
2 Classification Using Optimization: Application 3 to Credit Ratings of Bonds

4 Vladimir Bugera, Stan Uryasev, and Grigory Zrazhevsky

5 University of Florida, ISE, Risk Management and Financial Engineering Lab,
6 uryasev@ufl.edu

7 **Summary.** The classification approach, previously considered in credit card scor-
8 ing, is extended to multi class classification in application to credit rating of bonds.
9 The classification problem is formulated as minimization of a penalty constructed
10 with quadratic separating functions. The optimization is reduced to a linear pro-
11 gramming problem for finding optimal coefficients of the separating functions. Var-
12 ious model constraints are considered to adjust model flexibility and to avoid data
13 overfitting. The classification procedure includes two phases. In phase one, the clas-
14 sification rules are developed based on “in-sample” dataset. In phase two, the clas-
15 sification rules are validated with “out-of-sample” dataset. The considered methodol-
16 ogy has several advantages including simplicity in implementation and classification
17 robustness. The algorithm can be applied to small and large datasets. Although the
18 approach was validated with a finance application, it is quite general and can be
19 applied in other engineering areas.

20 **Key words:** Bond ratings, credit risk, classification.

21 1 Introduction

22 We consider a general approach for classifying objects into several classes and
23 applies it to a bond-rating problem. Classification problems become increas-
24 ingly important in the decision science. In finance, they are used for grouping
25 financial instruments according to their risk, or profitability characteristics.
26 In the bonds rating problem, for example, the debt instruments are arranged
27 according to the likelihood of debt issuer to default on the defined obligation.

28 Mangasarian et al. (1995) used a *utility function* for the failure discrim-
29 inant analysis (applications to breast cancer diagnosis). The utility function
30 was considered to be linear in control parameters and indicator variables and
31 it was found by minimizing the error of misclassification. Zopounidis and
32 Doumpos (1997), Zopounidis et al. (1998) and Pardalos et al. (1997) used
33 linear utility functions for trichotomous classifications of credit card applica-
34 tions. Konno and Kobayashi (2000) and Konno et al. (2000) considered utility

1 functions, quadratic in indicator parameters and linear in decision parame-
2 ters. The approach was tested with the classification of enterprises and breast
3 cancer diagnosis. Konno and Kobayashi (2000) and Konno et al. (2000) im-
4 posed convexity constraints on utility functions in order to avoid discontinuity
5 of discriminant regions. Similar to Konno and Kobayashi (2000) and Konno
6 et al. (2000), Bugera et al. (2003) applied a quadratic utility function to tri-
7 chotomous classification, but instead of convexity constraints, monotonicity
8 constraints reflecting experts' opinions were used. The approach by Bugera
9 et al. (2003) is closely related to ideas by Zopounidis and Doumpos (1997),
10 Zopounidis et al. (1998) and Pardalos et al. (1997); it considers a multi-class
11 classification with several levelsets of the utility function, where every levelset
12 corresponds to a separate class.

13 We extend the classification approach considered by Bugera et al. (2003).
14 Several innovations improving the efficiency of the algorithm are suggested:

- 15 • A *set* of utility functions, called *separating functions*, is used for classifi-
16 cation. The separating functions are quadratic in indicator variables and
17 linear in decision variables. A set of optimal separating functions is found
18 by minimizing the misclassification error. The problem is formulated as a
19 linear programming problem w.r.t. decision variables.
- 20 • To control flexibility of the model and avoid overfitting we impose various
21 *new constraints* on the separating functions.

22 Controlling flexibility of the model with constraints is crucially important
23 for the suggested approach. Quadratic separating functions (depending upon
24 the problem dimension) may have a very large number of free parameters.
25 Therefore, a tremendously large dataset may be needed to “saturate” the
26 model with data. Constraints reduce the number of degrees of freedom of the
27 model and adjust “flexibility” of the model to the size of the dataset.

28 This paper is focused on a numerical validation of the proposed algo-
29 rithm. We rated a set of international bonds using the proposed algorithm.
30 The dataset for the case study was provided by the research group of the
31 RiskSolutions branch of Standard and Poor's, Inc. We investigated the im-
32 pact of model flexibility on classification characteristics of the algorithm and
33 compared performance of several models with different types of constraints.
34 Experiments showed the importance of constraints adjusting the flexibility of
35 the model. We studied “in-sample” and “out-of-sample” characteristics of the
36 suggested algorithm. At the first stage of the algorithm, we minimized the em-
37 pirical risk, that is, the error of misclassification on a training set (in-sample
38 error). However, the real objective of the algorithm is to classify objects out-
39 side the training set with a minimal error (out-of-sample error). The in-sample
40 error is always not greater than the out-of-sample error. Similar issues were
41 studied in the Statistical Learning Theory (Vapnik, 1998). For validation, we
42 used the leave-one-out cross-validation scheme. This technique provides the
43 highest confidence level while effectively improving the predicting power of
44 the model.

1 The advantage of the considered methodology is its simplicity and consis-
 2 tency if compared to the proprietary models which are based on the combi-
 3 nation of various techniques, such as expert models, neural networks, classifi-
 4 cation trees etc.

5 The paper is organized as follows. Section 2 provides a general descrip-
 6 tion of the approach. Section 3 considers constraints applied to the studied
 7 optimization problem. Section 4 discusses techniques for choosing model flex-
 8 ibility. Section 5 explains how the errors estimation was done for the models.
 9 Section 6 discusses the bond-rating problem used for testing the methodology.
 10 Section 7 describes the datasets used in the study. Section 8 presents the re-
 11 sults of computational experiments and analyses obtained results. We finalize
 12 the paper with concluding remarks in Section 9.

13 2 Description of Methodology

14 The *object space* is a set of elements (objects) to be classified. Each element
 15 of the space has n quantitative characteristics describing properties of the
 16 considered objects.

17 We represent objects by n -dimensional vectors and the object space by an
 18 n -dimensional set $\Psi \subset \mathbb{R}^n$. The classification problem assigns elements of the
 19 object space Ψ to several classes so that each class consists of elements with
 20 similar properties. The classification is based on available prior information.
 21 In our approach, the prior information is provided by set of objects $S =$
 22 $\{\mathbf{x}_1, \dots, \mathbf{x}_m\}$ with known classification (*in-sample dataset*). The purpose of the
 23 methodology is to develop a classification algorithm that assigns a class to a
 24 new object based on the in-sample information.

25 Let us consider a classification problem with objects having n character-
 26 istics and J classes. Since each object is represented by a vector in a multi-
 27 dimensional space \mathbb{R}^n , the classification can be defined by an integer-valued
 28 function $f_0(\mathbf{x})$, $\mathbf{x} \in \mathbb{R}^n$. The value of the function defines the class of an ob-
 29 ject. We call $f_0(\mathbf{x})$ a *classification function*. This function splits the object
 30 space into J non-intersecting areas:

$$\mathbb{R}^n = \bigcup_{i=1}^J F_i, F_i \cap F_j = \emptyset, F_i \neq \emptyset, F_j \neq \emptyset, i \neq j \quad , \quad (1)$$

31 where each area F_i consists of elements belonging to the corresponding class
 32 i :

$$F_i = \{ \mathbf{x} \in \mathbb{R}^n \mid f_0(\mathbf{x}) = i \} \quad . \quad (2)$$

33 We can approximate the classification function $f_0(\mathbf{x})$ using optimization
 34 methods. Let $F(\mathbf{x})$ be a cumulative distribution function of objects in the
 35 object space \mathbb{R}^n , and Λ be a parameterized set of discrete-value approxim-
 36 ating functions. Then, the function $f_0(\mathbf{x})$ can be approximated by solving the
 37 following minimization problem:

$$\min_{f \in \Lambda} \int_{R^n} Q(f(\mathbf{x}) - f_0(\mathbf{x})) dF(\mathbf{x}), \quad (3)$$

1 where $Q(f(\mathbf{x}) - f_0(\mathbf{x}))$ is a penalty function defining the value of misclassification for a single object \mathbf{x} . The optimal solution $f(\mathbf{x})$ of optimization
 2 problem (3) gives an approximation of the function $f_0(\mathbf{x})$. The main difficulty
 3 in solving problem (3) is the discontinuity of functions $f(\mathbf{x})$ and $f_0(\mathbf{x})$
 4 leading to non-convex optimization. To circumvent this difficulty, we reformulate
 5 problem (3) in a convex optimization setting.

6 Let us consider a classification function $\bar{f}(\mathbf{x})$ defining a classification
 7 on the object space R^n . Suppose we have a set of continuous functions
 8 $U_0(\mathbf{x}), \dots, U_J(\mathbf{x})$. We call them *separating functions for the classification function*
 9 $\bar{f}(\mathbf{x})$ if for every object \mathbf{x}^* from class $i = \bar{f}(\mathbf{x}^*)$ values of functions with
 10 numbers lower than i are positive:
 11

$$U_0(\mathbf{x}^*), \dots, U_{i-1}(\mathbf{x}^*) > 0; \quad (4)$$

12 values of functions with numbers higher or equal to i are negative or zeros:

$$U_i(\mathbf{x}^*), \dots, U_J(\mathbf{x}^*) \leq 0. \quad (5)$$

13 If we know the functions $U_0(\mathbf{x}), \dots, U_J(\mathbf{x})$ we can classify objects according
 14 to the following rule:

$$\left\{ \begin{array}{l} \forall p = 0, \dots, i-1 : U_p(\mathbf{x}^*) > 0 \\ \forall q = i, \dots, J : U_q(\mathbf{x}^*) \leq 0 \end{array} \right\} \Leftrightarrow \{ \bar{f}(\mathbf{x}^*) = i \}. \quad (6)$$

15 The class number of an object \mathbf{x}^* can be determined by the number of
 16 positive values of the functions: an object in i class has exactly i positive
 17 separating functions. Figure 1 illustrates this property.

18 Suppose we can represent any classification function $f(\mathbf{x})$ from a parameterized
 19 set of functions Λ by a set of separating functions. By constructing a
 20 penalty for the classification $f_0(\mathbf{x})$ using separating functions, we can formulate
 21 optimization problem (3) with respect to the parameters of the separating
 22 functions.

23 Suppose the classification function $\bar{f}(\mathbf{x})$ is defined by a set of separating
 24 functions $U_0(\mathbf{x}), \dots, U_J(\mathbf{x})$. We denote by $D_{f_0}(U_0, \dots, U_J)$ an integral function
 25 measuring deviation of the classification (implied by the separating functions
 26 $U_0(\mathbf{x}), \dots, U_J(\mathbf{x})$) from the true classification defined by $f_0(\mathbf{x})$:

$$D_{f_0}(U_0, \dots, U_J) = \int_{R^n} Q(U_0(\mathbf{x}), \dots, U_J(\mathbf{x}), f_0(\mathbf{x})) dF(\mathbf{x}) \quad , \quad (7)$$

27 where $Q(U_0(\mathbf{x}), \dots, U_J(\mathbf{x}), f_0(\mathbf{x}))$ is a penalty function.

28 Further, we consider the following penalty function:

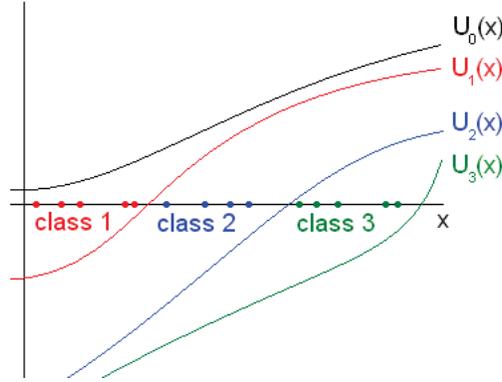


Fig. 1. Classification by separating functions. At each point (object), values of the separating functions with the indices lower than i are positive; the values of the remaining functions are negative or zero. Class of an object is determined by the number of functions with positive values.

$$Q(U_0, \dots, U_J, f_0(\mathbf{x})) = \sum_{k=0}^{f_0(\mathbf{x})-1} \lambda_k (-U_k(\mathbf{x}))^+ + \sum_{k=f_0(\mathbf{x})}^J \lambda_k (U_k(\mathbf{x}))^+, \quad (8)$$

1 where $(y)^+ = \max(0, y)$, and $\lambda_k, k = 0, \dots, J$ are positive parameters. The
 2 penalty function equals zero if $U_k(\mathbf{x}) \geq 0$ for $k = 0, \dots, f_0(\mathbf{x}) - 1$, and $U_k(\mathbf{x}) \leq$
 3 0 for $k = f_0(\mathbf{x}), \dots, J$. If the classification defined by the separating functions
 4 $U_0(\mathbf{x}), \dots, U_J(\mathbf{x})$ according to rule (6) coincides with the true classification
 5 $f_0(\mathbf{x})$, then, the value of penalty function (8) equals zero for any object \mathbf{x} . If
 6 the penalty is positive for an object \mathbf{x} then separating functions misclassify this
 7 object. Therefore, the penalty function defined by (8) provides the condition
 8 of the correct classification by the separating functions:

$$\begin{aligned} & \{Q(U_0, \dots, U_J, f_0(\mathbf{x})) = 0\} \\ & \quad \updownarrow \\ & \{U_k(\mathbf{x}) > 0 > U_l(\mathbf{x}), \forall k = 1, \dots, f_0(\mathbf{x}) - 1, \forall l = f_0(\mathbf{x}), \dots, J\}. \end{aligned} \quad (9)$$

9 The choice of penalty function (8) is motivated by the possibility of building
 10 an efficient algorithm for the optimization of integral (7) when the separa-
 11 rating functions are linear w.r.t. control parameters. In this case, optimization
 12 problem (3) can be reduced to linear programming.

13 After introducing the separating functions we reformulate optimization
 14 problem (3) as follows:

$$\min_{U_0, \dots, U_J} D_{f_0}(U_0, \dots, U_J). \quad (10)$$

1 This optimization problem finds optimal separating functions. We can ap-
 2 proximate the function $D_{f_0}(U_0, \dots, U_J)$ by sampling objects according to the
 3 cumulative distribution function $F(\mathbf{x})$. Assuming that $\mathbf{x}^1, \dots, \mathbf{x}^m$ are some sam-
 4 ple points, the approximation of $D_{f_0}(U_0, \dots, U_J)$ is given by:

$$\tilde{D}_{f_0}^m(U_0, \dots, U_J) = \frac{1}{m} \sum_{i=1}^m Q(U_0(\mathbf{x}^i), \dots, U_J(\mathbf{x}^i), f_0(\mathbf{x}^i)). \quad (11)$$

5 Therefore, the approximation of deviation function (7) becomes:

$$\begin{aligned} \tilde{D}_{f_0}^m(U_0, \dots, U_J) &= \\ &= \frac{1}{m} \sum_{i=1}^m \left(\sum_{k=0}^{f_0(\mathbf{x}^i)-1} \lambda_k (-U_k(\mathbf{x}^i))^+ + \sum_{k=f_0(\mathbf{x}^i)}^J \lambda_k (U_k(\mathbf{x}^i))^+ \right). \end{aligned} \quad (12)$$

6 To avoid possible ambiguity when the value of a separating function equals
 7 0, we introduced a small positive constant δ inside of each term in the penalty
 8 function (constant δ has to be chosen small enough in order not to have a
 9 significant impact on the final classification):

$$\begin{aligned} Q(U_0, \dots, U_J, f_0(\mathbf{x})) &= \\ &= \sum_{k=0}^{f_0(\mathbf{x})-1} \lambda_k (-U_k(\mathbf{x}) + \delta)^+ + \sum_{k=f_0(\mathbf{x})}^J \lambda_k (U_k(\mathbf{x}) + \delta)^+. \end{aligned} \quad (13)$$

10 The penalty function equals zero if $U_k(\mathbf{x}) \geq \delta$ for $k = 0, \dots, f_0(\mathbf{x}) - 1$, and
 11 $U_k(\mathbf{x}) \leq -\delta$ for $k = f_0(\mathbf{x}), \dots, J$. The approximation of deviation function (7)
 12 becomes

$$\begin{aligned} \tilde{D}_{f_0}^m(U_0, \dots, U_J) &= \\ &= \frac{1}{m} \sum_{i=1}^m \left(\sum_{k=0}^{f_0(\mathbf{x}^i)-1} \lambda_k (-U_k(\mathbf{x}^i) + \delta)^+ + \sum_{k=f_0(\mathbf{x}^i)}^J \lambda_k (U_k(\mathbf{x}^i) + \delta)^+ \right). \end{aligned} \quad (14)$$

13 Further, we parameterized each separating function by a K -dimensional
 14 vector $\boldsymbol{\alpha} \in A \subset \mathbb{R}^K$. Parameter K is defined by the design of the separ-
 15 ating functions. Therefore, a set of separating functions $U_{(\boldsymbol{\alpha}^0, \dots, \boldsymbol{\alpha}^J)}(\mathbf{x}) =$
 16 $\{U_{\boldsymbol{\alpha}^0}, \dots, U_{\boldsymbol{\alpha}^J}\}$ is determined by a set of vectors $\{\boldsymbol{\alpha}^0, \dots, \boldsymbol{\alpha}^J\}$. With this param-
 17 eterization we reformulated the problem (10) as minimization of the convex
 18 piece-wise linear functions:

$$\min_{\boldsymbol{\alpha}^0, \dots, \boldsymbol{\alpha}^J \in A} \sum_{i=1}^m \left(\sum_{k=0}^{f_0(\mathbf{x}^i)-1} \lambda_k (-U_{\boldsymbol{\alpha}^k}(\mathbf{x}^i) + \delta)^+ + \sum_{k=f_0(\mathbf{x}^i)}^J \lambda_k (U_{\boldsymbol{\alpha}^k}(\mathbf{x}^i) + \delta)^+ \right), \quad (15)$$

19 where $A \subset \mathbb{R}^K$. By introducing new variables σ_i^j , we reduced problem (15)
 20 to equivalent mathematical programming problem with the linear objective
 21 function:

$$\begin{aligned}
 & \min_{\alpha, \sigma} \sum_{i=1}^m \sum_{j=1}^J \lambda_j \sigma_i^j \\
 & \sigma_i^j \geq -U_{\alpha^j}(\mathbf{x}^i) + \delta, j = 0, \dots, f_0(\mathbf{x}^i) - 1 \\
 & \sigma_i^j \geq U_{\alpha^j}(\mathbf{x}^i) + \delta, j = f_0(\mathbf{x}^i), \dots, J \\
 & \alpha^0, \dots, \alpha^J \in A \\
 & \sigma_i^1, \dots, \sigma_i^J \geq 0
 \end{aligned} \tag{16}$$

1 Further, we consider that the separating functions are linear in control
 2 parameters $\alpha_1, \dots, \alpha_K$. In this case the separating functions can be represented
 3 in the following form:

$$U_{\alpha}(\mathbf{x}) = \sum_{k=1}^K \alpha_k g_k(\mathbf{x}). \tag{17}$$

4 In this case, optimization problem (16) for finding optimal separating func-
 5 tions can be reduced to the following linear programming problem:

$$\begin{aligned}
 & \min_{\alpha, \sigma} \sum_{i=1}^m \sum_{j=1}^J \lambda_j \sigma_i^j \\
 & \sigma_i^j + \sum_{k=1}^K \alpha_k^j g_k(\mathbf{x}^i) \geq \delta, j = 0, \dots, f_0(\mathbf{x}^i) - 1 \\
 & \sigma_i^j - \sum_{k=1}^K \alpha_k^j g_k(\mathbf{x}^i) \geq \delta, j = f_0(\mathbf{x}^i), \dots, J \\
 & \alpha^0, \dots, \alpha^J \in A \\
 & \sigma_i^1, \dots, \sigma_i^J \geq 0
 \end{aligned} \tag{18}$$

6 Further, we consider quadratic (in indicator variables) separating func-
 7 tions:

$$U_j(\mathbf{x}) = \sum_{k=1}^n \sum_{l=1}^n a_{kl}^j x_k x_l + \sum_{k=1}^n b_k^j x_k + c^j, j = 0, \dots, J. \tag{19}$$

8 Optimization problem (18) with quadratic separating functions is refor-
 9 mulated as follows:

$$\begin{aligned}
 & \min_{a, b, c, \sigma} \sum_{i=1}^m \sum_{j=1}^J \lambda_j \sigma_i^j \\
 & \sigma_i^j + \sum_{k=1}^n \sum_{l=1}^n a_{kl}^j x_k^i x_l^i + \sum_{k=1}^n b_k^j x_k^i + c^j \geq \delta, j = 0, \dots, f_0(\mathbf{x}^i) - 1 \\
 & \sigma_i^j - \sum_{k=1}^n \sum_{l=1}^n a_{kl}^j x_k^i x_l^i - \sum_{k=1}^n b_k^j x_k^i - c^j \geq \delta, j = f_0(\mathbf{x}^i), \dots, J \\
 & \sigma_i^1, \dots, \sigma_i^J \geq 0
 \end{aligned} \tag{20}$$

Although there are $J + 1$ separating functions in problem (20), only $J - 1$
 functions are essential for the classification. The functions

$$\sum_{k=1}^n \sum_{l=1}^n a_{kl}^0 x_k x_l + \sum_{k=1}^n b_k^0 x_k + c^0$$

and

$$\sum_{k=1}^n \sum_{l=1}^n a_{kl}^J x_k x_l + \sum_{k=1}^n b_k^J x_k + c^J$$

are boundary functions. For all the classified objects, the value of the first boundary function is positive, and the value of the second boundary function is negative. This can be easily achieved by setting $c^0 = M$ and $c^J = -M$, where M is a sufficiently large number. Thus, these functions can be removed from optimization problem (20). However, these boundary functions can be used for adjusting flexibility of the classification model. In the next section, we will show how to use these functions for imposing the so-called "squeezing" constraints.

For the case $J = 2$ with only two classes, problem (20) finds a quadratic surface $\sum_{k=1}^n \sum_{l=1}^n a_{kl}^1 x_k x_l + \sum_{k=1}^n b_k^1 x_k + c^1 = 0$ dividing the object space R^n into two areas. After solving optimization problem (20) we expect that a majority of objects from the first class will belong to the first area. On these points the function $\sum_{k=1}^n \sum_{l=1}^n a_{kl}^1 x_k x_l + \sum_{k=1}^n b_k^1 x_k + c^1$ is positive. Similar, for a majority of objects from the second class the function $\sum_{k=1}^n \sum_{l=1}^n a_{kl}^1 x_k x_l + \sum_{k=1}^n b_k^1 x_k + c^1$ is negative.

For the case with $J > 2$, the geometrical interpretation of optimization problem (20) refers to the partition of the object space R^n into J areas by $J - 1$ non-intersecting quadratic surfaces

$$\sum_{k=1}^n \sum_{l=1}^n a_{kl}^j x_k x_l + \sum_{k=1}^n b_k^j x_k + c^j = 0, \quad j = 1, \dots, J - 1.$$

Additional feasibility constraints, that assure non-intersection of the surfaces, will be discussed in the following section.

3 Constraints

The considered separating functions, especially the quadratic functions, may be too flexible (have too many degrees of freedom) for datasets with a small number of datapoints. Imposing additional constraints may reduce excessive model flexibility. In this section, we will discuss different types of constraints applied to the model.

Konno and Kobayashi (2000) and Konno et al. (2000) considered convexity constraints on indicator variables of a quadratic utility function. Bugera et al. (2003) imposed monotonicity constraints on the model to incorporate expert preferences.

Constraints play a crucial role in developing the classification model because they reduce excessive flexibility of a model for small training datasets.

1 Moreover, a classification with multiple separating functions may not be possible for the majority of objects if appropriate constraints are not imposed.

3 3.1 Feasibility Constraints (F-Constraints)

4 For classification with multiple separating functions we may potentially come to a possible intersection of separating surfaces. This may lead to inability of the approach to classify some objects. To circumvent this difficulty, we introduce feasibility constraints, that keep the separating functions apart from each other. It makes possible to classify any new point by rule (6). In general, these constraints have the form:

$$U_j(\mathbf{x}) \leq U_{j-1}(\mathbf{x}); j = 1, \dots, J; \mathbf{x} \in W \subset R^n, \quad (21)$$

10 where W is a set on which we want to achieve the feasibility. We do not consider $W = R^n$, because it leads to “parallel” separating surfaces, which were studied in the previous work by Bugera et al. (2003). In this section, we consider a set W with a finite number of elements. In particular, we consider the set W being the training set plus the set of objects we want to classify (out-of-sample dataset). In this case, constraints (21) can be rewritten as

$$U_j(\mathbf{x}_i) \leq U_{j-1}(\mathbf{x}_i); j = 1, \dots, J; i = 1, \dots, r, \quad (22)$$

16 where r is a number of objects in the set W . The fact that we use out-of-sample points does not cause any problems because we use the data without knowing their class numbers.

19 Since classification without feasibility constraints may lead to inability of classifying new objects (especially for small training sets), we will always include the feasibility constraints (22) to classification problem (16). For quadratic separating functions, classification problem (20) with feasibility constraints can be rewritten as

$$\begin{aligned} & \min_{a,b,c,\sigma} \sum_{i=1}^m \sum_{j=1}^J \lambda_j \sigma_i^j \\ & \sigma_i^j + \sum_{k=1}^n \sum_{l=1}^n a_{kl}^j x_k^i x_l^i + \sum_{i=1}^n b_k^j x_k^i + c^j \geq \delta, j = 0, \dots, f_0(\mathbf{x}^i) - 1 \\ & \sigma_i^j - \sum_{k=1}^n \sum_{l=1}^n a_{kl}^j x_k^i x_l^i - \sum_{k=1}^n b_k^j x_k^i - c^j \geq \delta, j = f_0(\mathbf{x}^i), \dots, J \\ & \sum_{k=1}^n \sum_{l=1}^n a_{kl}^{j-1} x_k^i x_l^i + \sum_{i=1}^n b_k^{j-1} x_k^i + c^{j-1} - \sum_{k=1}^n \sum_{l=1}^n a_{kl}^j x_k^i x_l^i + \sum_{i=1}^n b_k^j x_k^i + c^j \geq \delta \\ & i = 1, \dots, r, j = 1, \dots, J \\ & \sigma_i^1, \dots, \sigma_i^J \geq 0 \end{aligned} \quad (23)$$

24 3.2 Monotonicity Constraints (M-Constraints)

25 We use monotonicity constraints to incorporate the preference of greater values of indicator variables. In financial applications, monotonicity with respect

1 to some indicators follows from “engineering” considerations. For instance, in
 2 the bond rating problem, considered in this paper, we have indicators (return
 3 on capital, operating-cash-flow-to-debt ratio, total assets), greater values of
 4 which lead to the higher ratings of bonds. If we enforce the monotonicity of
 5 the separating functions with respect to these indicators, objects with greater
 6 values of the indicators will have higher ratings. For a smooth function $h(\mathbf{x})$,
 7 the monotonicity constraints can be written as the constraint on the non-
 8 negativity of the first partial derivative:

$$\frac{\partial h(\mathbf{x})}{\partial x_i} \geq 0, \quad i \in \{1, \dots, n\}. \quad (24)$$

9 For the case of linear separating functions,

$$U_\alpha(\mathbf{x}) = \sum_{k=1}^K a_k^\alpha x_k + c^\alpha, \quad (25)$$

10 the monotonicity constraints are

$$a_k \geq 0, k \in K \subset \{1, \dots, n\}. \quad (26)$$

11 3.3 Positivity Constraints for Quadratic Separating Functions 12 (P-Constraints)

13 Unlike the case with the linear functions, imposing exact monotonicity con-
 14 straints on a quadratic function

$$U_j(\mathbf{x}) = \sum_{k=1}^n \sum_{l=1}^n a_{kl}^j x_k x_l + \sum_{k=1}^n b_k^j x_k + c^j \quad (27)$$

15 is a more complicated issue (indeed, in general, the quadratic function is not
 16 monotonic in the whole space R^n). Instead of imposing exact monotonicity
 17 constraints we consider the following constraints (we call them “positivity
 18 constraints”):

$$a_{kl}^\gamma \geq 0, b_k^\gamma \geq 0, k, l = 1, \dots, n. \quad (28)$$

19 Bugera et al. (2003) demonstrated that the positivity constraints can be
 20 easily included into the linear programming formulation of the classification
 21 problem. They do not significantly increase the computational time of the
 22 classification procedure, but provide robust results for small training datasets
 23 and datasets with missing or erroneous information. These constraints impose
 24 monotonicity with respect to variables $x_i, i = 1, \dots, n$ on the positive part
 25 $R^+ = \{ \mathbf{x} \in R^n \mid x_k \geq 0, k = 1, \dots, n \}$ of the object space R^n .

3.4 Gradient Monotonicity Constraints (GM-Constraints)

Another way to enforce monotonicity of quadratic separating functions is to restrict the gradient of separating functions on some set of objects X^* (for example, on a set of objects combining in-sample and out-of-sample points):

$$\left. \frac{\partial U_\alpha(\mathbf{x})}{\partial x^s} \right|_{\mathbf{x}=\mathbf{x}^*} \geq \gamma_s, \alpha = 0, \dots, J, s = 1, \dots, K, \mathbf{x}^* \in X^*, \quad (29)$$

where

$$\begin{aligned} \left. \frac{\partial U_\alpha(\mathbf{x})}{\partial x^s} \right|_{\mathbf{x}=\mathbf{x}^*} &= \left. \frac{\partial \left(\sum_{k=1}^n \sum_{l=1}^n a_{kl}^\alpha x_k x_l + \sum_{k=1}^n b_k^\alpha x_k + c^\alpha \right)}{\partial x^s} \right|_{\mathbf{x}=\mathbf{x}^*} = \\ &= \sum_{k=1}^n (a_{ks}^\alpha + a_{sk}^\alpha) x_k^* + b_s^\alpha. \end{aligned} \quad (30)$$

In constraint (29), $\gamma_s, s = 1, \dots, K$ are nonnegative constants.

3.5 Risk Constraints (R-Constraints)

Another important constraint that we apply to the model is the risk constraint. The risk constraint restricts the average value of the penalty function for misclassified objects. For this purpose we use the concept of Conditional Value-at-Risk (CVaR). The optimization approach for CVaR introduced by Rockafellar and Uryasev (2000), was further developed in Rockafellar and Uryasev (2002) and Rockafellar et al. (2002). Suppose X is a training set and for each object from this set the class number is known. In other words, discrete-value function $f(x)$ assigns the class for each object from set X . Let J be a total number of classes. The function $f(x)$ splits the set X into a set of subsets $\{X_1, \dots, X_J\}$,

$$X_j = \{ x \mid x \in X, f(x) = j \}. \quad (31)$$

Let I_j be the number of elements (objects) in the set X_j . We define the CVaR constraints as follows:

$$\begin{aligned} &\left\{ \zeta_j^+ + \frac{1}{\kappa} \frac{1}{I_j} \sum_{\mathbf{x} \in X_j} (U_j(\mathbf{x}) - \zeta_j^+)^+ \right\} \leq \\ &\leq \left\{ -\zeta_{j+1}^- - \frac{1}{\kappa} \frac{1}{I_{j+1}} \sum_{\mathbf{x} \in X_{j+1}} (-U_j(\mathbf{x}) - \zeta_{j+1}^-)^+ \right\} - \delta \end{aligned} \quad j = 1, \dots, I - 1, \quad (32)$$

where ζ_j^+ and ζ_j^- are free variables; δ and κ are parameters.

For an optimal solution of the optimization problem (16) with the constraints (32), the left-hand part of the inequality is an average of $U_j(\mathbf{x})$ for the largest $\kappa\%$ of objects from the j^{th} class; the right-hand part of the inequality is an average of $U_j(\mathbf{x})$ for the smallest $\kappa\%$ of objects from the $(j+1)^{\text{th}}$ class. We will call these values κ -CVaR largest of the j^{th} class and κ -CVaR smallest of the $(j+1)^{\text{th}}$ class, correspondingly. The general sense of the constraints is the following: the κ -CVaR largest of the j^{th} class is smaller at least by δ than the κ -CVaR smallest of the $(j+1)^{\text{th}}$ class.

1 3.6 Monotonicity of Separating Function with Respect to Class 2 Number (MSF-Constraints)

3 By introducing these constraints, we move apart values of the separating
4 functions for the objects belonging to different classes. Moreover, the con-
5 straints imply monotonicity of the separating functions with respect to the
6 index denoting the class number. The constraint is set for every pair of the
7 objects belonging to the neighboring classes. For every p and q , so that
8 $f_0(\mathbf{x}^p) + 1 = f_0(\mathbf{x}^q)$, the following constraint is imposed on the optimiza-
9 tion model:

$$U_i(\mathbf{x}^p) - U_i(\mathbf{x}^q) \geq \delta_{ipq}, \quad i = 0, \dots, I, \quad (33)$$

10 where δ_{ipq} are non-negative constants, and I is the number of separating
11 functions. Another way to impose monotonicity on the separating functions
12 with respect to the number of class is to consider δ_{ipq} as variables, and include
13 these variables into the objective function:

$$\sum_{i=1}^m \sum_{j=1}^J \lambda_j \sigma_i^j - \sum_{i,p,q} \delta_{ipq}. \quad (34)$$

14 3.7 Model Squeezing Constraints (MSQ-Constraints)

15 Squeezing constraints (implemented together with feasibility constraints) effi-
16 ciently adjust the flexibility of the model and control the out-of-sample versus
17 in-sample performance of the algorithm. The constraints have the following
18 form:

$$U_0(\mathbf{x}^p) - U_J(\mathbf{x}^p) \leq S_{f(\mathbf{x}^p)}, \quad p = 1, \dots, m. \quad (35)$$

19 With these constraints we bound variation of values of different separating
20 functions on each object. Another way to squeeze the spread of the separating
21 functions is to introduce a penalty coefficient for the difference of the functions
22 in (35). In this case the objective function of problem (16) can be rewritten
23 as:

$$\sum_{i=1}^m \sum_{j=1}^J \lambda_j \sigma_i^j + \sum_{p=1}^m (m_{f(\mathbf{x}^p)} (U_0(\mathbf{x}^p) - U_J(\mathbf{x}^p))). \quad (36)$$

24 The advantage of the squeezing constraints is that it is easy to implement in
25 the linear programming framework.

26 3.8 Level Squeezing Constraints (LSQ-Constraints)

27 This type of constraints is very similar to the model squeezing constraints.
28 The difference is that instead of squeezing the boundary functions $U_0(\mathbf{x})$ and
29 $U_J(\mathbf{x})$, we bound the absolute deviation of values of the separating functions
30 from their mean values on each class of objects. The constraints have the
31 following form:

$$\sum_{p: f_0(\mathbf{x}^p)=k} \left| U_j(\mathbf{x}^p) - \frac{1}{I_k} \sum_{q: f_0(\mathbf{x}^q)=k} U_j(\mathbf{x}^q) \right| \leq \gamma_{jk}, \quad (37)$$

1 where γ_{jk} are positive constants, and I_k is the number of objects \mathbf{x} in class
 2 K . Similar to constraints (35), constants γ_{jk} can be considered as variables
 3 and be included into the objective function:

$$\sum_{i=1}^m \sum_{j=1}^J \lambda_j \sigma_i^j + \sum_{j,k} \gamma_{jk}. \quad (38)$$

4 Choosing Model Flexibility

5 This section explains our approach to adjusting flexibility of classification
 6 models without strict definitions and mathematical details.

7 The considered constraints form various models based on the optimization
 8 of quadratic separating functions. The models differ in the type of constraints
 9 imposed on the separating functions.

10 The major characteristics of a particular classification model are in-sample
 11 and out-of-sample errors. To find a classification model we solve optimization
 12 problem (16) constructed for a training dataset. The error of misclassification
 13 of a classification model on the training dataset is called the “in-sample” er-
 14 ror. The error achieved by the constructed classification model on the objects
 15 outside of the training set is called the “out-of-sample” error. The misclassifi-
 16 cation error can be expressed in various ways. We measure the misclassification
 17 (if it is not specified otherwise) by the percentage of the objects on which the
 18 computed model gives wrong classification.

19 Theoretically, classification models demonstrate the following characteris-
 20 tics (see Figure 2). For small training sets, the model fits the data with zero
 21 in-sample error, whereas the expected out-of sample error is large. As the size
 22 of the training set increases, the expected in-sample error diverges from zero
 23 (the model cannot exactly fit the training data). Increasing the size of the
 24 training set leads to a larger expected in-sample error and a smaller expected
 25 out-of-sample error. For sufficiently large datasets, the in-sample and out-of-
 26 sample errors are quite close. In this case, we say that the model is saturated
 27 with data.

28 We say that the class of models A is more flexible than the class of models B
 29 if the class A includes the class B. Imposing constraints reduces the flexibility
 30 of the model. Figure 3 illustrates theoretical in/out-of-sample characteristics
 31 for two classes of models with different flexibilities. For small training sets,
 32 the less flexible model gives a smaller expected out-of-sample error (compared
 33 to the more flexible model) and, consequently, predicts better than the more
 34 flexible model. However, the more flexible models outperform the less flexible
 35 models for large training sets (more flexible models have a smaller expected

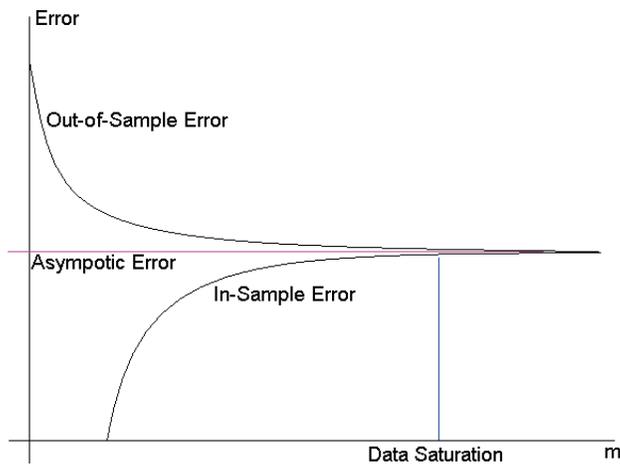


Fig. 2. In-sample and out-of-sample performance vs. size of training dataset.

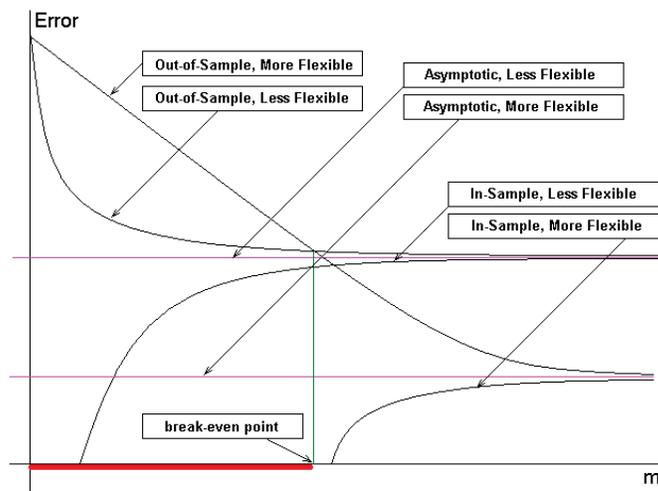


Fig. 3. In-sample and out-of-sample performance for different models.

- 1 out-of sample error compared to less flexible models). This is because the less
- 2 flexible model requires less data for saturation than the more flexible model.
- 3 A more flexible model may need a large training dataset to provide a good
- 4 prediction, while a less flexible model is saturated with data more quickly
- 5 than a more flexible model. Therefore, for small datasets, less flexible models

1 tend to outperform (out-of sample) more flexible models. However, for large
 2 datasets, more flexible models outperform (out-of-sample) less flexible ones.
 3 We demonstrate these considerations with the following example. For il-
 4 lustration purposes we consider four different models (A, B, C and D) with
 5 different levels of flexibility. We applied these models to the classification of
 6 a dataset containing 278 in-sample datapoints and 128 out-of-sample data-
 7 points. The in-sample dataset is used to construct linear programming prob-
 8 lem (20). The solution of this problem is a set of separating quadratic func-
 9 tions (19). The in-sample and out-of-sample classification is done using the obtained
 10 separating functions.

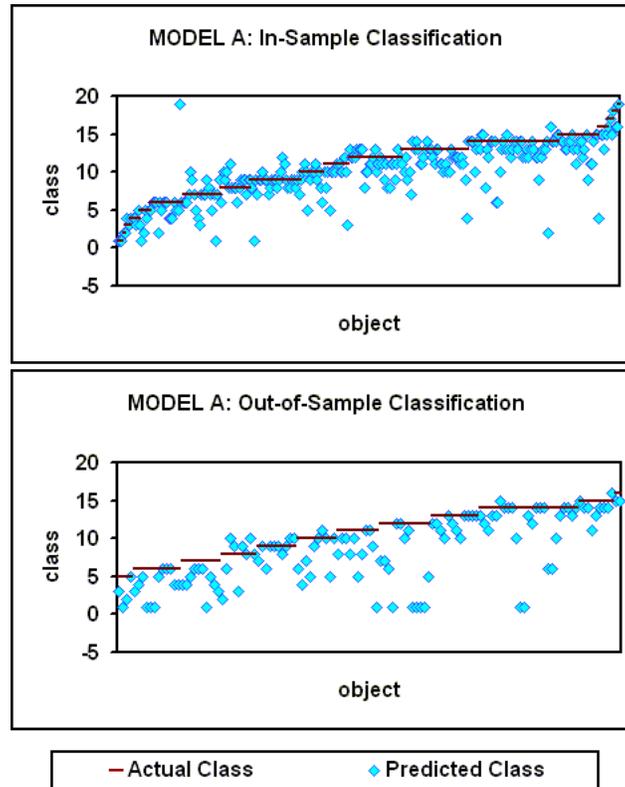


Fig. 4. In-sample and out-of-sample errors for model A. Model A uses parallel separating functions, and is the least flexible model among considered.

11 Figures 4, 5, 6, and 7 demonstrate the in-sample and out-of-sample per-
 12 formances of the considered models. Although the figures are based on actual
 13 computational results related to classification of bonds, they are presented
 14 in this paper for illustration purposes only. The horizontal axis of each graph

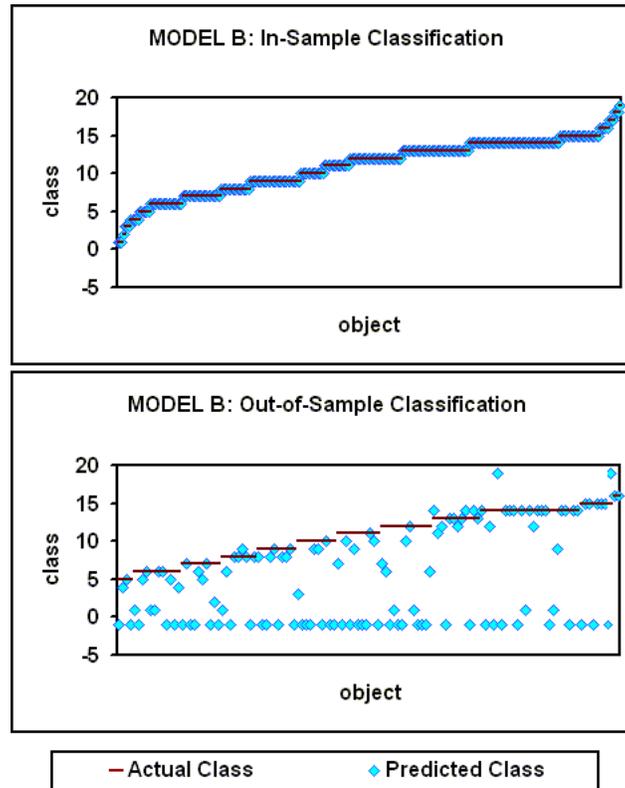


Fig. 5. In-sample and out-of-sample errors for model B. Model B uses separating functions without any constraints, and is the most flexible model. It has the best in-sample performance than Model A, but fails to assign classes for most of out-of-sample objects.

1 corresponds to the object number, which is ordered according to the object ac-
 2 tual class number. The vertical line corresponds to the calculated by the model
 3 class number. The solid line on the graph represents the actual class number
 4 of an object. The round-point line represents the calculated class. The left
 5 graphs show in-sample classification: the class is computed by the separating
 6 functions for the objects from the in-sample dataset. The right graphs show
 7 out-of-sample classification: the class is computed by the separating functions
 8 for the objects from the out-of-sample dataset.

9 Model A uses “parallel” quadratic separating functions. This case corre-
 10 sponds to the classification model considered by Bugera et al. (2003), where
 11 instead of multiple separating functions one utility function was used. The
 12 classification with the utility function can be interpreted as a classification
 13 with separating functions (15) with the following constraints on the functions

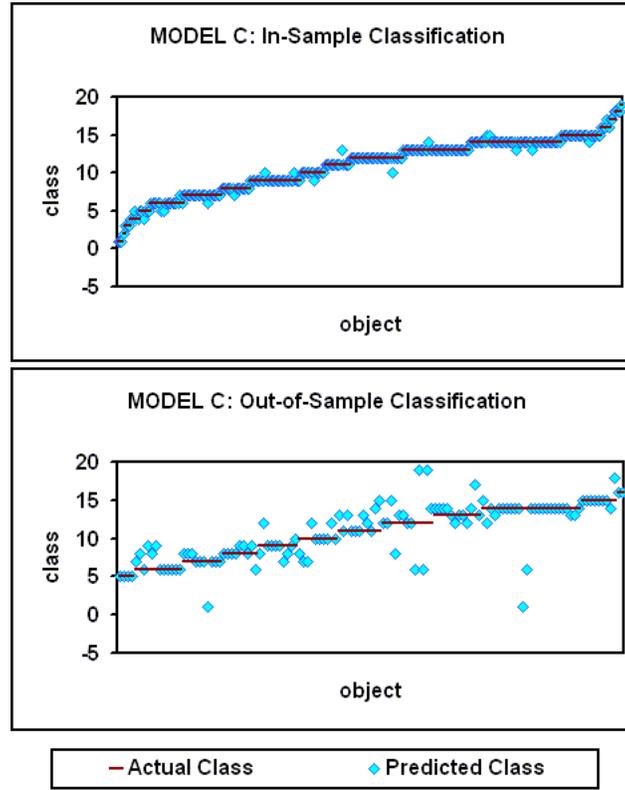


Fig. 6. In-sample and out-of-sample errors for model C. Model C is obtained from Model B by imposing feasibility constraints. This makes the out-of-sample prediction possible for all out-of-sample objects.

$$\begin{aligned}
 1 \quad U_\alpha(\mathbf{x}) &= \sum_{i=1}^n \sum_{j=1}^n a_{ij}^\alpha x_i x_j + \sum_{i=1}^n b_i^\alpha x_i + c^\alpha: \\
 &\forall \alpha_1, \alpha_2 \in \{0, \dots, J\} : \forall i, j \in \{0, \dots, n\} : a_{ij}^{\alpha_1} = a_{ij}^{\alpha_2}, b_i^{\alpha_1} = b_i^{\alpha_2} \quad (39)
 \end{aligned}$$

2 Constraints (39) make the separating functions parallel in the sense that
 3 the difference between the functions remains the same for all the points of the
 4 object space.

5 Figure 4 shows that Model A has large errors for both in-sample and
 6 out-of-sample calculations. According to Figure 8, when a model has approx-
 7 imately the same in-sample and out-of-sample errors, the data saturation occurs.
 8 Therefore, the out-of-sample performance cannot be improved by reduc-
 9 ing flexibility of the model. A more flexible model should be considered.

10 Model B uses quadratic separating functions without any constraints. This
 11 model is more flexible than Model A, because it has more control variables.

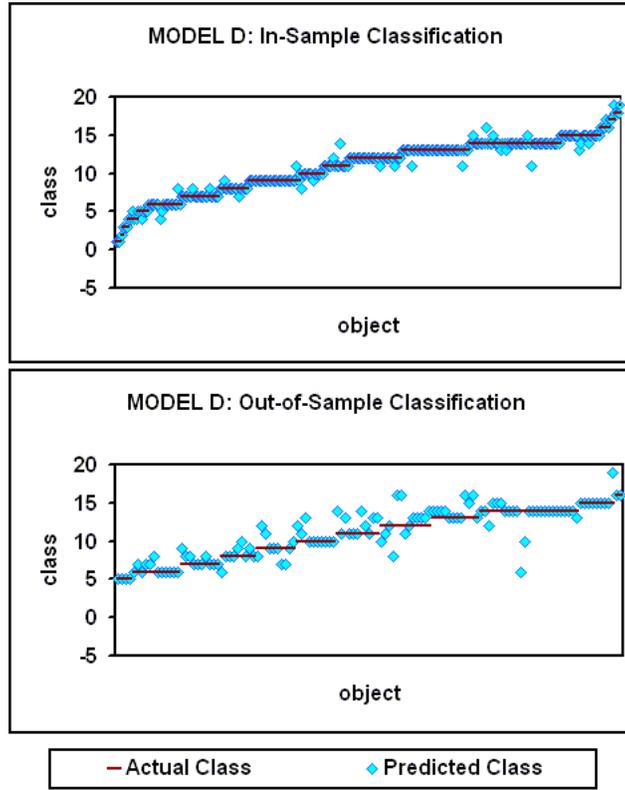


Fig. 7. In-sample and out-of-sample errors for model D. Model D is obtained from Model C by adding CVaR-risk constraints. Model D has the best out-of-sample performance among all the considered models (A, B, C, and D).

1 According to Figure 5, the in-sample error equals zero for Model B. However,
 2 in contrast to the in-sample behavior, the out-of-sample error is larger for
 3 Model B than for Model A. Moreover, Model B is not able to assign a class
 4 to many of the out-of-sample objects (on the picture the unpredictable result
 5 corresponds to the negative number -1 of the class). The separating functions
 6 may intersect, and some areas of the object are impossible to classify. For
 7 objects from these areas, the condition

$$U_0(\mathbf{x}^*), \dots, U_{j-1}(\mathbf{x}^*) > 0 \tag{40}$$

8 and

$$U_j(\mathbf{x}^*), \dots, U_J(\mathbf{x}^*) \leq 0 \tag{41}$$

9 are not satisfied for any j . It makes the result of the out-of-sample classification
 10 uninterpretable.

11 Model C is obtained by adding the feasibility constraints

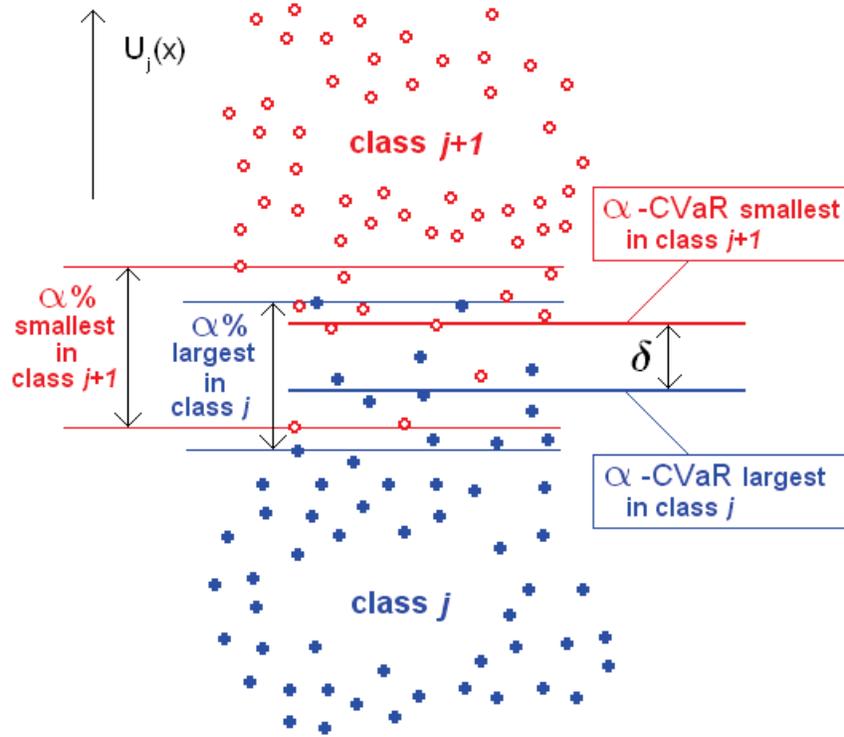


Fig. 8. Nature of risk constraint. For an optimal solution of the optimization problem with the CVaR constraints, the α -CVaR largest of the j^{th} class (the average of $U_j(\mathbf{x})$ for the largest $\alpha\%$ of objects from the j^{th} class) is smaller at least by δ than the α -CVaR smallest of the $(j + 1)^{th}$ class (the average of $U_{j+1}(\mathbf{x})$ for the smallest $\alpha\%$ of objects from the $(j + 1)^{th}$ class).

$$U_{j-1}(\mathbf{x}_i) \leq U_j(\mathbf{x}_i), j \in \{1, \dots, J\}, i = 1, \dots, m. \tag{42}$$

1 to Model B. Since Model C is less flexible than Model B, the in-sample error
 2 becomes greater than in Model B (see Figure 6). On the other hand, the model
 3 has a better out-of-sample performance than Models A and B. Moreover, a
 4 feasibility constraint makes the classification possible for any out-of-sample
 5 object. Since there is a significant discrepancy between in-sample and out-of-
 6 sample performances of the model, it is reasonable to impose more constraints
 7 on the model.

8 Model D is obtained by adding CVaR-risk constraints to Model C:

$$\left\{ \begin{array}{l} \left\{ \zeta_j^+ + \frac{1}{\alpha} \frac{1}{I_j} \sum_{\mathbf{x} \in X_j} (U_j(\mathbf{x}) - \zeta_j^+)^+ \right\} \leq \\ \left\{ -\zeta_{j+1}^- - \frac{1}{\alpha} \frac{1}{I_{j+1}} \sum_{\mathbf{x} \in X_{j+1}} (-U_j(\mathbf{x}) - \zeta_{j+1}^-)^+ \right\} - \delta \end{array} \right\} \quad j = 1, \dots, I - 1. \quad (43)$$

1 Figure 7 shows that this model has higher in-sample error compared to
 2 Model B, but the out-of-sample performance is the best among the considered
 3 models. Moreover, the risk constraint reduces the number of many-class mis-
 4 classification jumps. Whereas Model C has 8 out-of-sample objects for which
 5 the value of misprediction is more than 5 classes, Model D has only one object
 6 with a 5-class misprediction. So, Model D has the best out-of-sample perfor-
 7 mance among all the considered models. This has been achieved by choosing
 8 appropriate flexibility of the model.

9 A model with low flexibility may fit an in-sample dataset well. For models
 10 with a high in-sample error (such as Model A), more flexibility can be added
 11 by introducing more control variables to the model. In classification models
 12 with separating functions, feasibility constraints play a crucial role because
 13 they make classification always possible for out-of-sample objects. An exces-
 14 sive flexibility of the model may lead to “overfitting” and poor prediction
 15 characteristics of the model. To remove excessive flexibility of the considered
 16 models various types of constraints (such as risk constraints) can be imposed.
 17 The choice of the types of constraints, as well as the choice of the class of
 18 separating functions, plays a critical role for good out-of-sample performance
 19 of the algorithm.

20 5 Error Estimation

21 To estimate the performance of the considered models, we use the “leave-one-
 22 out” cross validation scheme. For the description of this scheme and other cross
 23 validation approaches, see, for instance, Efron and Tibshirani (1994). Let us
 24 denote by m a number of objects in the considered dataset. For each model
 25 (defined by the classes of the constraints imposed on optimization problem
 26 (20)), we performed m experiments. By excluding objects \mathbf{x}^i one by one from
 27 the set X , we constructed m training sets,

$$Y_i = X \setminus \{\mathbf{x}^i\}, \quad i = 1, \dots, m.$$

28 For each set Y_i , we solved (20) optimization problems with the appropriate
 29 set of constraints and found the optimal parameters of the separating func-
 30 tions $\{U_{\alpha^0}(\mathbf{x}), \dots, U_{\alpha^J}(\mathbf{x})\}$. Further, we computed the number of misclassified
 31 objects M_i from the set Y_i . Let us introduce the variable

$$P_i = \begin{cases} 0, & \text{if } \{U_{\alpha^0}(\mathbf{x}), \dots, U_{\alpha^J}(\mathbf{x})\} \text{ correctly classifies } \mathbf{x}^i, \\ 1, & \text{otherwise.} \end{cases} \quad (44)$$

1 The in-sample error estimate $E_{in-sample}$ is calculated by the following
 2 formula:

$$E_{in-sample} = \frac{1}{m} \sum_{i=1}^m \frac{M_i}{(m-1)} = \frac{1}{m(m-1)} \sum_{i=1}^m M_i, \quad (45)$$

3 where M_i is the number of misclassified objects in the set Y_i . In the last for-
 4 mula, the ratio $\frac{M_i}{m-1}$ estimates the probability of an in-sample misclassification
 5 in the set Y_i . We calculated the average of these probabilities to estimate the
 6 in-sample error.

7 The out-of-sample error estimate $E_{out-of-sample}$ is defined by the ratio of
 8 the total number of the misclassified out-of-sample objects in m experiments
 9 to the number of experiments:

$$E_{out-of-sample} = \frac{1}{m} \sum_{i=1}^m P_i \quad . \quad (46)$$

10 The considered leave-one-out cross-validation scheme provides the highest
 11 confidence level (by increasing the number of experiments) while effectively
 12 improving the predicting power of the model (by increasing the size of the
 13 in-sample datasets)

14 6 Bond Classification Problem

15 We have tested the approach with a bond classification problem. Bonds repre-
 16 sent the most liquid class of the fixed-income securities. A bond is an obligation
 17 of the bond issuer to pay cash flow to the bond holder according to the rules
 18 specified at the time the bond is issued. A bond pays a specific cash flow (face
 19 value) at the time of maturity. In addition, a bond may pay periodic coupon
 20 payments.

21 Although a bond generates a prespecified cash flow stream, it may default,
 22 if an issuer gets into financial difficulties or becomes a bankrupt. To charac-
 23 terize this risk, bonds are rated by several rating organizations (*Standard &*
 24 *Poor's* and *Moody's* are major rating companies). Bond rating evaluates the
 25 possibility of default of a bond issuer based on the issuer's financial condition
 26 and profits potential. The assignment of a rating class is mostly based on the
 27 issuer's financial status. A rating organization evaluates the status using ex-
 28 pert opinions and formal models based on various factors including financial
 29 ratios, such as the ratio of debt to equity, the ratio of current assets to current
 30 liabilities, and the ratio of cash flow to outstanding debt.

31 According to *Standard and Poor's*, bond ratings start at AAA for bonds
 32 having the highest investment quality, and end at D for bonds in payment
 33 default. The rating may be modified by a plus or minus to show relative
 34 standing within the category. Typically, a bond with a lower rating has a
 35 lower price than a bond generating the same cash flow, but having a higher
 36 rating.

1 In this study, we have replicated Standard & Poor's ratings of bonds using
 2 the suggested classification methodology. The rating replication problem is
 3 reduced to a classification problem with 20 classes.

4 7 Description of Data

5 For the computational verification of the proposed methodology we used sever-
 6 eral datasets (A, W, X, Y, Z) provided by the research group of the RiskSol-
 7 utions branch of Standard and Poor's, Inc. The datasets contain quantitative
 8 information about several hundred companies rated by Standard and Poor's,
 9 Inc. Each entry in the datasets has fourteen fields that correspond to certain
 10 parameters of a specific firm in a specific year. The first two fields are the com-
 11 pany name, and the year when the rating was calculated. We used these fields
 12 as identifiers of objects for classification. The next eleven fields contain quan-
 13 titative information about financial performance of the considered company.
 14 These fields are used for the decision-making in the classification process. The
 15 last field is the credit rating of the company assigned by Standard and Poor's,
 16 Inc. Table 1 represents the information used for classification.

Table 1. Data format of an entry of the dataset.

Identifier	Quantitative characteristics	Class
1) Company Name	1) Industry Sector	1) Credit rating
2) Year	2) EBIT interests coverage (times)	
	3) EBITDA interest coverage (times)	
	4) Return on capital (%)	
	5) Operating income/sales (%)	
	6) Free operating cash flow/total debt (%)	
	7) Funds form operations/total debt (%)	
	8) Total debt/capital (%)	
	9) Sales (mil)	
	10) Total equity (mil)	
	11) Total assets (mil)	

17 We preprocessed the data and rescaled all quantitative characteristics into
 18 $[-1,1]$ intervals: all quantitative characteristics are monotonic. Since the total
 19 number of rating classes equals 20, we used integer numbers from 1 to 20 to
 20 represent the credit rating of an object. The rating is arranged according to
 21 credit quality of the objects: the greater value corresponds to the better credit
 22 rating.

23 In the case study, we considered 5 different datasets (A, W, X, Y, Z) to
 24 verify the proposed methodology for different sizes of the input data. Each
 25 set was split into an in-sample set and an out-of-sample set. The first one was

1 used for developing a model, and the second one was used for the verification
 2 of the out-of-sample performance of the developed model. Table 2 contains
 3 information about the sizes of the considered datasets.

Table 2. Sizes of the considered datasets.

Dataset	A	W	X	Y	Z
Size of Set	406	205	373	187	108
Size of In-Sample	278	172	315	157	89
Size of Out-of-Sample	128	33	58	30	19

4 8 Numerical Experiments

5 For dataset A we performed computational experiments with 16 models gener-
 6 erated by all possible combinations of four different types of constraints. For
 7 each dataset W, X, Y, and Z we found a model with the best out-of-sample
 8 performance.

9 For dataset A we applied optimization problem (20) with Feasibility Con-
 10 straints (F). Besides, we applied all possible combinations of the following four
 11 types of constraints: Gradient Monotonicity Constraints (GM), Monotonicity
 12 of Separation Function w.r.t. Class Constraints (MSF), Risk Constraints (R),
 13 and Level Squeezing Constraints (LSQ). Each combination of the constraints
 14 corresponds to one of 16 different models, see Table 3. Although we considered
 15 other types of constraints (see Section 3.3), we report the results only for the
 16 models with F-, GM-, MSF-, R- and LSQ-constraints.

17 The numerical experiments were conducted with Pentium III , 1.2GHz
 18 in C/C++ environment. The linear programming problems were solved by
 19 CPLEX 7.0 package. The calculation results for 16 models for dataset A are
 20 presented in Table 4 and Figure 9.

21 Figure 9 represents in-sample and out-of-sample errors for different mod-
 22 els. For each model there are three columns representing the computational
 23 results. The first column corresponds to the percentage of in-sample classifica-
 24 tion error, the second column corresponds to the percentage of out-of-sample
 25 classification error, and the third column represents the average out-of-sample
 26 classification expressed in percents (100% misprediction corresponds to the
 27 case when the difference between actual and calculated classes equals 1).

28 On Figure 9, the models are arranged according to an average of out-of-
 29 sample error (see Table 4, column AVR). Models 0000 and 0010 have zero in-
 30 sample error. These two models are the most “flexible” – they fit the data with
 31 zero objective in (20). The smallest average out-of-sample error is obtained
 32 by 0101- model with MSF and LSQ constraints. Table 4 shows that this
 33 model has the lowest maximal misprediction. Models 0100 and 0110 have the

Table 3. List of the considered models.

Model Identifier	Types of Constraints Applied			
	GM	MSF	R	LSQ
0000				
0001				YES
0010			YES	
0011			YES	YES
0100		YES		
0101		YES		YES
0110		YES	YES	
0111		YES	YES	YES
1000	YES			
1001	YES			YES
1010	YES		YES	
1011	YES		YES	YES
1100	YES	YES		
1101	YES	YES		YES
1110	YES	YES	YES	
1111	YES	YES	YES	YES

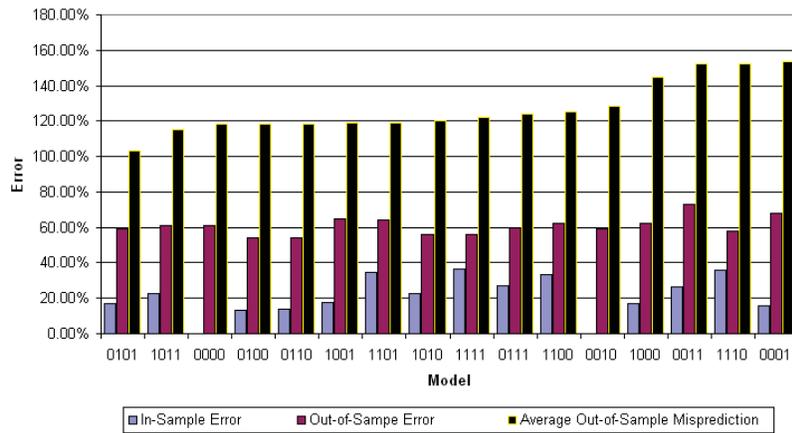


Fig. 9. In-sample and out-of-sample errors for various models.

1 minimal number of out-of-sample mispredictions (column Error(1)), but the
 2 maximal mispredictions (column MAX) for these models are high. It is worth
 3 mentioning that imposing constraints on the separating functions increases
 4 the computing time. However, for all the considered optimization problems,
 5 the classification procedure did not take more than 5 minutes: The CPLEX
 6 LP solver, that was used for this study, computed classification model 1111
 7 (all constraints were imposed) in about 5 minutes. Pentium III 1.2Ghz was

1 used for the study. As for all the other models, it took CPLEX LP solver less
 2 than 5 minutes.

Table 4. Calculation results for dataset A

Model	Error(In)	AVR	Correct	Error(1)	Error(2)	Error(3)	Error(4)	MAX
0000	0.00%	1.18	37.37%	62.63%	23.23%	9.09%	3.03%	18
0001	16.00%	1.54	30.30%	67.68%	35.35%	22.22%	13.13%	13
0010	0.00%	1.28	39.39%	60.61%	24.24%	10.10%	5.05%	17
0011	26.70%	1.53	25.25%	74.75%	37.37%	18.18%	12.12%	7
0100	13.50%	1.18	44.44%	55.56%	25.25%	12.12%	7.07%	18
0101	17.00%	1.03	39.39%	60.61%	25.25%	12.12%	4.04%	5
0110	13.60%	1.18	44.44%	55.56%	25.25%	12.12%	7.07%	18
0111	27.30%	1.24	38.38%	61.62%	30.30%	13.13%	5.05%	16
1000	16.90%	1.44	36.36%	63.64%	26.26%	11.11%	10.10%	16
1001	17.50%	1.19	33.33%	66.67%	29.29%	11.11%	7.07%	6
1010	22.60%	1.20	42.42%	57.58%	26.26%	12.12%	8.08%	14
1011	22.50%	1.15	37.37%	62.63%	29.29%	11.11%	7.07%	7
1100	33.10%	1.25	36.36%	63.64%	24.24%	14.14%	7.07%	18
1101	34.50%	1.19	34.34%	65.66%	24.24%	14.14%	7.07%	10
1110	35.70%	1.53	40.40%	59.60%	30.30%	19.19%	9.09%	18
1111	36.70%	1.22	42.42%	57.58%	28.28%	18.18%	8.08%	10

Error(In) = percentage of in-sample misclassifications; AVR = an average error (in number of classes); Correct = percentage of correct out-of-sample predictions; Error(1) = percentage of out-of-sample mispredictions; Error(i) = percentage of out-of-sample prediction errors for more than or equal to i classes; MAX = maximal error for out-of-sample classification.

3 For the datasets W, X, Y, and Z we imposed various constraints, and
 4 selected the models with best performance. For each dataset we choose the
 5 best performing models and compared performances of the models with refer-
 6 ence models used by *RiskSolutions group at Standard and Poor's*, see Table 5.
 7 For each data set the table contains the following performance information
 8 about the best chosen and the reference model: the average error (expressed
 9 in number of classes); the standard deviation; the percentage of correct out-
 10 of-sample predictions, the distribution of out-of-sample prediction (expressed
 11 in percentage of out-of-sample predictions with a deviation less than or equal
 12 to i classes for $i=1,2,3$); and the description of the constraints included in the
 13 best model.

14 The results of this analysis show that the proposed methodology of classi-
 15 fication with the separating functions is competitive with the reference model
 16 available from *RiskSolutions group at Standard and Poor's*. Although the pro-
 17 posed models do not provide better results compared to the reference model
 18 in all the cases, the proposed algorithms have a better or at least comparable
 19 performance, for small datasets Z and W.

Table 5. Comparison of best found model with reference model for sets W, X, Y and Z.

Model	W		X		Y		Z	
In-Sample	172		315		157		89	
Out-of-Sample	33		58		30		19	
Model	Best	Refer.	Best	Refer.	Best	Refer.	Best	Refer.
AVR	1.82	1.18	1.67	1.02	1.93	1.13	2.21	2.32
STDV	1.72	0.92	2.45	1.07	2.60	1.50	2.20	2.50
Correct	33.3%	21.2%	24.1%	37.9%	30.0%	50.0%	15.8%	21.1%
1-Class area	42.4%	72.7%	53.4%	70.7%	60.0%	70.0%	42.1%	63.2%
2-Class area	69.7%	87.9%	86.2%	94.8%	80.0%	83.3%	78.9%	68.4%
3-Class area	84.8%	100.0%	94.8%	96.6%	86.7%	86.7%	84.2%	68.4%
4-Class area	93.9%	100.0%	96.6%	98.3%	86.7%	96.7%	84.2%	78.9%
Constraints	Feasibility		Feasibility		Feasibility		Feasibility	
	MSF		MSF		MSF		CVaR	MSF

For each dataset the table has two columns. The first column corresponds to the best found model; the second column corresponds to the reference model from RiskSolutions. Each column contains information about out-of-sample performance of the corresponding model.

AVR = an average error (in number of classes); STDV = a standard deviation; Correct = percentage of correct out-of-sample predictions; i-Class area = percentage of out-of-sample predictions with a deviation less than or equal to i classes; Constraints = description of the constraints included in the best model

1 9 Concluding Remarks

2 We extended the methodology discussed in previous work (Bugera et al., 2003)
3 to the case of multi-class classification. The extended approach is based on
4 finding “optimal” separating functions belonging to a prespecified class of
5 functions. We have considered the models with quadratic separating func-
6 tions (in indicator parameters) with different types of constraints. As before,
7 selecting a proper class of constraints plays a crucial role in the success of the
8 suggested methodology: by imposing constraints, we adjust the flexibility of
9 a model to the size of the training dataset. The advantage of the considered
10 methodology is that the optimal classification model can be found by linear
11 programming techniques, which can be efficiently used for large datasets.

12 We have applied the developed methodology to a bond rating problem. We
13 have found that for the considered dataset, the best out-of-sample character-
14 istics (the smallest out-of sample error) are delivered by the model with MSF
15 and LSQ constraints on the control variables (coefficients of the separating
16 function).

17 The study was focused on the evaluation of computational characteristics
18 of the suggested approach. As before, despite the fact that the obtained results

1 are data specific, we conclude that the developed methodology is a robust
 2 classification technique suitable for a variety of engineering applications. The
 3 advantage of the proposed classification is its simplicity and consistency if
 4 compared to the proprietary models which are based on the combination of
 5 various techniques, such as expert models, neural networks, classification trees
 6 etc.

7 References

- 8 Bugera, V., Konno, H. and Uryasev, S.: 2003, Credit cards scoring with
 9 quadratic utility function, *Journal of Multi-Criteria Decision Analysis*
 10 **11**(4-5), 197–211.
- 11 Efron, B. and Tibshirani, R.J.: 1994, *An Introduction to the Bootstrap*, Chap-
 12 man & Hall, New York.
- 13 Konno, H., Gotoh, J. and Uho, T.: 2000, A cutting plane algorithm for semi-
 14 definite programming problems with applications to failure discrimination
 15 and cancer diagnosis, *Working Paper 00-5*, Tokyo Institute of Technology,
 16 Center for research in Advanced Financial Technology, Tokyo.
- 17 Konno, H. and Kobayashi, H.: 2000, Failure discrimination and rating of
 18 enterprises by semi-definite programming, *Asia-Pacific Financial Markets*
 19 **7**, 261–273.
- 20 Mangasarian, O., Street, W. and Wolberg, W.: 1995, Breast cancer diagnosis
 21 and prognosis via linear programming, *Operations Research* **43**, 570–577.
- 22 Pardalos, P.M., Michalopoulos, M. and Zopounidis, C.: 1997, On the use of
 23 multi-criteria methods for the evaluation of insurance companies in greece,
 24 in C. Zopounidis (ed.), *New Operational Approaches for Financial Modeling*,
 25 Physica, Berlin-Heidelberg, pp. 271–283.
- 26 Rockafellar, R.T. and Uryasev, S.: 2000, Optimization of conditional value-
 27 at-risk, *The Journal of Risk* **2**(3), 21–41.
- 28 Rockafellar, R.T. and Uryasev, S.: 2002, Conditional value-at-risk for general
 29 loss distributions, *Journal of Banking and Finance* **26**(7), 1443–1471.
- 30 Rockafellar, R.T., Uryasev, S. and Zabarankin, M.: 2002, Generalized devia-
 31 tions in risk analysis, *Finance and Stochastics* **10**, 51–74.
- 32 Vapnik, V.: 1998, *Statistical Learning Theory*, John Wiley & Sons, New York.
- 33 Zopounidis, C. and Doumpos, M.: 1997, Preference desegregation methodol-
 34 ogy in segmentation problems: The case of financial distress, in Zopouni-
 35 dis C. (ed.), *New Operational Approaches for Financial Modeling*, Physica,
 36 Berlin-Heidelberg, pp. 417–439.
- 37 Zopounidis, C., Pardalos, P., Doumpos, M. and Mavridou, T.: 1998, Multicri-
 38 teria decision aid in credit cards assessment, in C. Zopounidis and P. Parda-
 39 los (eds), *Managing in Uncertainty: Theory and Practice*, Kluwer Academic
 40 Publishers, New York, pp. 163–178.