

Statistical Model for Rough Set Approach to Multicriteria Classification

Krzysztof Dembczyński¹, Salvatore Greco², Wojciech Kotłowski¹,
and Roman Słowiński^{1,3}

¹ Institute of Computing Science, Poznań University of Technology,
60-965 Poznań, Poland

{[kdembczynski](mailto:kdembczynski@cs.put.poznan.pl),[wkotlowski](mailto:wkotlowski@cs.put.poznan.pl),[rslowinski](mailto:rslowinski@cs.put.poznan.pl)}@cs.put.poznan.pl

² Faculty of Economics, University of Catania, 95129 Catania, Italy
salgreco@unict.it

³ Institute for Systems Research, Polish Academy of Sciences, 01-447 Warsaw, Poland

Abstract. In order to discover interesting patterns and dependencies in data, an approach based on rough set theory can be used. In particular, Dominance-based Rough Set Approach (DRSA) has been introduced to deal with the problem of multicriteria classification. However, in real-life problems, in the presence of noise, the notions of rough approximations were found to be excessively restrictive, which led to the proposal of the Variable Consistency variant of DRSA. In this paper, we introduce a new approach to variable consistency that is based on maximum likelihood estimation. For two-class (binary) problems, it leads to the isotonic regression problem. The approach is easily generalized for the multi-class case. Finally, we show the equivalence of the variable consistency rough sets to the specific risk-minimizing decision rule in statistical decision theory.

1 Introduction

In decision analysis, a multicriteria classification problem is considered that consists in assignment of objects to m decision classes Cl_t , $t \in T = \{1, \dots, m\}$. The classes are preference ordered according to an increasing order of class indices, i.e. for all $r, s \in T$, such that $r > s$, the objects from Cl_r are strictly preferred to objects from Cl_s . Objects are evaluated on a set of *condition criteria*, i.e. attributes with preference ordered value sets. It is assumed that a better evaluation of an object on a criterion, with other evaluations being fixed, should not worsen its assignment to a decision class. In order to construct a preference model, one can induce it from a *reference (training)* set of objects U already assigned to decision classes. Thus, multicriteria classification problem resembles typical classification problem considered in machine learning [6,11] under monotonicity constraints: the expected decision value increases with increasing values on condition attributes. However, it still may happen that in U , there exists an object x_i not worse than another object x_k on all condition attributes, however, x_i is assigned to a worse class than x_k ; such a situation violates the

monotone nature of data, so we shall call objects x_i and x_k *inconsistent with respect to dominance principle*.

Rough set theory [13] has been adapted to deal with this kind of inconsistency and the resulting methodology has been called *Dominance-based Rough Set Approach* (DRSA) [7,8]. In DRSA, the classical indiscernibility relation has been replaced by a dominance relation. Using the rough set approach to the analysis of multicriteria classification problem, we obtain lower and upper (rough) approximations of unions of decision classes. The difference between upper and lower approximations shows inconsistent objects with respect to the dominance principle. It can happen that due to the presence of noise, the data is so inconsistent, that too much information is lost, thus making the DRSA inference model not accurate. To cope with the problem of excessive inconsistency the *variable consistency* model within DRSA has been proposed (VC-DRSA) [9].

In this paper, we look at DRSA from a different point of view, identifying its connections with statistics and statistical decision theory. Using the maximum likelihood estimation we introduce a new variable consistency variant of DRSA. It leads to the statistical problem of isotonic regression [14], which is then solved by the optimal object reassignment problem [5]. Finally, we explain the approach as being a solution to the problem of finding a decision minimizing the empirical risk [1].

Notation. We assume that we are given a set $U = \{(x_1, y_1), \dots, (x_\ell, y_\ell)\}$, consisting of ℓ training objects, with their decision values (class assignments), where each $y_i \in T$. Each object is described by a set of n condition criteria $Q = \{q_1, \dots, q_n\}$ and by $\text{dom}q_i$ we mean the set of values of attribute q_i . For each i , $\text{dom}q_i$ is ordered by some weak preference relation, here we assume for simplicity $\text{dom}q_i \subseteq \mathbb{R}$ and the order relation is a linear order \geq . We denote the evaluation of object x_i on attribute q_j by $q_j(x_i)$. Later on we will abuse a bit the notation, identifying each object x with its evaluations on all the condition criteria, $x \equiv (q_1(x), \dots, q_n(x))$ and denote $X = \{x_1, \dots, x_\ell\}$. By *class* $Cl_t \subset X$, we mean a set of objects, such that $y_i = t$, i.e. $Cl_t = \{x_i \in X: y_i = t, 1 \leq i \leq \ell\}$.

2 Classical Variable Precision Rough Set Approach

The classical rough set approach [13] (which does not take into account any monotonicity constraints) is based on the assumption that objects having the same description are indiscernible (similar) with respect to the available information [13,8]. The indiscernibility relation I is defined as:

$$I = \{(x_i, x_j) \in X \times X: q_k(x_i) = q_k(x_j) \quad \forall q_k \in Q\} \quad (1)$$

The equivalence classes of I (denoted $I(x)$ for some object $x \in X$) are called *granules*. The lower and upper approximations of class Cl_t are defined, respectively, by:

$$\underline{Cl}_t = \{x_i \in X: I(x_i) \subseteq Cl_t\} \quad \overline{Cl}_t = \bigcup_{x_i \in Cl_t} I(x_i) \quad (2)$$

For application to the real-life data, a less restrictive definition was introduced under the name of *variable precision rough set model* (VPRS) [16] and is expressed in the probabilistic terms. Let $\Pr(Cl_t|I(x))$ be a probability that an object x_i from granule $I(x)$ belongs to the class Cl_t . The probabilities are unknown, but are estimated by frequencies $\Pr(Cl_t|I(x)) = \frac{|Cl_t \cap I(x)|}{|I(x)|}$. Then, the lower approximation of class Cl_t is defined as:

$$\underline{Cl}_t = \bigcup_{I(x):x \in X} \{I(x): \Pr(Cl_t|I(x)) \geq u\} \tag{3}$$

so it is the sum of all granules, for which the probability of class Cl_t is at least equal to some threshold u .

It can be shown that frequencies used for estimating probabilities are the maximum likelihood (ML) estimators under assumption of common class probability distribution for every object within each granule. The sketch of the derivation is the following. Let us choose some granule $G = I(x)$. Let n_G be the number of objects in G , and for each class Cl_t , let n_G^t be the number of objects from this class in G . Then the decision value y has a multinomial distribution when conditioned on granule G . Let us denote those probabilities $\Pr(y = t|G)$ by p_G^t . Then, the conditional probability of observing n_G^1, \dots, n_G^t objects in G (conditional likelihood) is given by $L(p; n_G|G) = \prod_{t=1}^m (p_G^t)^{n_G^t}$, so that the log-likelihood is given by $\mathcal{L}(p; n_G|G) = \ln L(n; p, G) = \sum_{t=1}^m n_G^t \ln p_G^t$. The maximization of $\mathcal{L}(p; n_G|G)$ with additional constraint $\sum_{t=1}^m p_G^t = 1$ leads to the well-known formula for ML estimators \hat{p}_G^t in multinomial distribution:

$$\hat{p}_G^t = \frac{n_G^t}{n_G} \tag{4}$$

which are exactly the frequencies used in VPRS. This observation will lead us in section 4 to the definition of the variable consistency for dominance-based rough set approach.

3 Dominance-Based Rough Set Approach (DRSA)

Within DRSA [7,8], we define the *dominance* relation D as a binary relation on X in the following way: for any $x_i, x_k \in X$ we say that x_i *dominates* x_k , $x_i D x_k$, if on every condition criterion from Q , x_i has evaluation not worse than x_k , $q_j(x_i) \geq q_j(x_k)$, for $j = 1, \dots, n$. The dominance relation D is a partial pre-order on X , i.e. it is reflexive and transitive. The *dominance principle* can be expressed as follows:

$$x_i D x_j \implies y_i \geq y_j \tag{5}$$

for any $x_i, x_j \in X$. We say that two objects $x_i, x_j \in X$ are consistent if they satisfy the dominance principle. We say that object x_i is consistent, if it is consistent with every other object from X .

The rough approximations concern granules resulting from information carried out by the decisions. The decision granules can be expressed by upward and downward unions of decision classes, respectively:

$$Cl_t^{\geq} = \{x_i \in X : y_i \geq t\} \quad Cl_t^{\leq} = \{x_i \in X : y_i \leq t\} \quad (6)$$

The condition granules are dominating and dominated sets defined, respectively, for each $x \in X$, as:

$$D^+(x) = \{x_i \in X : x_i D x\} \quad D^-(x) = \{x_i \in X : x D x_i\} \quad (7)$$

Lower approximations of Cl_t^{\geq} and Cl_t^{\leq} are defined as:

$$\underline{Cl}_t^{\geq} = \{x_i \in X : D^+(x_i) \subseteq Cl_t^{\geq}\} \quad \underline{Cl}_t^{\leq} = \{x_i \in X : D^-(x_i) \subseteq Cl_t^{\leq}\} \quad (8)$$

Upper approximations of Cl_t^{\geq} and Cl_t^{\leq} are defined as:

$$\overline{Cl}_t^{\geq} = \{x_i \in X : D^-(x_i) \cap Cl_t^{\geq} \neq \emptyset\} \quad \overline{Cl}_t^{\leq} = \{x_i \in X : D^+(x_i) \cap Cl_t^{\leq} \neq \emptyset\} \quad (9)$$

4 Statistical Model of Variable Consistency in DRSA

In this section, we introduce a new model of variable consistency DRSA (VC-DRSA), by miming the ML estimation shown in section 2. The name *variable consistency* instead of *variable precision* is used in this chapter only to be consistent with the already existing theory [9].

In section 2, although it was not mentioned straightforward, while estimating the probabilities, we have made the assumption that in a single granule $I(x)$, each object $x \in G$ has the same conditional probability distribution, $\Pr(y = t|I(x)) \equiv p_G^t$. This is due to the property of indiscernibility of objects within a granule. In case of DRSA, indiscernibility is replaced by a dominance relation, so that a different relation between the probabilities must hold. Namely, we conclude from the dominance principle that:

$$x_i D x_j \implies p_i^t \geq p_j^t \quad \forall t \in T, \quad \forall x_i, x_j \in X \quad (10)$$

where p_i^t is a probability (conditioned on x_i) of decision value at least t , $\Pr(y \geq t|x_i)$. In other words, if object x_i dominates object x_j , probability distribution conditioned at point x_i *stochastically dominates* probability distribution conditioned at x_j . Equation (10) will be called *stochastic dominance principle*.

In this section, we will restrict the analysis to two-class (binary) problem, so we assume $T = \{0, 1\}$ (indices start with 0 for simplicity). Notice, that \underline{Cl}_0^{\geq} and \underline{Cl}_1^{\leq} are trivial, so that only \underline{Cl}_1^{\geq} and \underline{Cl}_0^{\leq} are used and will be denoted simply by \underline{Cl}_1 and \underline{Cl}_0 , respectively. We relax the definition of lower approximations for $T = \{0, 1\}$ in the following way (in analogy to the classical variable precision model):

$$\underline{Cl}_t = \{x_i \in X : p_i^t \geq \alpha\}, \quad (11)$$

where $\alpha \in (0.5, 1]$ is a chosen *consistency level*. Since we do not know probabilities p_i^t , we will use instead their ML estimators \hat{p}_i^t . The conditional likelihood function (probability of decision values with X being fixed) is a product of binomial distributions and is given by $\prod_{i=1}^{\ell} (p_i^1)^{y_i} (p_i^0)^{1-y_i}$, or using $p_i \equiv p_i^1$ (since $p_i^0 = 1 - p_i$), is given by $\prod_{i=1}^{\ell} (p_i)^{y_i} (1 - p_i)^{1-y_i}$. The log-likelihood is then

$$\mathcal{L}(p; y|X) = \sum_{i=1}^{\ell} (y_i \ln(p_i) + (1 - y_i) \ln(1 - p_i)) \tag{12}$$

The stochastic dominance principle (10) simplifies to:

$$x_i D x_j \implies p_i \geq p_j \quad \forall x_i, x_j \in X \tag{13}$$

To obtain probability estimators \hat{p}_i , we need to maximize (12) subject to constraints (13). This is exactly the problem of statistical inference under the order restriction [14]. Before investigating properties of the problem, we state the following theorem:

Theorem 1. *Object $x_i \in X$ is consistent with respect to the dominance principle if and only if $\hat{p}_i = y_i$.*

Using Theorem 1 we can set $\hat{p}_i = y_i$ for each consistent object $x_i \in X$ and optimize (12) only for inconsistent objects, which usually gives a large reduction of the problem size (number of variables). In the next section, we show that solving (12) boils down to the isotonic regression problem.

5 Isotonic Regression

For the purpose of this paper we consider the simplified version of the *isotonic regression problem* (IRP) [14]. Let $X = \{x_1, \dots, x_{\ell}\}$ be a finite set with some pre-order relation $D \subseteq X \times X$. Suppose also that $y: X \rightarrow \mathbb{R}$ is some function on X , where $y(x_i)$ is shortly denoted y_i . A function $y^*: X \rightarrow \mathbb{R}$ is an *isotonic regression* of y if it is the optimal solution to the problem:

$$\begin{aligned} & \text{minimize } \sum_{i=1}^{\ell} (y_i - p_i)^2 \\ & \text{subject to } x_i D x_j \implies p_i \geq p_j \quad \forall 1 \leq i, j \leq \ell \end{aligned} \tag{14}$$

so that it minimizes the squared error in the class of all *isotonic* functions p (where we denoted $p(x_i)$ as p_i in (14)). In our case, the ordering relation D is the dominance relation, the set X and values of function y on X , i.e. $\{y_1, \dots, y_{\ell}\}$ will have the same meaning as before. Although squared error in (14) seems to be arbitrarily chosen, it can be shown that minimizing many other error functions leads to the same function y^* as in the case of (14). Suppose that Φ is a convex function, finite on an interval I , containing the range of function y on X , i.e. $y(X) \subseteq I$ and Φ has value $+\infty$ elsewhere. Let ϕ be a nondecreasing function on

I such that, for each $u \in I$, $\phi(u)$ is a subgradient of Φ . For each $u, v \in I$ define the function $\Delta_\Phi(u, v) = \Phi(u) - \Phi(v) - (u - v)\phi(v)$. Then the following theorem holds:

Theorem 2. [14] *Let y^* be an isotonic regression of y on X , i.e. y^* solves (14). Then it holds:*

$$\sum_{x_i \in X} \Delta_\Phi(y_i, f(x_i)) \geq \sum_{x_i \in X} \Delta_\Phi(y_i, y^*(x_i)) + \sum_{x_i \in X} \Delta_\Phi(y^*(x_i), f(x_i)) \quad (15)$$

for any isotonic function f with the range in I , so that y^* minimizes

$$\sum_{x_i \in X} \Delta_\Phi(y_i, f(x_i)) \quad (16)$$

in the class of all isotonic functions f with range in I . The minimizing function is unique if Φ is strictly convex.

It was shown in [14] that by using the function:

$$\Phi(u) = \begin{cases} u \ln u + (1 - u) \ln(1 - u) & \text{for } u \in (0, 1) \\ 0 & \text{for } u \in \{0, 1\} \end{cases} \quad (17)$$

in Theorem 2, we end up with the problem of maximizing (12) subject to constraints (13). Thus, we can find solution to the problem (12) subject to (13) by solving the IRP (14).

Suppose A is a subset of X and $f: X \rightarrow \mathbb{R}$ is any function. We define $Av(f, A) = \frac{1}{|A|} \sum_{x_i \in A} f(x_i)$ to be an average of f on a set A . Now suppose y^* is the isotonic regression of y . By a *level set* of y^* , $[y^* = a]$ we mean the subset of X , on which y^* has constant value a , i.e. $[y^* = a] = \{x \in X: y^*(x) = a\}$. The following theorem holds:

Theorem 3. [14] *Suppose y^* is the isotonic regression of y . If a is any real number such that the level set $[y^* = a]$ is not empty, then $a = Av(y, [y^* = a])$.*

Theorem 3 states, that for a given x , $y^*(x)$ equal to the average of y over all the objects having the same value $y^*(x)$. Since there is a finite number of divisions of X into level sets, we conclude there are only finite number of values that y^* can possibly take. In our case, since $y_i \in \{0, 1\}$, all values of y^* must be of the form $\frac{r}{r+s}$, where r is the number of objects from class Cl_1 in the level set, while s is the number of objects from Cl_0 .

6 Minimal Reassignment Problem

In this section we briefly describe the *minimal reassignment problem* (MRP), introduced in [5]. We define the reassignment of an object $x_i \in X$ as changing its decision value y_i . Moreover, by minimal reassignment we mean reassigning the smallest possible number of objects to make the set X consistent (with respect

to the dominance principle). One can see, that such a reassignment of objects corresponds to indicating and correcting possible errors in the dataset. To find minimal reassignment, one can formulate a linear program. Such problems were already considered in [3] (under the name *isotonic separation*, in the context of binary and multi-class classification) and also in [2] (in the context of boolean regression).

Assume $y_i \in \{0, 1\}$. For each $x_i \in X$ we introduce a binary variable d_i which is to be a new decision value for x_i . The request that the new decision values must be consistent with respect to the dominance principle implies:

$$x_i D x_j \implies d_i \geq d_j \quad \forall 1 \leq i, j \leq \ell \tag{18}$$

Notice, that (18) has the form of the stochastic dominance principle (13). The reassignment of an object x_i takes place if $y_i \neq d_i$. Therefore, the number of reassigned objects (which is also the objective function for MRP) is given by $\sum_{i=1}^{\ell} |y_i - d_i| = \sum_{i=1}^{\ell} (y_i(1 - d_i) + (1 - y_i)d_i)$, where the last equality is due to the fact, that both $y_i, d_i \in \{0, 1\}$ for each i . Finally notice that the matrix of constraints (18) is totally unimodular, so we can relax the integer condition for d_i reformulating it as $0 \leq d_i \leq 1$, and get a linear program [3,12]. Moreover, constraint $0 \leq d_i \leq 1$ can be dropped, since if there were any $d_i > 1$ (or $d_i < 0$) in any feasible solution, we could decrease their values down to 1 (or increase up to 0), obtaining a new feasible solution with smaller value of the objective function. Finally, for the purpose of the paper, we rewrite the problem in the following form:

$$\begin{aligned} & \text{minimize} \quad \sum_{i=1}^{\ell} |y_i - d_i| \\ & \text{subject to} \quad x_i D x_j \implies d_i \geq d_j \quad \forall 1 \leq i, j \leq \ell \end{aligned} \tag{19}$$

Comparing (19) with (14), we notice that, although both problems emerged in different context, they look very similar and the only difference is in the objective function (L_1 -norm in MRP instead of L_2 -norm in IRP). In fact, both problems are closely connected, which will be shown in the next section.

7 Connection Between IRP and MRP

To show the connection between IRP and MRP we consider the latter to be in more general form, allowing the cost of reassignment to be different for different classes. The *weighted* minimal reassignment problem (WMRP) is given by

$$\begin{aligned} & \text{minimize} \quad \sum_{i=1}^{\ell} w_{y_i} |y_i - d_i| \\ & \text{subject to} \quad x_i D x_j \implies d_i \geq d_j \quad \forall 1 \leq i, j \leq \ell \end{aligned} \tag{20}$$

where w_{y_i} are arbitrary, positive weights associated with decision classes. The following results hold:

Theorem 4. Suppose $\hat{p} = \{\hat{p}_1, \dots, \hat{p}_\ell\}$ is an optimal solution to IRP (14). Choose some value $\alpha \in [0, 1]$ and define two functions:

$$l(p) = \begin{cases} 0 & \text{if } p \leq \alpha \\ 1 & \text{if } p > \alpha \end{cases} \quad (21)$$

and

$$u(p) = \begin{cases} 0 & \text{if } p < \alpha \\ 1 & \text{if } p \geq \alpha \end{cases} \quad (22)$$

Then the solution $\hat{d}^l = \{\hat{d}_1^l, \dots, \hat{d}_\ell^l\}$ such that $\hat{d}_i^l = l(\hat{p}_i)$ for each $i \in \{1, \dots, \ell\}$, and the solution $\hat{d}^u = \{\hat{d}_1^u, \dots, \hat{d}_\ell^u\}$ such that $\hat{d}_i^u = u(\hat{p}_i)$ for each $i \in \{1, \dots, \ell\}$, are the optimal solutions to WMRP (20) with weights:

$$\begin{aligned} w_0 &= p \\ w_1 &= 1 - p \end{aligned} \quad (23)$$

Moreover, if $\hat{d} = \{\hat{d}_1, \dots, \hat{d}_\ell\}$ is an optimal integer solution to WMRP with weights (23), it must hold $\hat{d}_i^l \leq \hat{d}_i \leq \hat{d}_i^u$, for all $i \in \{1, \dots, \ell\}$. In particular, if $\hat{d}^l \equiv \hat{d}^u$, the solution to the WMRP is unique.

Theorem 4 clearly states, that if the optimal value for a variable \hat{p}_i in IRP (14) is greater (or smaller) than α , then the optimal value for the corresponding variable \hat{d}_i in the WMRP (20) with weights (23) is 1 (or 0). In particular, for $\alpha = \frac{1}{2}$ we have $w_0 = w_1 = 1$, so we obtain MRP (19). It also follows from Theorem 4, that if α cannot be taken by any \hat{p}_i in the optimal solution \hat{p} to the IRP (14), the optimal solution to the WMRP (20) is unique. It follows from Theorem 3 (and discussion after it), that \hat{p} can take only finite number of values, which must be of the form $\frac{r}{r+s}$, where $r < \ell_1$ and $s < \ell_1$ are integer (ℓ_0 and ℓ_1 are numbers of objects from class, respectively, 0 and 1). Since it is preferred to have a unique solution to the reassignment problem, from now on, we always assume that α was chosen not to be of the form $\frac{r}{r+s}$ (in practice it can easily be done by choosing α to be some simple fraction, e.g. $2/3$ and adding some small number ϵ). We call such value of α to be *proper*.

It is worth noticing that WMRP is easier to solve than IRP. It is linear, so that one can use linear programming, it can also be transformed to the network flow problem [3] and solved in $O(n^3)$. In the next section, we show, that to obtain lower and upper approximations for the VC-DRSA, it is enough to solve IRP and solves two reassignment problems instead.

8 Summary of the Statistical Model for DRSA

We begin with reminding the definitions of lower approximations of classes (for two-class problem) for consistency level α :

$$\underline{Cl}_t = \{x_i \in X : p_i^t \geq \alpha\} \quad (24)$$

for $t \in \{0, 1\}$. The probabilities p^t are estimated using the ML approach and from the previous analysis it follows that the set of estimators \hat{p} is the optimal solution to the IRP.

As it was stated in the previous section we choose α to be proper, so that the definition (24) can be equivalently stated as:

$$\begin{aligned} \underline{Cl}_1 &= \{x_i \in X : \hat{p}_i > \alpha\} \\ \underline{Cl}_0 &= \{x_i \in X : 1 - \hat{p}_i > \alpha\} = \{x_i \in X : \hat{p}_i < 1 - \alpha\} \end{aligned} \tag{25}$$

where we replaced the probabilities by their ML estimators. It follows from Theorem 4, that to obtain \underline{Cl}_0 and \underline{Cl}_1 we do not need to solve IRP. Instead we solve two weighted minimal reassignment problems (20), first one with weights $w_0 = \alpha$ and $w_1 = 1 - \alpha$, second one with $w_0 = 1 - \alpha$ and $w_1 = \alpha$. Then, objects with new decision value (optimal assignment) $\hat{d}_i = 1$ in the first problem form \underline{Cl}_1 , while objects with new decision value $\hat{d}_i = 0$ in the second problem form \underline{Cl}_0 . It is easy to show that the boundary between classes (defined as $X - (\underline{Cl}_1 \cup \underline{Cl}_0)$) is composed of objects, for which new decision values are different in those two problems.

9 Extension to the Multi-class Case

Till now, we focused on binary classification problems considered within DRSA. Here we show, how to solve the general problem with m decision classes.

We proceed as follows. We divide the problem into $m - 1$ binary problems. In t th binary problem, we estimate the lower approximations of upward union for class $t+1$, $\underline{Cl}_{t+1}^{\geq}$, and the lower approximation of downward union for class t , \underline{Cl}_t^{\leq} using the theory stated in the section 8 for two-class problem with $Cl_0 = Cl_t^{\leq}$ and $Cl_1 = Cl_{t+1}^{\geq}$. Notice, that for the procedure to be consistent, it must hold if $t' > t$ than $\underline{Cl}_{t'}^{\geq} \subseteq \underline{Cl}_t^{\geq}$ and $\underline{Cl}_t^{\leq} \subseteq \underline{Cl}_{t'}^{\leq}$. In other words, the solution has to satisfy the property of inclusion that is one of the fundamental properties considered in rough set theory. Fortunately, we have:

Theorem 5. *For each $t = 1, \dots, m - 1$, let \underline{Cl}_t^{\leq} and $\underline{Cl}_{t+1}^{\geq}$ be the sets obtained from solving two-class isotonic regression problem with consistency level α for binary classes $Cl_0 = Cl_t^{\leq}$ and $Cl_1 = Cl_{t+1}^{\geq}$. Then, we have:*

$$t' \geq t \implies \underline{Cl}_t^{\leq} \subseteq \underline{Cl}_{t'}^{\leq} \tag{26}$$

$$t' \geq t \implies \underline{Cl}_{t'+1}^{\geq} \subseteq \underline{Cl}_{t+1}^{\geq} \tag{27}$$

10 Decision-Theoretical View

In this section we look at the problem of VPRS and VC-DRSA from the point of view of statistical decision theory [1,11]. A decision-theoretic approach has already been proposed in [15] (for VPRS) and in [10] (for DRSA). The theory

presented here for VPRS is slightly different than in [15], while the decision-theoretic view for DRSA proposed in this section is completely novel.

Suppose, we seek for a function (classifier) $f(x)$ which, for a given input vector x , predicts value y as well as possible. To assess the goodness of prediction, the loss function $L(f(x), y)$ is introduced for penalizing the prediction error. Since x and y are random variables, the overall measure of the classifier $f(x)$ is the expected loss or risk, which is defined as a functional:

$$R(f) = E[L(y, f(x))] = \int L(y, f(x))dP(y, x) \tag{28}$$

for some probability measure $P(y, x)$. Since $P(y, x)$ is unknown in almost all the cases, one usually minimize the empirical risk, which is the value of risk taken for the points from a training sample U :

$$R_e(f) = \sum_{i=1}^{\ell} L(y_i, f(x_i)). \tag{29}$$

Function f is usually chosen from some restricted family of functions. We now show that the rough set theory leads to the classification procedures, which are naturally suited for dealing with problems when the classifiers are allowed to abstain from giving answer in some cases.

Let us start with VPRS. Assume, that we allow the classifier to give no answer, which is denoted as $f(x) = ?$. The loss function suitable for the problem is the following:

$$L_c(f(x), y) = \begin{cases} 0 & \text{if } f(x) = y \\ 1 & \text{if } f(x) \neq y \\ a & \text{if } f(x) = ? \end{cases} \tag{30}$$

There is a penalty a for giving no answer. To be consistent with the classical rough set theory, we assume, that any function must be constant within each granule, i.e. for each $G = I(x)$ for some $x \in X$, we have:

$$x_i, x_j \in G \implies f(x_i) = f(x_j) \quad \forall x_i, x_j \in X \tag{31}$$

which is in fact the principle of indiscernibility. We now state:

Theorem 6. *The function f^* minimizing the empirical risk (29) with loss function (30) between all functions satisfying (31) is equivalent to the VPRS in the sense, that $f^*(G) = t$ if and only if granule G belongs to the lower approximation of class t with the precision threshold $u = 1 - a$, otherwise $f^*(G) = ?$.*

Concluding, the VPRS can be derived by considering the class of functions constant in each granule and choosing the function f^* , which minimizes the empirical risk (29) for loss function (30) with parameter $a = 1 - u$. As we see, classical rough set theory suits well for considering the problems when the classification procedure is allowed not to give predictions for some x .

We now turn back to DRSA. Assume, that to each point x , the classifier f assigns the interval of classes, denoted $[l(x), u(x)]$. The lower and upper ends of each interval are supposed to be consistent with the dominance principle:

$$\begin{aligned} x_i D x_j &\implies l(x_i) \geq l(x_j) && \forall x_i, x_j \in X \\ x_i D x_j &\implies u(x_i) \geq u(x_j) && \forall x_i, x_j \in X \end{aligned} \tag{32}$$

The loss function $L(f(x), y)$ is composed of two terms. First term is a penalty for the size of the interval (degree of imprecision) and equals to $a(u(x) - l(x))$. Second term measures the accuracy of the classification and is zero, if $y \in [l(x), u(x)]$, otherwise $f(x)$ suffers additional loss equal to distance of y from the closer interval range:

$$L(f(x), y) = a(u(x) - l(x)) + I(y \notin [l(x), u(x)]) \min\{|y - l(x)|, |y - u(x)|\} \tag{33}$$

where $I(\cdot)$ is an indicator function. We now state:

Theorem 7. *The function f^* minimizing the empirical risk (29) with loss function (33) between all interval functions satisfying (32) is equivalent to the statistical VC-DRSA with consistency level $\alpha = 1 - a$ in the sense, that for each $x \in X$, $x \in \underline{Cl}_t^{\geq}$ or $x \in \underline{Cl}_t^{\leq}$ if and only if $t \in f^*(x)$.*

Concluding, the statistical VC-DRSA, can be derived by considering the class of interval functions, for which the lower and upper ends of interval are isotonic (consistent with the dominance principle) and choosing the function f^* , which minimizes the empirical risk (29) with loss function (33) with parameter $a = 1 - \alpha$.

11 Conclusions

The paper introduced a new variable consistency theory for Dominance-based Rough Set Approach. Starting from the general remarks about the estimation of probabilities in the classical rough set approach (which appears to be maximum likelihood estimation), we used the same statistical procedure for DRSA, which led us to the isotonic regression problem. The connection between isotonic regression and minimal reassignment solutions was considered and it was shown that in the case of the new variable consistency model, it is enough to solve minimal reassignment problem (which is linear), instead of the isotonic regression problem (quadratic). The approach has also been extended to the multi-class case by solving $m - 1$ binary subproblems for the class unions. The proposed theory has an advantage of basing on well investigated maximum likelihood estimation method – its formulation is clear and simple, it unites seemingly different approaches for classical and dominance-based case.

Finally notice that a connection was established between statistical decision theory and rough set approach. It follows from the analysis that rough set theory can serve as a tools for constructing classifiers, which can abstain from assigning

a new object to a class in case of doubt (in classical case) or give imprecise prediction in the form of interval of decision values (in DRSA case). However, rough set theory itself has a rather small generalization capacity, due to its nonparametric character, which was shown in section 10. The plans for further research are to investigate some restricted classes of functions which would allow to apply rough set theory directly for classification.

References

1. Berger, J.: *Statistical Decision Theory and Bayesian Analysis*. Springer, New York (1993)
2. Boros, E., Hammer, P.L., Hooker, J.N.: Boolean regression. *Annals of Operations Research* 58, 3 (1995)
3. Chandrasekaran, R., Ryu, Y.U., Jacob, V., Hong, S.: Isotonic separation. *INFORMS J. Comput.* 17, 462–474 (2005)
4. Dembczyński, K., Greco, S., Kotłowski, W., Słowiński, R.: Quality of Rough Approximation in Multi-Criteria Classification Problems. In: Greco, S., Hata, Y., Hirano, S., Inuiguchi, M., Miyamoto, S., Nguyen, H.S., Słowiński, R. (eds.) *RSCTC 2006*. LNCS (LNAI), vol. 4259, pp. 318–327. Springer, Heidelberg (2006)
5. Dembczyński, K., Greco, S., Kotłowski, W., Słowiński, R.: Optimized Generalized Decision in Dominance-based Rough Set Approach. LNCS. Springer, Heidelberg (2007)
6. Duda, R., Hart, P.: *Pattern Classification*. Wiley-Interscience, New York (2000)
7. Greco, S., Matarazzo, B., Słowiński, R.: Rough approximation of a preference relation by dominance relations. *European Journal of Operational Research* 117, 63–83 (1999)
8. Greco, S., Matarazzo, B., Słowiński, R.: Rough sets theory for multicriteria decision analysis. *European Journal of Operational Research* 129(1), 1–47 (2001)
9. Greco, S., Matarazzo, B., Słowiński, R., Stefanowski, J.: In: Ziarko, W., Yao, Y. (eds.) *RSCTC 2000*. LNCS (LNAI), vol. 2005, pp. 170–181. Springer, Heidelberg (2001)
10. Greco, S., Słowiński, R., Yao, Y.: Bayesian Decision Theory for Dominance-based Rough Set Approach. *Lecture Notes in Computer Science* 4481, 134–141 (2007)
11. Hastie, T., Tibshirani, R., Friedman, J.: *The Elements of Statistical Learning*. Springer, Heidelberg (2003)
12. Papadimitriou, C.H., Steiglitz, K.: *Combinatorial Optimization*. Dover Publications, New York (1998)
13. Pawlak, Z.: Rough sets. *International Journal of Information & Computer Sciences* 11, 341–356 (1982)
14. Robertson, T., Wright, F.T., Dykstra, R.L.: *Order Restricted Statistical Inference*. John Wiley & Sons, Chichester (1998)
15. Yao, Y., Wong, S.: A decision theoretic Framework for approximating concepts. *International Journal of Man-machine Studies* 37(6), 793–809 (1992)
16. Ziarko, W.: Probabilistic Rough Sets. In: Ślęzak, D., Wang, G., Szczuka, M., Düntsch, I., Yao, Y. (eds.) *RSFDGrC 2005*. LNCS (LNAI), vol. 3641, pp. 283–293. Springer, Heidelberg (2005)