carried out taking into account the expected quality indicators of the final product. To reduce the dimensionality of the problem, a given set of executed technological routes is divided into discrete classes, in each of which an algorithm for constructing a decision tree can be implemented. The paper gives a formal description of the developed algorithm for the node of the decision tree and an example of implementation.

Introduction

In work [Gainanov & Berenov, 2017] it has been developed the general approach for solving the problem of the quality of the product of metallurgical production which is based on the idea of realization of Big Data technologies. In the process of production activity there is a process of continuous accumulation of information about the completed technological routes, while the accumulated information has all the signs of big data, namely:

1. the accumulated information has significant volumes measured by many terabytes of information,
2. the accumulation of information occurs in the streaming mode at a high speed,
3. the accumulated information is characterized by a great variety and contains the values of several thousand and even tens of thousands of different parameters.

As a result of processing the obtained data the research reduces to the classical problem of pattern recognition in a geometric formulation, when it is necessary to construct a decision rule for assigning the input vector a_i to one of m classes. To solve the problem of pattern recognition in a geometric formulation, various algorithms can be applied, described, for example, in [Gainanov, 2014], [Gainanov, 2016], [Mazurov, 1990], [Mazurov, 2004]. Also, to solve the problem, an alternative cover algorithm from [Gainanov, 1992] can be applied.

Copyright © by the paper's authors. Copying permitted for private and academic purposes.

In: Yu. G. Evtushenko, M. Yu. Khachay, O. V. Khamisov, Yu. A. Kochetov, V.U. Malkova, M.A. Posypkin (eds.): Proceedings of the OPTIMA-2017 Conference, Petrovac, Montenegro, 02-Oct-2017, published at <http://ceur-ws.org>

In this paper, we present an algorithm for solving the problem of pattern recognition in a geometric formulation, the implementation of which determines the classification of the expected products. The algorithm is based on the principles of constructing logical decision trees and reducing the dimension of the problem by preliminary clustering the source data space.

1 Basic Definitions

Let $\mathcal{A} = \{A_1, \dots, A_n\}$ be the set of technological aggregates involved in production.

Definition 1.1 The directed graph $\vec{G} = (\mathcal{A}, E)$ with the set of vertices \mathcal{A} and the set of arcs $E \subseteq \mathcal{A}^2$ is called the infrastructural graph if $(A_1, A_2) \in E$ if and only if the output EP of the aggregate A_1 can serve as the input EP for the aggregate A_2 .

Definition 1.2 The technological route $P = (A_{i_1}, A_{i_2}, \dots, A_{i_k})$ is any directed path in \vec{G} .

The set of all technological routes $\mathcal{P} = \{P_1, \dots, P_k\}$ is the technological base of the production under consideration. We denote by $EP = \{ep_1, \dots, ep_n\}$ the set of all possible product units of production under consideration.

Let each product unit ep_i be characterized by a set of parameters $P_i = \{p_{i1}, p_{i2}, \dots, p_{in_i}\}$, $i \in [1, K]$.

Definition 1.3 Sequence

$$AI_i = (A_{i1}, P_{i1}(AI_i), \dots, A_{is}, P_{is}(AI_i)) ,$$

where $P_{ij}(AI_i)$ is the set of parameter values for ep_{ij} in a particular implementation of the technological route AI_i , is called the executed technological route (ETR).

As a result of the production activity of the production under consideration, the set of executed technological routes will be generated at the current time moment t :

$$\mathcal{P}_{\text{ETR}}(t) = \{AI_i : i \in [1, q(t)]\}. \quad (1)$$

Definition 1.4 The terminal vertex of a graph (subgraph) is a vertex from which no arc leaves in this graph (subgraph).

Let V' be the set of terminal vertices of the subgraph $\langle v \cup \vec{G}(v) \cup \vec{G}^2(v) \cup \dots \cup \vec{G}^k(v) \rangle_{\vec{G}}$. Here $\vec{G}^k(v)$ denotes the set of all vertices v' of \vec{G} such that there exists a simple path from the vertex v to the vertex v' of the length $(k - 1)$. For each terminal vertex $v_i \in V'$ there is a certain output unit ep_i , which is the output EP for this vertex, and there may be several such ep_i depending on the type of ETR as a result of which this product unit was received.

Definition 1.5 An ETR is called a productive ETR if the output EP of the terminal vertex A_{is} of this ETR AI_i — denote this vertex as $\text{term}(AI_i)$ — is one of the types of the final product delivered to the market.

Definition 1.6 The vertex $v' \in (v \cup \vec{G}(v) \cup \dots \cup \vec{G}^k(v))$ is called a fork-vertex if $|\vec{G}(v')| > 1$.

In the framework of this paper, the generalized problem of assigning a technological route is considered, in which it is supposed that it is possible to control the choice of the further passage to processing a product unit in the fork-vertices of the graph \vec{G} . This means that the technological route in the process of its execution can be reassigned in order to increase production efficiency and reduce the level of rejection.

2 Formulation of the Problem

Consider the set (1) of all productive ETRs. For each productive ETR AI_i you can define two parameters for EP of its terminal vertex $ep_i = ep(\text{term}(AI_i))$: Price(ep_i) is the market price of the unit of measure ep_i , C(ep_i) is the production cost of the product unit ep_i .

Let there be a set of productive ETRs such that the initial sections to the fork-vertex are coincident in the part of the passage of aggregates. Then each productive ETR can be represented as a sequence:

$$AI_i = (BI_i, CI_i), i \in [1, q(t)] ,$$

where BI_i is the ETR from the initial vertex v_1 to the considered fork-vertex v' and CI_i — ETR from the vertex v' to the terminal vertex $\text{term}(AI_i)$. Since each ETR AI_i passages a certain technological route — we denote such route as $P(AI_i)$ — then the set of all ETRs can be divided into several classes

$$\mathcal{P}_{\text{ETR}}(v, t) = \mathcal{P}_{\text{ETR}}^{(1)}(v, t) \cup \mathcal{P}_{\text{ETR}}^{(2)}(v, t) \cup \dots \cup \mathcal{P}_{\text{ETR}}^{(l)}(v, t) \quad (2)$$

such as AI_i and AI_j belong to the same class if and only if $P(AI_i) = P(AI_j)$.

We denote by $P_i = P(\mathcal{P}_{\text{ETR}}^{(i)}(v, t))$ a technological route which is common for all ETRs from $\mathcal{P}_{\text{ETR}}^{(i)}(v, t)$. Then the problem is to determine which of the technological routes P_i should be chosen for further passage when reaching the fork-vertex v' .

For each class $\mathcal{P}_{\text{ETR}}^{(i)}(v, t)$ from (2) we compile a training sample

$$Z(P_i) : \left(\text{Ef}(AI_j) = \frac{\text{Price}(\text{term}(AI_j)) - C(\text{term}(AI_j))}{C(\text{term}(AI_j))}, i, BI_j \right), \quad (3)$$

where $AI_j \in \mathcal{P}_{\text{ETR}}^{(i)}(v, t)$, and separate the set of values of efficiency $\text{Ef}(AI_j)$ into several discrete intervals E_1, E_2, \dots, E_m .

Next, we represent the sample $Z(P_i)$ in the form of the corresponding multidimensional vectors

$$\bar{a}_j = (a_{j0}, a_{j1}, a_{jm_1}, \dots, a_{jn}),$$

where $a_{j0} \in \{E_1, \dots, E_m\}$ is the value of the ETR's efficiency, a_{j1} is the identifier of the technological route of this ETR. We assign the vector $a_j = (a_{j1}, \dots, a_{jn})$ to the class K_i , if $a_{j0} \in E_i$. Then each a_j vector will be assigned to one of the classes K_1, \dots, K_m and the well-known problem of pattern recognition in geometric formulation arises.

3 Algorithm for Constructing a Decision Rule

A set of n -dimensional vectors is given

$$A = \{ (a_{i_1}, \dots, a_{i_n}) : i \in [1, N] \},$$

and its partition into m classes

$$A = A_1 \cup A_2 \cup \dots \cup A_m.$$

It is required to construct a decision rule for assigning the vector a_i to one of the classes. The solution will be sought in the class of logical decision trees given by a directed binary tree $\vec{G} = (V, E)$ with root vertex $v_0 \in V$.

The binary tree $\vec{G} = (V, E)$ defines the process of sequentially separating of the sample A into two subsamples at the vertices of degree 2 so that each terminal vertex v_i corresponds to a subset $A_{v_i} \subseteq A$, which can be assigned to one of the classes $class_{v_i} \in [1, m]$. In the case under consideration, linear functions will be used to separate the subsample at each vertex of the decision tree.

If v is a vertex of degree 2 in the graph \vec{G} , then a vector n_v and a scalar variable \mathcal{E}_v are given for it, such that A_v is separated into two subsamples of A'_v and A''_v according to the following rule:

$$\begin{aligned} A'_v &= \{ a_i \in A_v : \langle n_v, a_i \rangle \leq \mathcal{E}_v \}, \\ A''_v &= \{ a_i \in A_v : \langle n_v, a_i \rangle > \mathcal{E}_v \}, \end{aligned}$$

and for the root-vertex v_0 we have:

$$A_{v_0} = A.$$

It is required to construct the decision tree $\vec{G} = (V, E)$ with minimal number of vertices, and at each terminal vertex $v \in V$ we have:

$$p(v) = \frac{|\{a_i \in A_v : a_i \in class_v\}|}{|A_v|} \geq p_{\min}, \quad (4)$$

that is, the fraction of vectors belonging to the some class $class_v$ is not less than a given value p_{\min} . If $p_{\min} = 1$ then each terminal vertex corresponds to the vectors of one particular class.

The rule (4) acts if $|A_v| \geq K_{\min}$. If $|A_v| < K_{\min}$ then the process of further separating of the sample A_v is not performed and the vertex v is declared terminal, and the (4) rule may not be executed. In other words, for $|A_v| < K_{\min}$ the sample A_v is not representative enough for constructing a further decision rule and the probability $p(v)$ of the vector from this sample belongs to the class $class_v$ can be less than p_{\min} .

3.1 Algorithm for Constructing the Decision Function for the Node

Suppose that we have a vertex $v \in V$ for which A_v is given. Suppose we have a partition

$$A_v = (A_v \cap A_1) \cup \dots \cup (A_v \cap A_m) , \quad (5)$$

in which there are m' non-empty sets. If $m' = 1$ then the vertex v is terminal and $p(v) = 1$, if $2 \leq m' \leq m$ then sequentially calculate the values:

$$p_i(v) = \frac{|\{a_i \in A_v \cap A_i\}|}{|A_v|} , \quad i \in [1, m] .$$

If there exists $i_0 \in [1, m]$ such that $p_{i_0}(v) \geq p_{\min}$ then the vertex v is terminal and the class $class_v = i_0$, if $|A_v| < K_{\min}$ then the vertex v is terminal and

$$class_v = \arg \max_i \{ |p_i(v)| : i \in [1, m'] \} .$$

Consider the case

$$\begin{cases} 2 \leq m' \leq m , \\ |A_v| \geq K_{\min} , \\ p_i(v) < p_{\min} \quad \forall i \in [1, m'] , \end{cases}$$

and denote by

$$I = \{i : A_v \cap A_i \neq \emptyset , i \in [1, m]\} .$$

Let some vector n_v and a scalar value \mathcal{E}_v be assigned. Then the vertex v is associated with two vertices v_1 and v_2 that are descendants of the vertex v in the constructed decision tree such that:

$$\begin{aligned} A_{v_1} &= \{a_j \in A_v : \langle n_v , a_j \rangle \leq \mathcal{E}_v\} , \\ A_{v_2} &= \{a_j \in A_v : \langle n_v , a_j \rangle > \mathcal{E}_v\} . \end{aligned}$$

Let

$$\begin{aligned} p(A_{v_1}) &= \left(\frac{|A_{v_1} \cap A_1|}{|A_{v_1}|} , \dots , \frac{|A_{v_1} \cap A_n|}{|A_{v_1}|} \right) , \\ p(A_{v_2}) &= \left(\frac{|A_{v_2} \cap A_1|}{|A_{v_2}|} , \dots , \frac{|A_{v_2} \cap A_n|}{|A_{v_2}|} \right) . \end{aligned}$$

In this case the vector n_v and the quantity \mathcal{E}_v require the sets $A_{v_1} \neq \emptyset$ and $A_{v_2} \neq \emptyset$.

Consider the following value

$$\text{discrim}(A_v , n_v , \mathcal{E}_v) = \sum_{i \in I} \left| \frac{|A_{v_1} \cap A_i|}{|A_{v_1}|} - \frac{|A_{v_2} \cap A_i|}{|A_{v_2}|} \right| .$$

The value $\text{discrim}(A_v , n_v , \mathcal{E}_v)$ will be called the separating force of the function

$$f(a) = a \cdot n_v - \mathcal{E}_v$$

concerning the subsample A_v . The meaning of this notion is that the stronger the vectors from the classes A_i of the training sample are separated in the half-space obtained by dividing the space by a hyperplane

$$f(a) = a \cdot n_v - \mathcal{E}_v = 0 ,$$

the more the function $f(a)$ separates vectors from the training sample into classes. Thus, the formulation is natural, where it is required to find $n_v \in \mathbb{R}^{n+1}$ and $\mathcal{E}_v \in \mathbb{R}$ for the sample (5) such that the value of quantity $\text{discrim}(A_v , n_v , \mathcal{E}_v)$ reaches its maximum. The naturalness of such formulation is also confirmed by the fact that for $m = 2$ the best solution is achieved for $\text{discrim}(A_v , n_v , \mathcal{E}_v) = 2$, which corresponds to a linear separation into classes $A_v \cap A_1$ and $A_v \cap A_2$ by the hyperplane $f(a) = n_v \cdot a - \mathcal{E} = 0$. The problem in this formulation always

has a solution, since $A_{v_1} \neq \emptyset$ and $A_{v_2} \neq \emptyset$ for each vertex v its descendants correspond to subsamples of lower power and when the condition (4) or the condition $|A_v| < K_{\min}$ are reached the vertex v becomes terminal.

For an arbitrary subset $A' \subseteq A$ we introduce the notation for the center of the subsample

$$C(A') = \frac{1}{|A'|} \sum_{i \in A'} \{a_i : a_i \in A'\} ,$$

and

$$A(I) = \{a_i : a_i \in A, i \in I\} .$$

Let be given the partition $I = I_1 \cup I_2$, where $I_1 \neq \emptyset, I_2 \neq \emptyset$. Consider the interval $[C(I_1), C(I_2)] \subset \mathbb{R}^{n+1}$. Let $n(I_1, I_2)$ be the normal vector

$$\frac{C(I_2) - C(I_1)}{\|C(I_2) - C(I_1)\|} ,$$

then we divide the interval $[C(I_1), C(I_2)]$ into M parts, where the length of each part is

$$\frac{\|C(I_2) - C(I_1)\|}{M} .$$

We consider the $(M - 1)$ separating functions $f_j(a) = a \cdot n_v - \mathcal{E}_j$, which are passing sequential through all $(M - 1)$ dividing points of the interval $[C(I_1), C(I_2)]$. We will search the best option for the separating force among these functions:

$$j_0 = \arg \max \{ \text{discrim}(A_v, n_v, \mathcal{E}_j) : j \in [1, M - 1] \} .$$

It is easy to see that for j_0 we have $A_{v_1} \neq \emptyset, A_{v_2} \neq \emptyset$.

We denote by

$$\text{discrim}(I_1, I_2) = \text{discrim}(A_v, n_v, \mathcal{E}_{j_0}) .$$

In the case of $C(A(I_1)) = C(A(I_2))$ any two most distant points from the sample A_v are choosed and for the interval which connects these points it is used the same procedure for constructing $(M - 1)$ separating planes and choosing the best of them. In the general case it is assumed that all partitions of the form $I = I_1 \cup I_2$ are searched, and the chosen partition is such that $\text{discrim}(I_1, I_2)$ reaches its maximum. In practical implementation instead of a search a sequential algorithm can be used, in which $I_1 = I, I_2 = \emptyset$ is initially assigned. Further among all partitions of the form $I = I_1 \setminus \{i\} \cup I_2 \cup \{i\}$, where $i \in I$, the best is chosen by the criterion $\text{discrim}(I_1, I_2)$ and so on. In the end, the best result is also chosen from the entire row obtained.

The algorithm of decision tree constructing for the node:

1. Let us consider the node v of decision tree with training sample A_v .
2. The vertex v is declared terminal if there exists a class $i \in [1, m]$ such that $p_i(A_v) \geq p_{\min}$ or $|A_v| \leq K_{\min}$.
3. We consider all partitions $I = I_1 \cup I_2$.
4. For each partition suppose:

$$C_{I_1} = \frac{1}{|I_1|} \sum_{i \in I_1} C_i \quad C_{I_2} = \frac{1}{|I_2|} \sum_{i \in I_2} C_i .$$

5. If the length of the interval $[C_{I_1}, C_{I_2}] \neq 0$ then the normal vector n_{I_1, I_2} is constructed and the sample \mathcal{E}_{I_1, I_2} such that $\text{discrim}(A_v, n_{I_1, I_2}, \mathcal{E}_{I_1, I_2})$ is maximal (search through all hyperplanes perpendicular to n_{I_1, I_2} and bypassing $[C_{I_1}, C_{I_2}]$ for M steps from the point C_{I_1} to the point C_{I_2}).

The process is finite since some separation takes place at the each step.

Example 3.1 Consider an example in which four classes of vectors are given, and each class contains 20 vectors.

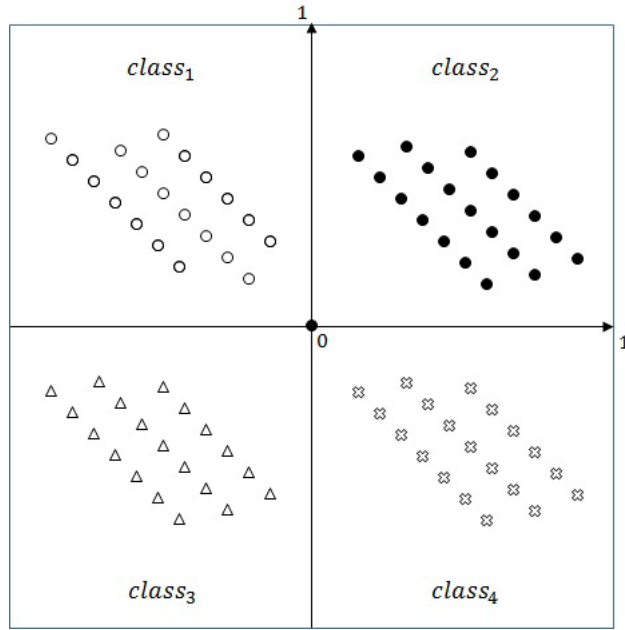


Figure 1: Classes of vectors

1. $A_v = \text{class}_1 \cup \text{class}_2 \cup \text{class}_3 \cup \text{class}_4$,
 $p(A_v) = (0.25, 0.25, 0.25, 0.25)$,
 $I = I_1 \cup I_2 = \{1, 2\} \cup \{3, 4\}$,
 $n_v = (0, 1)$, $\mathcal{E}_v = 0$,
 $p(A_{v_1}) = (0.5, 0.5, 0, 0)$,
 $p(A_{v_2}) = (0, 0, 0.5, 0.5)$,
 $\text{discrim}(A_v, n_v, \mathcal{E}_v) = 0.5 + 0.5 + 0.5 + 0.5 = 2$.

2. $A_{v_1} = \text{class}_1 \cup \text{class}_2$,
 $p(A_{v_1}) = (0.5, 0.5, 0, 0)$,
 $I = \{1\} \cup \{2\}$,
 $n_{v_1} = (1, 0)$, $\mathcal{E}_{v_1} = 0$,
 $p(A_{v_{1,1}}) = (1, 0, 0, 0)$,
 $p(A_{v_{1,2}}) = (0, 1, 0, 0)$,
 $\text{discrim}(A_{v_1}, n_{v_1}, \mathcal{E}_{v_1}) = 1 + 1 + 0 + 0 = 2$.
The vertices $A_{v_{1,1}}$ and $A_{v_{1,2}}$ are terminal.

3. $A_{v_2} = \text{class}_3 \cup \text{class}_4$,
 $p(A_{v_2}) = (0, 0, 0.5, 0.5)$,
 $I = \{3\} \cup \{4\}$,
 $n_{v_2} = (1, 0)$, $\mathcal{E}_{v_2} = 0$,
 $p(A_{v_{2,1}}) = (0, 0, 1, 0)$,
 $p(A_{v_{2,2}}) = (0, 0, 0, 1)$,
 $\text{discrim}(A_{v_2}, n_{v_2}, \mathcal{E}_{v_2}) = 0 + 0 + 1 + 1 = 2$.
The vertices $A_{v_{2,1}}$ and $A_{v_{2,2}}$ are terminal.

Consider the partition $I = I_1 \cup I_2$ and denote by

$$a(I_1, I_2) = \sum \{(a_s - a_t) : a_s \in A_v \cap I_i \forall i \in I_1, a_t \in A_v \cap I_j \forall j \in I_2\} .$$

The most practical efficiently seems the algorithm for separating the sample A_v in the vertex v which chooses the partition $I = I_1 \cup I_2$, for which $|a(I_1, I_2)|$ is maximal among all partitions $I = I_1 \cup I_2$. Then the choice of the vector $n(I_1, I_2)$ and the scalar value $\mathcal{E}(I_1, I_2)$ can be made according to the procedure described above for the obtained fixed partition $I = I_1 \cup I_2$.

We introduce the notation:

$$a_{ij} = \sum \{(a_s - a_t) : a_s \in A_i, a_t \in A_j\}, i, j \in [1, m] ,$$

then for $I = I_1 \cup I_2$ we have

$$a(I_1, I_2) = \sum \{a_{ij} : i \in I_1, j \in I_2\} . \quad (6)$$

Proceeding from the relation (6) it is possible to significantly reduce the amount of computation when choosing the optimal partition $I = I_1 \cup I_2$ for the sample A_v using the previously calculated values $a_{ij} \forall i, j \in I$.

Conclusion

To solve the problem of the quality of the product of metallurgical production, an approach has been developed in which the research is reduced to solving the problem of pattern recognition in a geometric formulation. A heuristic algorithm for constructing a decision rule is developed which is aiming to find the best possible discrimination of training sample corresponding to each vertex of decision tree. Substantial reducing of the computational complexity of considered algorithm is proposed. Thus, for the production under consideration, a large number of problems of pattern recognition are constructed for each of which a decision tree is constructed.

An important feature of the approach is that the set $\mathcal{P}(t)$ of ETRs is continuously expanding, thus providing all the new data to improve the decision rule. To achieve the effectiveness of the proposed approach in practice it is necessary to carry out additional training as soon as a new portion of the ETRs arrives. This will ensure the continuous improvement of the decision rules and consequently the improvement of production efficiency.

References

- [Gainanov & Berenov, 2017] Gainanov, D. N., & Berenov, D. A. (2017). Big Data Technologies in Metallurgical Production Quality Control Systems. *Proceedings of the Conference on Big Data and Advanced Analytics*. (pp. 65-70). Minsk, Belarus: Minsk State University Press.
- [Gainanov, 2014] Gainanov, D. N. (2014). *Combinatorial geometry and graphs in the analysis of infeasible systems and pattern recognition*. Moscow: Nauka.
- [Gainanov, 2016] Gainanov, D. N. (2016). *Graphs for Pattern Recognition. Infeasible Systems of Linear Inequalities*. DeGruyter.
- [Gainanov, 1985] Gainanov, D. N. (1985). Combinatorial properties of infeasible systems of linear inequalities and convex polyhedra. *Math notices*, 3(38), 463-474.
- [Mazurov, 1990] Mazurov, Vl. D. (1990). *Committees method in problems of optimization and classification*. Moscow: Nauka.
- [Mazurov, 2004] Mazurov, Vl. D., Khachai, M. Yu. (2004). Committees of systems of linear inequalities. *Automation and Remote Control*, 2, 43-54.
- [Khachai, 1997] Khachai, M. Yu. (1997). On the estimate of the number of members of the minimal committee of a system of linear inequalities. *Jour. of Computational Math. and Math. Physics*, 11(37), 1399-1404.
- [Gainanov, 1992] Gainanov, D. N. (1992). Alternative Covers and Independence Systems in Pattern Recognition. *Pattern Recognition and Image Analysis*, 2(2), 147-160.
- [Gainanov, 1991] Gainanov, D. N., & Matveev, A. O. (1991). Lattice Diagonals and Geometric Pattern Recognition Problems // Pattern Recognition and Image Analysis. *Pattern Recognition and Image Analysis*, 3(1), 277-282.