

# Construction of $\alpha$ -Decision Trees for Tables with Many-Valued Decisions

Mikhail Moshkov<sup>1</sup> and Beata Zielosko<sup>1,2</sup>

<sup>1</sup> Mathematical and Computer Sciences & Engineering Division  
King Abdullah University of Science and Technology  
Thuwal 23955-6900, Saudi Arabia

{mikhail.moshkov, beata.zielosko}@kaust.edu.sa

<sup>2</sup> Institute of Computer Science, University of Silesia  
39, Będzińska St., 41-200 Sosnowiec, Poland

**Abstract.** The paper is devoted to the study of greedy algorithm for construction of approximate decision trees ( $\alpha$ -decision trees). This algorithm is applicable to decision tables with many-valued decisions where each row is labeled with a set of decisions. For a given row, we should find a decision from the set attached to this row. We consider bound on the number of algorithm steps, and bound on the algorithm accuracy relative to the depth of decision trees.

**Keywords:** decision table with many-valued decisions, decision tree, greedy algorithm.

## 1 Introduction

We consider here one more extension of the notion of decision table – decision table with many-valued decisions. In a table with many-valued decisions, each row is labeled with a nonempty finite set of decisions, and for a given row, we should find a decision from the set of decisions attached to this row.

Such tables arises in problems of discrete optimization, pattern recognition, computational geometry, etc. However, the main source of decision tables with many-valued decisions are datasets filled by statistical or experimental data. In such datasets, we often have groups of objects with equal values of conditional attributes but, probably, different values of the decision attribute. Instead of a group of objects, we can consider one object given by values of conditional attributes. We attach to this object a set of decisions: either all decisions for objects from the group, or  $k$  the most frequent decisions for objects from the group, etc. As a result we obtain a decision table with many-valued decisions.

The rough set theory [2,4] is devoted mainly to the investigation of inconsistent decision tables which have equal rows with different decisions (see, in particular, the notion of generalized decision in rough set theory [3]). So, the study of decision tables with many-valued decisions can give us new tools for the rough set theory.

The paper is devoted to the study of a greedy algorithm for construction of approximate decision trees for decision tables with many-valued decisions. To define the notion of approximate decision tree we fix some uncertainty measure  $B(T)$  for decision tables with many-valued decisions. We consider so-called  $\alpha$ -decision trees where  $\alpha$  is a real number such that  $0 \leq \alpha < 1$ . For a given row  $r$  of a table  $T$ , an  $\alpha$ -decision tree localizes it in a subtable of  $T$  with uncertainty at most  $\alpha B(T)$ . The notion of 0-decision tree for  $T$  coincides with the notion of exact decision tree for  $T$ .

We prove new bound on accuracy of the greedy algorithm relative to the depth of decision trees. As a corollary, we obtain a bound on accuracy of the greedy algorithm presented in [1] without proof. We obtain also an upper bound on the number of steps of the considered algorithm. From this bound it follows that, for an arbitrary natural  $t$ , the greedy algorithm has polynomial time complexity on tables which have at most  $t$  decisions in each set of decisions attached to rows.

As we know, the considered algorithm is the only algorithm for construction of decision trees for decision tables with many-valued decisions that has nontrivial bound on accuracy.

We discuss also a problem of recognition of colors of points in the plain which illustrates the obtained bound on accuracy.

In this paper, we consider only binary decision tables with many-valued decisions. However, the obtained results can be extended to the decision tables filled by numbers from the set  $\{0, \dots, k-1\}$ , where  $k \geq 3$ .

This paper consists of six sections. In Sect. 2, main notions are discussed. In Sect. 3, four lemmas are proved which are used in Sect. 4, that is devoted to the study of greedy algorithm for decision tree construction. In Sect. 5, we discuss the problem of recognition of colors of points in the plain. Section 6 contains conclusions.

## 2 Main Notions

In this section, we consider definitions of notions corresponding to decision tables with many-valued decisions.

A *(binary) decision table with many-valued decisions* is a rectangular table  $T$  filled by numbers from the set  $\{0, 1\}$ . Columns of this table are labeled with attributes  $f_1, \dots, f_n$ . Rows of the table are pairwise different, and each row is labeled with a nonempty finite set of natural numbers (set of decisions). Note that each (binary) decision table with one-valued decisions can be interpreted also as a decision table with many-valued decisions. In such table, each row is labeled with a set of decisions which has one element.

We correspond a *game* of two players to  $T$ . The first player chooses a row  $r$  of  $T$ . The second player should find a decision from the set of decisions attached to  $r$ . To this end, he can choose columns (attributes) of  $T$  and ask the first player what is at the intersection of the row  $r$  and these columns.

We will say that  $T$  is a *degenerate* table if either  $T$  has no rows, or the intersection of sets of decisions attached to rows of  $T$  is nonempty.

A decision which belongs to the maximum number of sets of decisions attached to rows in  $T$  is called the *most common decision for  $T$* . If we have more than one such decisions we choose the minimum one. If  $T$  is empty then 1 is the most common decision for  $T$ .

A table obtained from  $T$  by removal of some rows is called a *subtable* of  $T$ . A subtable  $T'$  of  $T$  is called *boundary subtable* if  $T'$  is not degenerate but each proper subtable of  $T'$  is degenerate. We denote by  $B(T)$  the number of boundary subtables of the table  $T$ . It is clear that  $T$  is a degenerate table if and only if  $B(T) = 0$ . The value  $B(T)$  will be interpreted as *uncertainty* of  $T$ .

Let  $f_{i_1}, \dots, f_{i_m} \in \{f_1, \dots, f_n\}$  and  $\delta_1, \dots, \delta_m \in \{0, 1\}$ . We denote by

$$T(f_{i_1}, \delta_1) \dots (f_{i_m}, \delta_m)$$

the subtable of the table  $T$  which consists of all rows that at the intersection with columns  $f_{i_1}, \dots, f_{i_m}$  have numbers  $\delta_1, \dots, \delta_m$  respectively.

A *decision tree over  $T$*  is a finite tree with root in which each terminal node is labeled with a decision (a natural number), and each nonterminal node is labeled with an attribute from the set  $\{f_1, \dots, f_n\}$ . Two edges start in each nonterminal node. These edges are labeled with 0 and 1 respectively.

Let  $\Gamma$  be a decision tree over  $T$  and  $v$  be a node of  $\Gamma$ . We correspond to the node  $v$  a subtable  $T(v)$  of the table  $T$ . If  $v$  is the root of  $\Gamma$  then  $T(v) = T$ . Otherwise, let nodes and edges in the path from the root to  $v$  be labeled with attributes  $f_{i_1}, \dots, f_{i_m}$  and numbers  $\delta_1, \dots, \delta_m$  respectively. Then  $T(v)$  is the subtable  $T(f_{i_1}, \delta_1) \dots (f_{i_m}, \delta_m)$  of the table  $T$ .

It is clear that for any row  $r$  of  $T$  there exists exactly one terminal node  $v$  in  $\Gamma$  such that  $r$  belongs to  $T(v)$ . The decision attached to  $v$  will be considered as the result of  $\Gamma$  work on the row  $r$ .

Let  $\alpha$  be a real number such that  $0 \leq \alpha < 1$ . We will say that  $\Gamma$  is an  $\alpha$ -*decision tree for  $T$*  if for any terminal node  $v$  of  $\Gamma$ , the inequality  $B(T(v)) \leq \alpha B(T)$  holds, and the node  $v$  is labeled with the most common decision for  $T(v)$ . An  $\alpha$ -decision tree  $\Gamma$  for  $T$  can be considered as an approximate strategy for the second player: for any row  $r$  of  $T$ , as a result of  $\Gamma$  work, the row  $r$  will be localized in a subtable  $T(v)$  of  $T$  ( $v$  is a terminal node of  $\Gamma$ ) such that the uncertainty of  $T(v)$  is at most  $\alpha$  times to the uncertainty of the initial table  $T$ .

We denote by  $h(\Gamma)$  the *depth* of a decision tree  $\Gamma$  which is the maximum length of a path from the root to a terminal node. We denote by  $h_\alpha(T)$  the minimum depth of an  $\alpha$ -decision tree for the table  $T$ .

Let  $\alpha, \beta$  be real numbers such that  $0 \leq \alpha \leq \beta < 1$ . It is not difficult to show that each  $\alpha$ -decision tree for  $T$  is also a  $\beta$ -decision tree for  $T$ . Thus,  $h_\alpha(T) \geq h_\beta(T)$ .

### 3 Auxiliary Statements

We denote by  $Tab(t)$ , where  $t$  is a natural number, the set of decision tables with many-valued decisions such that each row in the table has at most  $t$  decisions (is labeled with a set of decisions which cardinality is at most  $t$ ).

**Lemma 1.** *Let  $T'$  be a boundary subtable with  $m$  rows. Then each row of  $T'$  is labeled with a set of decisions which cardinality is at least  $m - 1$ .*

*Proof.* Let rows of  $T'$  be labeled with sets of decisions  $D_1, \dots, D_m$  respectively. Then  $D_1 \cap \dots \cap D_m = \emptyset$  and for any  $i \in \{1, \dots, m\}$ , the set  $D_1 \cap \dots \cap D_{i-1} \cap D_{i+1} \cap \dots \cap D_m$  contains a number  $d_i$ . Assume that  $i \neq j$  and  $d_i = d_j$ . Then  $D_1 \cap \dots \cap D_m \neq \emptyset$  which is impossible. Therefore  $d_1, \dots, d_m$  are pairwise different numbers. It is clear that for  $i = 1, \dots, m$ , the set  $\{d_1, \dots, d_m\} \setminus \{d_i\}$  is a subset of the set  $D_i$ . □

**Corollary 1.** *Each boundary subtable of a table  $T \in Tab(t)$  has at most  $t + 1$  rows.*

Therefore, for tables from  $Tab(t)$ , there exists a polynomial algorithm for the computation of parameter  $B(T)$ . For example, for any decision table  $T$  with one-valued decisions (really, for any table from  $Tab(1)$ ) the equality  $B(T) = P(T)$  holds where  $P(T)$  is the number of unordered pairs of rows of  $T$  with different decisions.

**Lemma 2.** *Let  $T$  be a decision table with many-valued decisions,  $T'$  be a subtable of the table  $T$ ,  $f_i$  be an attribute attached to a column of  $T$ , and  $\delta \in \{0, 1\}$ . Then*

$$B(T) - B(T(f_i, \delta)) \geq B(T') - B(T'(f_i, \delta)).$$

*Proof.* Denote by  $J$  (respectively by  $J'$ ) the set of boundary subtables of  $T$  (respectively of  $T'$ ) in each of which at least one row has at the intersection with column  $f_i$  a number which is not equal to  $\delta$ . One can show that  $J' \subseteq J$ ,  $|J'| = B(T') - B(T'(f_i, \delta))$  and  $|J| = B(T) - B(T(f_i, \delta))$ . □

Let  $T$  be a decision table with many-valued decisions, which has  $n$  columns labeled with attributes  $\{f_1, \dots, f_n\}$ . We define now a parameter  $M(T)$  of the table  $T$ . If  $T$  is a degenerate table then  $M(T) = 0$ . Let now  $T$  be a nondegenerate table. Let  $\bar{\delta} = (\delta_1, \dots, \delta_n) \in \{0, 1\}^n$ . Then  $M(T, \bar{\delta})$  is the minimum natural  $m$  such that there exist attributes  $f_{i_1}, \dots, f_{i_m} \in \{f_1, \dots, f_n\}$  for which  $T(f_{i_1}, \delta_{i_1}) \dots (f_{i_m}, \delta_{i_m})$  is a degenerate table. We denote  $M(T) = \max\{M(T, \bar{\delta}) : \bar{\delta} \in \{0, 1\}^n\}$ .

**Lemma 3.** *Let  $T$  be a decision table with many-valued decisions, and  $T'$  be a subtable of  $T$ . Then*

$$M(T) \geq M(T').$$

*Proof.* Let  $T$  have  $n$  columns labeled with attributes  $f_1, \dots, f_n, f_{i_1}, \dots, f_{i_m} \in \{f_1, \dots, f_n\}$  and  $\delta_1, \dots, \delta_m \in \{0, 1\}$ . If  $T(f_{i_1}, \delta_1) \dots (f_{i_m}, \delta_m)$  is a degenerate table then  $T'(f_{i_1}, \delta_1) \dots (f_{i_m}, \delta_m)$  is a degenerate table too. From here and from the definition of parameter  $M$  the statement of lemma follows. □

**Lemma 4.** *Let  $T$  be a decision table. Then*

$$h_0(T) \geq M(T).$$

*Proof.* If  $T$  is a degenerate table then  $h_0(T) = 0$  and  $M(T) = 0$ . Let  $T$  be a nondegenerate table having  $n$  columns labeled with attributes  $f_1, \dots, f_n$ , and  $\Gamma$  be a 0-decision tree for the table  $T$  such that  $h(\Gamma) = h_0(T)$ . Let  $\bar{\delta} = (\delta_1, \dots, \delta_n) \in \{0, 1\}^n$  be a  $n$ -tuple for which  $M(T, \bar{\delta}) = M(T)$ . We consider a path  $\tau = v_1, d_1, \dots, v_m, d_m, v_{m+1}$  from the root  $v_1$  to a terminal node  $v_{m+1}$  in  $\Gamma$  which satisfies the following condition: if nodes  $v_1, \dots, v_m$  are labeled with attributes  $f_{i_1}, \dots, f_{i_m}$  then edges  $d_1, \dots, d_m$  are labeled with numbers  $\delta_{i_1}, \dots, \delta_{i_m}$ . Since  $\Gamma$  is a 0-decision tree for the table  $T$ , the subtable  $T(v_{m+1}) = T(f_{i_1}, \delta_{i_1}) \dots (f_{i_m}, \delta_{i_m})$  is a degenerate table. Therefore  $m \geq M(T, \bar{\delta})$  and  $h(\Gamma) \geq M(T, \bar{\delta})$ . Since  $h(\Gamma) = h_0(T)$  and  $M(T, \bar{\delta}) = M(T)$ , we have  $h_0(T) \geq M(T)$ .  $\square$

### 4 Algorithm $U_\alpha$ for $\alpha$ -Decision Tree Construction

Let  $\alpha$  be a real number such that  $0 \leq \alpha < 1$ . We now describe an algorithm  $U_\alpha$  which for a given decision table with many-valued decisions  $T$  constructs an  $\alpha$ -decision tree  $U_\alpha(T)$  for the table  $T$ . Let  $T$  have  $n$  columns labeled with attributes  $f_1, \dots, f_n$ .

#### *Greedy Algorithm $U_\alpha$*

*Step 1.* Construct a tree consisting of a single node labeled with the table  $T$  and proceed to the second step.

Suppose  $t \geq 1$  steps have been made already. The tree obtained at the step  $t$  will be denoted by  $G$ .

*Step  $(t + 1)$ .* If no one node of the tree  $G$  is labeled with a table then we denote by  $U_\alpha(T)$  the tree  $G$ . The work of the algorithm  $U_\alpha$  is completed.

Otherwise, we choose a node  $v$  in the tree  $G$  which is labeled with a subtable of the table  $T$ . Let the node  $v$  be labeled with the table  $T'$ . If  $B(T') \leq \alpha B(T)$  then instead of  $T'$  we mark the node  $v$  with the most common decision for  $T'$  and proceed to the step  $(t + 2)$ . Let  $B(T') > \alpha B(T)$ . Then for  $i = 1, \dots, n$ , we compute the value

$$Q(f_i) = \max\{B(T'(f_i, 0)), B(T'(f_i, 1))\}.$$

Instead of  $T'$  we mark the node  $v$  with the attribute  $f_{i_0}$  where  $i_0$  is the minimum  $i$  for which  $Q(f_i)$  has the minimum value. For each  $\delta \in \{0, 1\}$ , we add to the tree  $G$  the node  $v(\delta)$ , mark this node with the subtable  $T'(f_{i_0}, \delta)$ , draw the edge from  $v$  to  $v(\delta)$ , and mark this edge with  $\delta$ . Proceed to the step  $(t + 2)$ .

First, we obtain an upper bound on the number of algorithm  $U_\alpha$  steps.

**Theorem 1.** *Let  $\alpha$  be a real number such that  $0 \leq \alpha < 1$ , and  $T$  be a decision table with many-valued decisions. Then during the construction of the tree  $U_\alpha(T)$  the algorithm  $U_\alpha$  makes at most  $2N(T) + 1$  steps where  $N(T)$  is the number of rows in  $T$ .*

*Proof.* One can show that for each terminal node  $v$  of the tree  $U_\alpha(T)$ , there exists a row  $r(v)$  of  $T$  such that  $r(v)$  belongs to  $T(v)$ . It is clear that  $r(v_1) \neq r(v_2)$  if  $v_1 \neq v_2$ . Therefore the number of terminal nodes in  $U_\alpha(T)$  is at most  $N(T)$ . It is not difficult to prove that the number of nonterminal nodes in  $U_\alpha(T)$  is equal to the number of terminal nodes minus 1. Simple analysis of the algorithm  $U_\alpha$  work shows that the number of steps of  $U_\alpha$  in the process of the tree  $U_\alpha(T)$  construction is equal to the number of nodes in  $U_\alpha(T)$  plus 2. Therefore the number of steps is bounded from above by  $2N(T) + 1$ .  $\square$

From this theorem it follows that for any natural  $t$  the algorithm  $U_\alpha$  has polynomial time complexity on the set  $Tab(t)$ .

Now we obtain a bound on the algorithm  $U_\alpha$ ,  $0 < \alpha < 1$ , accuracy (relative to the depth of decision trees) which does not depend on  $B(T)$ .

**Theorem 2.** *Let  $\alpha$  be a real number such that  $0 < \alpha < 1$ , and  $T$  be a nondegenerate decision table with many-valued decisions. Then*

$$h(U_\alpha(T)) \leq M(T) \ln \frac{1}{\alpha} + 1.$$

*Proof.* Let  $T$  be a table with  $n$  columns labeled with attributes  $f_1, \dots, f_n$ . For  $i = 1, \dots, n$ , denote by  $\sigma_i$  a number from  $\{0, 1\}$  such that  $B(T(f_i, \sigma_i)) = \max\{B(T(f_i, \sigma)) : \sigma \in \{0, 1\}\}$ . It is clear that the root of the tree  $U_0(T)$  is labeled with the attribute  $f_{i_0}$  where  $i_0$  is the minimum  $i$  for which  $B(T(f_i, \sigma_i))$  has the minimum value. Of course,  $Q(f_i) = B(T(f_i, \sigma_i))$ .

Let us show that

$$B(T(f_{i_0}, \sigma_{i_0})) \leq \left(1 - \frac{1}{M(T)}\right) B(T).$$

It is clear that there exist attributes  $f_{i_1}, \dots, f_{i_m} \in \{f_1, \dots, f_n\}$  such that

$$T(f_{i_1}, \sigma_{i_1}) \dots (f_{i_m}, \sigma_{i_m})$$

is a degenerate table and  $m \leq M(T)$ . Evidently,  $B(T(f_{i_1}, \sigma_{i_1}) \dots (f_{i_m}, \sigma_{i_m})) = 0$ . Then

$$\begin{aligned} & B(T) - [B(T) - B(T(f_{i_1}, \sigma_{i_1}))] - [B(T(f_{i_1}, \sigma_{i_1})) - B(T(f_{i_1}, \sigma_{i_1})(f_{i_2}, \sigma_{i_2}))] \\ & \quad - \dots - [B(T(f_{i_1}, \sigma_{i_1}) \dots (f_{i_{m-1}}, \sigma_{i_{m-1}})) - B(T(f_{i_1}, \sigma_{i_1}) \dots (f_{i_m}, \sigma_{i_m}))] \\ & \quad = B(T(f_{i_1}, \sigma_{i_1}) \dots (f_{i_m}, \sigma_{i_m})) = 0. \end{aligned}$$

From Lemma 2 it follows that, for  $j = 1, \dots, m - 1$ ,

$$\begin{aligned} & B(T(f_{i_1}, \sigma_{i_1}) \dots (f_{i_j}, \sigma_{i_j})) - B(T(f_{i_1}, \sigma_{i_1}) \dots (f_{i_j}, \sigma_{i_j})(f_{i_{j+1}}, \sigma_{i_{j+1}})) \\ & \quad \leq B(T) - B(T(f_{i_{j+1}}, \sigma_{i_{j+1}})). \end{aligned}$$

Therefore  $B(T) - \sum_{j=1}^m (B(T) - B(T(f_{i_j}, \sigma_{i_j}))) \leq 0$ . Since  $B(T(f_{i_0}, \sigma_{i_0})) \leq B(T(f_{i_j}, \sigma_{i_j}))$ ,  $j = 1, \dots, m$ , we have  $B(T) - m(B(T) - B(T(f_{i_0}, \sigma_{i_0}))) \leq 0$  and

$B(T(f_{i_0}, \sigma_{i_0})) \leq (1 - 1/m)B(T)$ . Taking into account that  $m \leq M(T)$  we obtain  $B(T(f_{i_0}, \sigma_{i_0})) \leq (1 - 1/M(T)) B(T)$ .

Assume that  $M(T) = 1$ . From the considered inequality and from the description of algorithm  $U_\alpha$  it follows that  $h(U_\alpha(T)) = 1$ . So if  $M(T) = 1$  then the statement of theorem is true.

Let now  $M(T) \geq 2$ . Consider a longest path in the tree  $U_\alpha(T)$  from the root to a terminal node. Let its length be equal to  $k$ , nonterminal nodes of this path be labeled with attributes  $f_{j_1}, \dots, f_{j_k}$ , where  $f_{j_1} = f_{i_0}$ , and edges be labeled with numbers  $\delta_1, \dots, \delta_k$ . For  $t = 1, \dots, k$ , we denote by  $T_t$  the table  $T(f_{j_1}, \delta_1) \dots (f_{j_t}, \delta_t)$ . From Lemma 3 it follows that  $M(T_t) \leq M(T)$  for  $t = 1, \dots, k$ . We know that  $B(T_1) \leq B(T)(1 - 1/M(T))$ . In the same way, it is possible to prove that  $B(T_t) \leq B(T)(1 - 1/M(T))^t$  for  $t = 2, \dots, k$ .

Let us consider the table  $T_{k-1}$ . For this table,

$$B(T_{k-1}) \leq B(T)(1 - 1/M(T))^{k-1}.$$

Using the description of the algorithm  $U_\alpha$  we obtain  $B(T_{k-1}) > \alpha B(T)$ . Therefore  $\alpha < (1 - 1/M(T))^{k-1}$  and  $(1 + 1/(M(T) - 1))^{k-1} < 1/\alpha$ . If we take natural logarithm of both sides of this inequality we obtain  $(k - 1) \ln(1 + 1/(M(T) - 1)) < \ln(1/\alpha)$ . It is known that for any natural  $p$  the inequality  $\ln(1+1/p) > 1/(p+1)$  holds. Since  $M(T) \geq 2$ , we obtain  $(k-1)/M(T) < \ln(1/\alpha)$  and  $k < M(T) \ln(1/\alpha) + 1$ . Taking into account that  $k = h(U_\alpha(T))$  we have  $h(U_\alpha(T)) < M(T) \ln(1/\alpha) + 1$ . □

Using Lemma 4 we obtain

**Corollary 2.** *For any real  $\alpha$ ,  $0 < \alpha < 1$ , and for any nondegenerate decision table with many-valued decisions  $T$*

$$h(U_\alpha(T)) < h_0(T) \ln \frac{1}{\alpha} + 1.$$

## 5 Problem of Recognition of Colors of Points in the Plain

Let we have a finite set  $S = \{(a_1, b_1), \dots, (a_n, b_n)\}$  of points in the plane and a mapping  $\mu$  which corresponds to each point  $(a_p, b_p)$  a nonempty subset  $\mu(a_p, b_p)$  of the set  $\{green, yellow, red\}$ . Colors are interpreted as decisions, and for each point from  $S$  we need to find a decision (color) from the set of decisions attached to this point. We denote this problem by  $(S, \mu)$ .

For the problem  $(S, \mu)$  solving, we use attributes corresponding to straight lines which are given by equations of the kind  $x = \beta$  or  $y = \gamma$ . These attributes are defined on the set  $S$  and take values from the set  $\{0, 1\}$ . Consider the line given by equation  $x = \beta$ . Then the value of corresponding attribute is equal to 0 on a point  $(a, b) \in S$  if and only if  $a < \beta$ . Consider the line given by equation  $y = \gamma$ . Then the value of corresponding attribute is equal to 0 if and only if  $b < \gamma$ .

We now choose a finite set of straight lines which allow us to construct a decision tree with the minimum depth for the problem  $(S, \mu)$ . It is possible that  $a_i = a_j$  or  $b_i = b_j$  for different  $i$  and  $j$ . Let  $a_{i_1}, \dots, a_{i_m}$  be all pairwise different numbers from the set  $\{a_1, \dots, a_n\}$  which are ordered such that  $a_{i_1} < \dots < a_{i_m}$ . Let  $b_{j_1}, \dots, b_{j_t}$  be all pairwise different numbers from the set  $\{b_1, \dots, b_n\}$  which are ordered such that  $b_{j_1} < \dots < b_{j_t}$ .

One can show that there exists a decision tree with minimum depth which use only attributes corresponding to the straight lines defined by equations  $x = a_{i_1} - 1, x = (a_{i_1} + a_{i_2})/2, \dots, x = (a_{i_{m-1}} + a_{i_m})/2, x = a_{i_m} + 1, y = b_{j_1} - 1, y = (b_{j_1} + b_{j_2})/2, \dots, y = (b_{j_{t-1}} + b_{j_t})/2, y = b_{j_t} + 1$ .

We now describe a decision table  $T(S, \mu)$  with  $m + t + 2$  columns and  $n$  rows. Columns of this table are labeled with attributes  $f_1, \dots, f_{m+t+2}$ , corresponding to the considered  $m + t + 2$  lines. Attributes  $f_1, \dots, f_{m+1}$  correspond to lines defined by equations  $x = a_{i_1} - 1, x = (a_{i_1} + a_{i_2})/2, \dots, x = (a_{i_{m-1}} + a_{i_m})/2, x = a_{i_m} + 1$  respectively. Attributes  $f_{m+2}, \dots, f_{m+t+2}$  correspond to lines defined by equations  $y = b_{j_1} - 1, y = (b_{j_1} + b_{j_2})/2, \dots, y = (b_{j_{t-1}} + b_{j_t})/2, y = b_{j_t} + 1$  respectively. Rows of the table  $T(S, \mu)$  correspond to points  $(a_1, b_1), \dots, (a_n, b_n)$ . At the intersection of the column  $f_l$  and row  $(a_p, b_p)$  the value  $f_l(a_p, b_p)$  stays. For  $p = 1, \dots, n$ , the row  $(a_p, b_p)$  is labeled with the set of decisions  $\mu(a_p, b_p)$ .

Let us evaluate the parameter  $M(T(S, \mu))$ .

**Proposition 1.**  $M(T(S, \mu)) \leq 4$ .

*Proof.* Denote  $T = T(S, \mu)$ . Let  $\bar{\delta} = (\delta_1, \dots, \delta_{m+t+2}) \in \{0, 1\}^{m+t+2}$ . If  $\delta_1 = 0$ , or  $\delta_{m+1} = 1$ , or  $\delta_{m+2} = 0$ , or  $\delta_{m+t+2} = 1$ , then  $T(f_1, \delta_1)$ , or  $T(f_{m+1}, \delta_{m+1})$ , or  $T(f_{m+2}, \delta_{m+2})$ , or  $T(f_{m+t+2}, \delta_{m+t+2})$  is empty table and  $M(T, \bar{\delta}) \leq 1$ . Let  $\delta_1 = 1, \delta_{m+1} = 0, \delta_{m+2} = 1$  and  $\delta_{m+t+2} = 0$ . One can show that in this case there exist  $i \in \{1, \dots, m\}$  and  $j \in \{m+2, \dots, m+t+1\}$  such that  $\delta_i = 1, \delta_{i+1} = 0, \delta_j = 1$ , and  $\delta_{j+1} = 0$ . It is clear that the table  $T(f_i, \delta_i)(f_{i+1}, \delta_{i+1})(f_j, \delta_j)(f_{j+1}, \delta_{j+1})$  contains exactly one row. So  $M(T, \bar{\delta}) \leq 4$  and  $M(T) \leq 4$ .  $\square$

From Corollary 1 it follows that each boundary subtable of the table  $T(S, \mu)$  has at most three rows. Thus, the following statement holds:

**Proposition 2.**  $B(T(S, \mu)) \leq |S|^3$ .

Note that there are two types of boundary subtables of the table  $T(S, \mu)$ :

- With two rows labeled with disjoint sets of decisions, for example,  $\{yellow\}$  and  $\{green, red\}$ , or  $\{yellow\}$  and  $\{green\}$ .
- With three rows labeled with sets of decisions  $\{green, yellow\}, \{green, red\}$ , and  $\{yellow, red\}$ .

Using this fact we can easily compute the value  $B(T(S, \mu))$ :

$$B(T(S, \mu)) = N(g)N(y) + N(g)N(r) + N(y)N(r) + N(g)N(y, r) + N(y)N(g, r) + N(r)N(g, y) + N(g, y)N(g, r)N(y, r),$$



where  $N(g)$  is the number of columns in  $T(S, \mu)$  labeled with the set of decisions  $\{green\}$ ,  $N(g, r)$  is the number of columns in  $T(S, \mu)$  labeled with the set of decisions  $\{green, red\}$ , etc.

From Theorem 2 and Proposition 1 the next statement follows:

**Corollary 3.** *For any real  $\alpha$ ,  $0 < \alpha < 1$ ,*

$$h(U_\alpha(T(S, \mu))) < 4 \ln \frac{1}{\alpha} + 1.$$

## 6 Conclusions

We studied algorithm  $U_\alpha$ ,  $0 \leq \alpha < 1$ , which allows us to construct  $\alpha$ -decision trees for decision tables with many-valued decisions. We proved that for an arbitrary natural  $t$ , the considered algorithm has polynomial time complexity on tables which have at most  $t$  decisions in each set of decisions attached to rows. We obtained bound on accuracy of this algorithm which does not depend on the uncertainty of decision tables.

**Acknowledgements.** The authors wish to express their thanks to anonymous reviewers for useful comments.

## References

1. Moshkov, M.: Greedy algorithm for decision tree construction in context of knowledge discovery problems. In: Tsumoto, S., Słowiński, R., Komorowski, J., Grzymała-Busse, J.W. (eds.) RSCCTC 2004. LNCS (LNAI), vol. 3066, pp. 192–197. Springer, Heidelberg (2004)
2. Pawlak, Z.: Rough Sets – Theoretical Aspects of Reasoning about Data. Kluwer Academic Publishers, Dordrecht (1991)
3. Pawlak, Z., Skowron, A.: Rudiments of rough sets. Information Sciences 177(1), 3–27 (2007); Rough sets: Some extensions. Information Sciences 177(1), 28–40 (2007); Rough sets and boolean reasoning. Information Sciences 177(1), 41–73 (2007)
4. Skowron, A., Rauszer, C.: The discernibility matrices and functions in information systems. In: Slowinski, R. (ed.) Intelligent Decision Support. Handbook of Applications and Advances of the Rough Set Theory, pp. 331–362. Kluwer Academic Publishers, Dordrecht (1992)