

Multi-group Support Vector Machines with measurement costs: a biobjective approach^{*}

Emilio Carrizosa

Facultad de Matemáticas. Universidad de Sevilla (Spain).

Belen Martin-Barragan

Universidad Carlos III de Madrid (Spain)

Dolores Romero Morales

Saïd Business School. University of Oxford (United Kingdom)

Abstract

Support Vector Machine has shown to have good performance in many practical classification settings. In this paper we propose, for multi-group classification, a biobjective optimization model in which we consider not only the generalization ability (modelled through the margin maximization), but also costs associated with the features. This cost is not limited to an economical payment, but can also refer to risk, computational effort, space requirements, etc. We introduce a biobjective mixed integer problem, for which Pareto optimal solutions are obtained. Those Pareto optimal solutions correspond to different classification rules, among which the user would choose the one yielding the most appropriate compromise between the cost and the expected misclassification rate.

Key words: Multi-group Classification, Pareto Optimality, Biobjective Mixed Integer Programming, Feature Cost, Support Vector Machines.

^{*} This research was partially supported by projects MTM2004-22566-E, MTM2005-24550-E, TIN2004-21343-E, BFM2002-11282-E and MTM2005-09362-C03-01, of Ministerio de Ciencia y Tecnología (Spain) and FQM-329 of Plan Andaluz de Investigación (Andalucía, Spain).

Email addresses: `ecarrizosa@us.es` (Emilio Carrizosa),
`belen.martin@uc3m.es` (Belen Martin-Barragan),
`dolores.romero-morales@sbs.ox.ac.uk` (Dolores Romero Morales).

22 1 Introduction

23 In the last years operations researchers have made significant contributions to
24 problems related with Data Mining (e.g. [2,4,9,8,23,28]), such as Supervised
25 Classification. Roughly speaking, supervised classification consists of building
26 a rule to predict the class-membership of new objects from the same population
27 than those in a given database. Support Vector Machine (SVM), e.g. [11,13,20],
28 has shown to be a powerful tool for Supervised Classification. When only two
29 groups exist, this method attempts to build a hyperplane with maximal margin
30 that separates the two groups. Margin can be seen as a value that is zero when
31 there are misclassified objects and otherwise it measures the confidence in the
32 prediction, [1]. It has been shown (e.g. [11,32,33]) that this method enjoys good
33 generalization properties, in the sense that one can expect the good behavior
34 obtained in the available data to be generalized to the population which data
35 come from, since the probability of misclassifying a forthcoming individual
36 can be bounded by a function which is decreasing in the margin.

37 Generalization ability, addressed via margin maximization, will be our first
38 goal. However, in real-world classification problems it is very convenient to
39 obtain classification rules that, not only achieve good classification behavior,
40 but are also cheap or quick. A typical example is medical diagnosis, where
41 some tests are much more expensive or take much longer than others. If the
42 classification rule does not use variables based on the most expensive tests,

43 classifying new patients will be much cheaper or quicker, perhaps without
44 deteriorating significantly the quality of classification.

45 Together with misclassification costs, which are related with the generalization
46 ability of the rule, other costs, linked to the variables or attributes, can be
47 defined. In the simplest model we associate equal costs to each feature; keeping
48 the total cost below a given level amounts to stating an upper bound on the
49 number of features to be used. Turney [31] proposed other types of nontrivial
50 cost, for instance the test cost, also called measurement cost, where each test
51 (attribute, measurement, feature) has an associated cost, such as economical
52 payment, computational effort or some kind of complexity.

53 The aim of minimizing such costs has been mentioned before in the literature
54 as a desirable consequence of feature selection, see e.g. [18], but hardly directly
55 addressed.

56 In this paper, we address classification problems in which both misclassifica-
57 tion rate and measurement costs are relevant. To do this, we formulate a biob-
58 jective program of simultaneous minimization of misclassification rate, via the
59 maximization of the margin (the natural measure in SVM), and measurement
60 costs. Pareto-optimal solutions, i.e. classifiers that cannot be improved at the
61 same time in both objectives, are sought. The set of Pareto-optimal solutions
62 of the biobjective program gives us a finite set of classification rules, in such a
63 way that any rule which is not Pareto-optimal should be discarded, since it is

64 beaten in terms of margin and cost by another rule. Choosing one out of the
65 set of Pareto-optimal rules is done by choosing an appropriate compromise
66 between the two criteria involved.

67 We have structured the paper as follows. In Section 2 the problem is formally
68 introduced. In Section 3 we model the first goal: the measurement cost. Max-
69 imizing the margin, as a surrogate of minimizing the misclassification rate,
70 will be our second goal. Formal definitions of margin are given in Section 4,
71 by generalizing the concept of margin for two groups. A Biobjective Mixed
72 Integer Program formulation is given in Section 5, where a method to find the
73 Pareto-optimal classifiers, the Two-Phase Method [34], is proposed. In Sec-
74 tion 6, such biobjective formulations are modified to allow some points in the
75 training sample to be misclassified. Doing this we avoid the problem called
76 overfitting. Finally, some numerical results are presented in Section 7.

77 **2 The problem**

78 We have a finite set of classes $\mathcal{C} = \{1, 2, \dots, C\}$, and a set of objects Ω , each
79 object u having two components (x^u, c^u) . The first component x^u is called the
80 *predictor vector* and takes values in a set X . The set X is usually assumed to
81 be a subset of \mathbb{R}^p , and then, the components x_l , $l = 1, 2, \dots, p$, of the predictor
82 vector x are called *predictor variables*. The other component c^u , with values
83 in the set of classes \mathcal{C} , is called the *class-membership* of object u . Object u is
84 said to belong to class c^u .

85 In general, class-membership of objects in Ω is known only for a subset I ,
 86 called the *training sample*: both predictor vector and class-membership are
 87 known for $u \in I$, whereas only x^u is known for $u \in \Omega \setminus I$.

88 For any $c \in \mathcal{C}$, denote by I_c the set of objects in I belonging to class c :
 89 $I_c = \{u \in I : c^u = c\}$. We assume that each class is represented in the training
 90 sample, i.e., $I_c \neq \emptyset \forall c \in \mathcal{C}$.

91 We use a classification model in which a *score function*, $f = (f_c)_{c \in \mathcal{C}}$ with
 92 $f_c : X \rightarrow \mathbb{R}$, enables us to classify (allocate) any $z \in X$ as member of one
 93 of the classes as follows

$$94 \quad z \text{ is allocated to the class } c \text{ if } f_c(z) > f_j(z), \forall j \neq c, \quad (1)$$

95 i.e. z is allocated to the class c^* whose score function is highest:

$$96 \quad c^* = \arg \max_{c \in \mathcal{C}} f_c(z). \quad (2)$$

97 Notice that in case of ties, the object will be unclassified by this rule, and
 98 can be later allocated randomly or by a prefixed order to some class in
 99 $\arg \max_{c \in \mathcal{C}} f_c(z)$. Following a worst-case approach, we will consider those ob-
 100 jects as misclassified throughout the paper. Score functions f_c are assumed to
 101 have the form

$$102 \quad f_c(x) = \sum_{k=1}^N \alpha_k^c \phi_k(x) + \beta^c, \quad (3)$$

103 where $\alpha^c \in \mathbb{R}^N$, $\beta^c \in \mathbb{R}$, and $\mathcal{G} = \{\phi_1, \phi_2, \dots, \phi_N\}$ is a finite set of real-valued
 104 functions on X . Hence, each f_c belongs to a vector space \mathcal{F} , generated by \mathcal{G} .
 105 For instance, linear classifiers correspond to scores generated by

$$106 \quad \mathcal{G} = \{x_1, x_2, \dots, x_p\}, \quad (4)$$

107 whereas quadratic classifiers, [15,16], are obtained by setting

$$108 \quad \mathcal{G} = \{x_1, x_2, \dots, x_p\} \cup \{x_i x_j : 1 \leq i \leq j \leq p\} \quad (5)$$

109 i.e., the set of monomials of degree up to 2.

110 This framework also includes voting classifiers, such as boosting, e.g. [14,17],
 111 in which $\mathcal{C} = \{1, 2\}$ and a set of primitive classifiers $\phi_k : X \rightarrow \{0, 1\}$

$$112 \quad \phi_k(x) = 1 \text{ iff } x \text{ is allocated to class 1 via the } k\text{-th classifier,} \quad (6)$$

113 are combined linearly into a single score function of the form (3). For a very
 114 promising strategy for generating such primitive classifiers see e.g. [7].

115 Denote the coefficients of the score function by $A = (\alpha^1, \dots, \alpha^C)$ and $b =$
 116 $(\beta^1, \dots, \beta^C)$. The problem of choosing f is reduced to the choice of its coeffi-
 117 cients (A, b) .

118 **Definition 1** $f = (f_c)_{c \in \mathcal{C}}$ with $f_c : X \rightarrow \mathbb{R}$, is said to separate $\{I_c : c \in \mathcal{C}\}$
 119 if

$$120 \quad f_{c^u}(x^u) > f_j(x^u) \quad \forall j \neq c^u, \quad \forall u \in I. \quad (7)$$

121 Moreover, $\{I_c : c \in \mathcal{C}\}$ is said to be separable by the space \mathcal{F} if there exists
 122 $f = (f_c)_{c \in \mathcal{C}}$, with $f_c \in \mathcal{F}$, separating $\{I_c : c \in \mathcal{C}\}$.

123 Now we compare the definition of separability given in Definition 1 with those
 124 existing in the literature, [1,19,20,32].

125 For the two-group case, $\mathcal{C} = \{1, 2\}$, our definition is equivalent to the classical
 126 definition of separability stating that the convex hulls of $\{\phi(x^u) : u \in I_1\}$ and
 127 $\{\phi(x^u) : u \in I_2\}$ are contained in open halfspaces with a common hyperplane
 128 as boundary.

129 **Property 2** Let $\mathcal{C} = \{1, 2\}$. $\{I_c : c \in \{1, 2\}\}$ is separable iff there exists $(\omega, \gamma) \in$
 130 $(\mathbb{R}^N \setminus \{0\}) \times \mathbb{R}$ such that

$$\begin{aligned} & \omega^\top \phi(x^u) + \gamma > 0 \quad \forall u \in I_1 \\ & \omega^\top \phi(x^u) + \gamma < 0 \quad \forall u \in I_2. \end{aligned} \tag{8}$$

132 **Proof.** Take $\omega = \alpha^1 - \alpha^2$, $\gamma = \beta^1 - \beta^2$ and conversely, given (ω, γ) , satisfying
 133 (8), setting $\alpha^1 = \omega$, $\beta^1 = \gamma$, $\alpha^2 = 0$ and $\beta^2 = 0$, we have a score function that
 134 correctly classifies $\{I_c : c \in \{1, 2\}\}$. □

135 For the multi-group case, $|\mathcal{C}| > 2$, we have that, together with the concept of
 136 separability given in Definition 1, a natural alternative exists: we will say that
 137 $\{I_c : c \in \mathcal{C}\}$ is one-against-rest separable (OAR-separable) iff for all $c_1 \in \mathcal{C}$,
 138 $\{I_{c_1}, \bigcup_{c \in \mathcal{C} \setminus \{c_1\}} I_c\}$ is separable.

Property 3 One has

$$\text{OAR-separability} \Rightarrow \text{separability}$$

139 **Proof.** Let $\{I_c : c \in \mathcal{C}\}$ be OAR-separable. It means that, for each class
 140 c_1 , we have two score functions: f_{c_1} associated with I_{c_1} , and $f_{\bar{c}_1}$, associated
 141 with the objects in the remaining classes $\bigcup_{c \in \mathcal{C} \setminus \{c_1\}} I_c$. Since $(f_{c_1}, f_{\bar{c}_1})$ separates
 142 $\{I_{c_1}, \bigcup_{c \in \mathcal{C} \setminus \{c_1\}} I_c\}$, then

$$f_{c_1}(x^u) > f_{\bar{c}_1}(x^u) \quad \forall u \in I_{c_1} \tag{9}$$

$$f'_{c_1}(x^u) > f_{\bar{c}_1}(x^u) \quad \forall u \in \bigcup_{c \in \mathcal{C} \setminus \{c_1\}} I_c$$

143
 144 Set $g_c = f_c - f_{\bar{c}}$, for each $c \in \mathcal{C}$. Then $g_c(x^u) > 0$ iff $u \in I_c$. The function
 145 $g = (g_1, g_2, \dots, g_C)$ trivially separates $\{I_c : c \in \mathcal{C}\}$. Hence, OAR-separability
 146 implies separability. □

147 Notice that the converse implication does not hold: for instance, in Figure 1
 148 we have three classes 1, 2, 3 with elements denoted respectively by crosses
 149 (points $(4, -3)$, $(1, 0)$ and $(4, 3)$), stars (points $(-1, -1)$ and $(3, -4)$) and cir-
 150 cles (points $(-1, 1)$ and $(3, 4)$), which, as one can see in Figure 1, are not
 151 OAR-separable, but they are separable by the following score function,

$$\begin{aligned} f_1(x_1, x_2) &= x_1 \\ f_2(x_1, x_2) &= -x_2 \\ f_3(x_1, x_2) &= x_2. \end{aligned}$$

152 The definition of separability, as given in Definition 1, depends on the gener-
 153 ator \mathcal{G} . Under weak assumptions, there exists a generator, \mathcal{G} , rich enough to
 154 enable separability of $\{I_c : c \in \mathcal{C}\}$.

155 **Property 4** *If X is a subset of \mathbb{R}^p and $x^u \neq x^v, \forall u, v \in I$ with $c^u \neq c^v$, then*
 156 *there exists a finite generator \mathcal{G} such that $\{I_c : c \in \mathcal{C}\}$ is separable in the space*
 157 *\mathcal{F} generated by \mathcal{G} .*

Proof. For each $c \in \mathcal{C}$, consider the function

$$f_c(x) = - \prod_{u \in I_c} d(x, x^u)^2$$

where $d(\cdot, \cdot)$ stands for the Euclidean distance. This function is zero for all x^u with $u \in I_c$ and strictly negative otherwise. Then, for $u \in I_c$, and $c' \neq c$,

$$f_c(x^u) - f_{c'}(x^u) = -f_{c'}(x^u) > 0,$$

158 thus, such set of functions separates $\{I_c : c \in \mathcal{C}\}$.

159 Moreover, each f_c is a polynomial in the variables x_1, x_2, \dots, x_p , then it can
 160 be written as

$$f_c(x) = \sum_{\mathbf{h}=(h_1, \dots, h_p) \in \{0, 1, \dots, 2|I_c|\}^p} \alpha_{\mathbf{h}}^c \prod_{k=1}^p (x_k)^{h_k}, \quad (10)$$

162 belonging to the space \mathcal{F} generated by \mathcal{G} the set of monomials of degree up
 163 to $2|I|$. □

164 Suppose that \mathcal{F} is rich enough to enable separability, which ensures the exist-
 165 tence of separating functions f . However, uniqueness never holds. Indeed, it is
 166 easy to see that given $(\hat{\alpha}, \hat{\beta}) \in \mathbb{R}^{N+1}$ the classification rules obtained by (A, b)
 167 and (\tilde{A}, \tilde{b}) with $\tilde{\alpha}^c = \lambda \alpha^c + \hat{\alpha}$ and $\tilde{\beta}^c = \lambda \beta^c + \hat{\beta}$ for all $c \in \mathcal{C}$, are equivalent
 168 for all $\lambda > 0$, in the sense that both allocate objects to the same classes.

Moreover, there are also more than one score function that separates $\{I_c : c \in \mathcal{C}\}$ and they are not equivalent. For instance, given a score function separating $\{I_c : c \in \mathcal{C}\}$, let ε be any number satisfying:

$$0 < \varepsilon < \min_{u \in I} \min_{j \neq c^u} \{f_{c^u}(x^u) - f_j(x^u)\}.$$

169 The function $f^\varepsilon = (f_1 + \varepsilon, f_2, \dots, f_C)$ also separates $\{I_c : c \in \mathcal{C}\}$. We need a
 170 criterion for choosing one of them. Following Vapnik's publications in gener-
 171 alization ability, e.g. [32], we will use the margin maximization criterion, as
 172 will be explained in Section 4.

173 **3 Measurement costs**

174 Finding classifiers separating conveniently the groups is a plausible criterion
 175 when obtaining the predictor vector x^u is costless. When this is not the case,
 176 we should also take into account the cost associated with the evaluation of the
 177 classification rule.

178 In many practical applications, as medical diagnosis, the predictor variables
 179 of the data may be some diagnosis test (such as blood test, ...) that have
 180 associated a cost, either money, or risk/damage incurred to the patient. If the
 181 classifier built does not depend on some of these variables, we could avoid their
 182 measurement (and the corresponding cost) in the diagnosis of new patients.
 183 In this situation, we should seek a classifier that enjoys good generalization
 184 properties, and at the same time, has low cost.

185 Obtaining cheaper or quicker classification rules have been mentioned as one
186 of the desirable consequences of feature selection, where the aim is to reduce
187 the number of variables or features used by the classification rule. However
188 costs associated with such variables or features have seldom been considered.

189 Several authors have addressed measurement cost issues related with classi-
190 fication. For instance, [24,25,30] consider classification trees whose branching
191 rule takes such costs into account. See [31] for a comparison of such meth-
192 ods and [3,31] and the references therein for other proposals. In most cases,
193 the unique goal is to minimize some surrogate of the expected misclassification
194 cost, and, since the algorithm takes somehow into account measurement costs,
195 it is hoped that the measurement cost of individuals with the rule obtained
196 this way is not too high.

197 In this paper, however, we explicitly consider the minimization of measure-
198 ment costs as one criterion, whose trade-off with margin optimization is to be
199 determined by the user.

200 Costs are modelled as follows: Denote by Π_k the cost associated with evaluat-
201 ing the feature $\phi_k \in \mathcal{G}$ at a given x . For instance, if we are following a linear
202 approach, as given by (4), Π_l represents the cost of measuring the predictor
203 variable l in a new object.

Given the parameter $A = (\alpha^1, \dots, \alpha^C)$, define

$$S(A) = \{k \mid \exists c \in \mathcal{C} : \alpha_k^c \neq 0, 1 \leq k \leq N\}.$$

204 In other words, $S(A)$ represents the set of features we have to use in order to
 205 classify new objects. In principle, these are the features we have to pay for,
 206 so a score function with coefficients (A, b) will have associated a measurement
 207 cost equal to

$$208 \quad \pi(A, b) = \sum_{k \in S(A)} \Pi_k. \quad (11)$$

209 Pure linearity, as assumed in (11), may be unrealistic in some practical sit-
 210 uations. For instance, it may be the case that, once we have incurred a cost
 211 for obtaining some feature ϕ_k , some other features may be given for free or
 212 at reduced cost. This may happen, for example, in a medical context when
 213 the measurement of a variable requires a blood extraction, and some other
 214 variables can be measured using the same blood test. Another context where
 215 one encounters this, is the case in which some features are functions of other
 216 features: In model (5), feature $\phi(x) = x_i x_j$ is obtained for free if both features
 217 $\phi(x) = x_i$ and $\phi(x) = x_j$ have been previously inspected.

218 In Table 1 one can see the costs of a simple example with two classes $C = 2$,
 219 and $\mathcal{G} = \{\phi_1, \dots, \phi_5\}$ with different costs.

220 The score function given by $f_1 = \phi_1 + 4\phi_5$ and $f_2 = 3\phi_1 + 2$ incurs a cost of
 221 $2 + 2 = 4$.

222 Suppose that precedence constraints, in the form of a partial order \preceq between
 223 the features, is given. This means that if $h \preceq k$, the use of the feature ϕ_k

224 requires also the payment for feature ϕ_h . Moreover, in computing the total
 225 cost, the cost for every feature has to be summed at most once. In order to
 226 formalize this, define an auxiliary variable $z_k \in \{0, 1\}$ for each $k = 1, \dots, N$,
 227 representing

$$z_k = \begin{cases} 1 & \text{if payment of } \Pi_k \text{ is needed} \\ 0 & \text{otherwise} \end{cases} \quad (12)$$

228
 229 in other words:

$$z_k = \begin{cases} 1 & \text{if } h \in S(A) \text{ for some } h \text{ with } k \preceq h \\ 0 & \text{otherwise} \end{cases} \quad (13)$$

230
 231 Thus, cost associated with a score function with coefficients (A, b) will be

$$\pi(A, b) = \sum_{k=1}^N z_k \Pi_k. \quad (14)$$

232
 233 Particular cases already suggested in the literature can be easily accommo-
 234 dated into our framework. For instance, in [26] variables are grouped in a way
 235 that, if one variable from a group is requested, then all the others in the same
 236 group are available for zero additional cost. To model this case in our setting,
 237 define the cost of one variable from each group to be equal to the cost of the
 238 group it belongs to, and set the remaining variables to have zero cost. More-
 239 over, choose a partial order \prec for which $h \prec j$ iff variables h and j are in the

240 same group and h has nonzero cost.

241 Moreover, this modelling technique allows us to use, but it is not limited to,
242 polynomial kernels. Indeed, suppose a kernel $k(x, y) = \Phi(x)^\top \Phi(y)$ for some
243 $\Phi : X \rightarrow F$. If Φ holds

- 244 • F is a finite dimensional *feature space*, $F \subset \mathbb{R}^N$,
- 245 • for any component ϕ_k , $k = 1, 2, \dots, N$ of $\Phi = (\phi_1, \phi_2, \dots, \phi_N)$, the informa-
246 tion about what original variables are needed to calculate ϕ_k is available,

247 then, the cost associated to a score function can be modelled using the method-
248 ology explained in this section.

249 We will show in Sections 5 and 6, that this modelling technique allows for-
250 mulations as Biobjective Mixed Integer Programs. For these models there ex-
251 ist suitable techniques for finding their Pareto-optimal solutions. Biobjective
252 problems for more general problems, such as e.g. measurement cost minimiza-
253 tion using kernels which are not of polynomial type, [36], can also be formu-
254 lated. However, they yield combinatorial problems which are much harder to
255 solve in practice.

256 Minimizing (14) will be one of our goals. However, our main goal is finding
257 classifiers with good generalization properties. This, the second objective in
258 our model, will be discussed in detail in the following section.

259 **4 Margin optimization**

260 Throughout this section, unless explicitly stated, we assume that \mathcal{F} is rich
 261 enough to enable separability:

262 **Assumption 1** $\{I_c : c \in \mathcal{C}\}$ is separable by \mathcal{F} .

263 We may observe that we can always consider F as in Property 4, and therefore
 264 Assumption 1 will be hold. However we expect in practice to attain separability
 265 with smaller generators.

266 Since by Assumption 1 objects in I will be correctly classified, the substan-
 267 tial matter is the classification of objects $u \in \Omega \setminus I$. Hence, we are interested
 268 in obtaining classifiers with good generalization properties, via margin max-
 269 imization,[11,32,33]). The concepts of functional and geometrical margin, in-
 270 troduced in Cristianini and Shawe-Taylor [13] for the case of two groups, are
 271 extended below to the multi-group case.

272 **Definition 5** The functional margin of an object u with respect to the score
 273 function f , with coefficients (A, b) , is the quantity

$$274 \quad \hat{\theta}^u(A, b) = \min_{j \neq c^u} \{f_{c^u}(x^u) - f_j(x^u)\} \quad (15)$$

275 The functional margin of a score function f , with coefficients (A, b) with respect
 276 to a training sample I is equal to

$$277 \quad \hat{\theta}^I(A, b) = \min_{u \in I} \hat{\theta}^u. \quad (16)$$

278 We immediately obtain

279 **Property 6** A score function f with coefficients (A, b) separates $\{I_c : c \in \mathcal{C}\}$
 280 if and only if, the margin $\hat{\theta}^I(A, b)$ is strictly positive.

281 The choices (A, b) and $(\lambda A, \lambda b)$ yield the same classification rule, but have
 282 different functional margins. Hence, as in the two-group case, we need to
 283 normalize this quantity in order to be able to compare score functions.

284 The normalization done here is made dependent on a norm $\|\cdot\|$, which can be
 285 different from the standard choice of the Euclidean norm, [13]. This will allow
 286 us, as shown in Section 5, to formulate the resulting optimization problems as
 287 mixed integer linear problems, solvable with existing commercial software.

288 **Definition 7** Let $\|\cdot\|$ be a norm in $\mathbb{R}^{C \times N}$. The geometrical margin of an
 289 object u with respect to the score function (A, b) , with $A \neq 0$, is the quantity

$$290 \quad \theta^u(A, b) = \frac{\hat{\theta}^u}{\|A\|}. \quad (17)$$

291 The geometrical margin of a score function (A, b) with respect to a training
 292 sample I is the minimum:

$$293 \quad \theta^I(A, b) = \min_{u \in I} \theta^u. \quad (18)$$

294 Now, we consider the problem of maximizing the geometrical margin

$$295 \quad \max_{A \neq 0, b \in \mathbb{R}^C} \frac{\min_{u \in I} \hat{\theta}^u(A, b)}{\|A\|}. \quad (19)$$

296 We have an alternative formulation, in terms of the functional margin, as given
 297 by the following proposition.

298 **Proposition 8** Problem (19) is equivalent to:

$$299 \quad \begin{aligned} & \max \min_{u \in I} \hat{\theta}^u(A, b) \\ & s.t: \quad \|A\| \leq 1, \end{aligned} \quad (20)$$

300 in the sense that any optimal solution of (20) is also optimal for (19), and for
 301 any optimal solution (A^*, b^*) of (19),

$$302 \quad (\hat{A}, \hat{b}) = \frac{1}{\|A^*\|} (A^*, b^*) \quad (21)$$

303 is an optimal solution of (20).

304 **Property 9** Problem (20) has finite optimal value.

305 **Proof.** Let $(A, b) = (\alpha^1, \dots, \alpha^C; \beta^1, \dots, \beta^C)$ be a feasible solution of (20).

306 Let $u \in I$ and $j \neq c^u$, then

$$\begin{aligned} & |\alpha^{c^u} \phi(x^u) + \beta^{c^u} - \alpha^j \phi(x^u) - \beta^j| \\ &= |(\alpha^{c^u} - \alpha^j) \phi(x^u) + \beta^{c^u} - \beta^j| \\ &\leq |(\alpha^{c^u} - \alpha^j) \phi(x^u)| + |\beta^{c^u} - \beta^j| \end{aligned} \quad (22)$$

307 To bound the first term, observe that, since all norms are equivalent, there

308 exists K such that $|\alpha_k^c| \leq K$ for all $k = 1, 2, \dots, N$, $c \in \mathcal{C}$.

309 Hence,

$$\begin{aligned} & |(\alpha^{c^u} - \alpha^j) \phi(x^u)| \\ &\leq \sum_{k=1}^N |\alpha_k^{c^u} - \alpha_k^j| |\phi_k(x^u)| \\ &\leq 2KN \max_{1 \leq k \leq N, u \in I} |\phi_k(x^u)| = K' < \infty \end{aligned}$$

310 Now, we will bound the term $|\beta^{c^u} - \beta^j|$. Since each class is represented, $I_j \neq \emptyset$,

311 let $v \in I_j$. Solution (A, b) feasible for (20) implies both u and v are correctly

312 classified,

$$\alpha^{c^u} \phi(x^u) + \beta^{c^u} - (\alpha^j \phi(x^u) + \beta^j) > 0$$

$$\alpha^{c^u} \phi(x^v) + \beta^{c^u} - (\alpha^j \phi(x^v) + \beta^j) < 0$$

313 yielding,

$$314 \quad (\alpha^{c^u} - \alpha^j) \phi(x^v) < \beta^j - \beta^{c^u} < (\alpha^{c^u} - \alpha^j) \phi(x^u). \quad (23)$$

315 Thus

$$\begin{aligned} & |\beta^j - \beta^{c^u}| \\ & \leq \max\{ |(\alpha^{c^u} - \alpha^j) \phi(x^u)|, |(\alpha^{c^u} - \alpha^j) \phi(x^v)| \} \\ & \leq \max_{v \in I} \{ |(\alpha^{c^u} - \alpha^j) \phi(x^v)| \} \\ & \leq 2KN \max_{1 \leq k \leq N, v \in I} |\phi_k(x^v)| = K'. \end{aligned}$$

316 Hence the objective function is bounded by

$$317 \quad \min_{u \in I} \theta^u = \min_{u \in I} \min_{j \neq c^u} |\alpha^{c^u} \phi(x^u) + \beta^{c^u} - \alpha^j \phi(x^u) - \beta^j| \leq 2K'. \quad (24)$$

318

□

319 We have assumed that \mathcal{F} is rich enough to enable separability of $\{I_c : c \in \mathcal{C}\}$.

320 However, it may be useful to have a method to check such separability. In case

321 we do not know if $\{I_c : c \in \mathcal{C}\}$ is separable in a space \mathcal{F} , solving Problem (20)

322 allow us to check it. Indeed we have the property:

323 **Property 10** $\{I_c : c \in \mathcal{C}\}$ is separable if and only if Problem (20) has strictly
324 positive optimal value.

325 Another reduction of Problem (20) is even possible. For all $\lambda \in \mathbb{R}$ the score

326 functions defined by (A, b) and (A, \tilde{b}) , with $\tilde{b}^c = b^c + \lambda$ for all $c \in \mathcal{C}$, are

327 equivalent in the sense that both classify objects to the same classes, and

328 both have the same margins. Then, we can restrict the coefficients β^c to be
 329 nonnegative, yielding the problem:

$$\begin{aligned} \max \quad & \min_{u \in I} \theta^u(A, b) \\ \text{s.t. :} \quad & \|A\| \leq 1 \end{aligned} \tag{25}$$

$$(A, b) \in \mathbb{R}^{NC} \times \mathbb{R}_+^C.$$

330

331 **Property 11** *Problems (20) and (25) are equivalent in the sense that every*
 332 *optimal solution of (25) is also optimal for (20), and, for any optimal solution*
 333 *of (20), there exists a feasible solution of (25) that is also optimal in both*
 334 *problems.*

335 5 A biobjective approach

336 In the last sections we have described the two objectives of our problem,
 337 namely, maximizing the margin and minimizing the measurement cost. Hence
 338 we have the following biobjective problem:

$$\begin{aligned} \max \quad & \theta(A, b) \\ \min \quad & \pi(A, b) \\ \text{s.t. :} \quad & \|A\| \leq 1 \end{aligned} \tag{26}$$

$$(A, b) \in \mathbb{R}^{NC} \times \mathbb{R}_+^C.$$

339

340 **Property 12** *The set of Pareto-optimal outcomes of the biobjective problem*
 341 *(26) is finite.*

342 **Proof.** The set of all outcomes of (26) can be calculated by solving the problem

$$\begin{aligned}
 \max \quad & \theta(A, b) \\
 \text{s.t. :} \quad & \|A\| \leq 1
 \end{aligned} \tag{27}$$

$$\pi(A, b) \leq \pi$$

$$(A, b) \in \mathbb{R}^{NC} \times \mathbb{R}_+^C$$

343

for any π in the set of possible costs:

$$\{\pi(A, b) : (A, b) \in \mathbb{R}^{NC} \times \mathbb{R}_+^C\},$$

344 which is contained in the finite set $\{\sum_{k \in S} \Pi_k : S \subseteq \{1, 2, \dots, N\}\}$. \square

345 Using the notation of section above, (26) can also be reformulated as

$$\max y$$

$$\min \sum_{k=1}^N \Pi_k z_k$$

$$s.t. : \sum_{k=1}^N \phi_k(x^u) (\alpha_k^i - \alpha_k^j) + \beta^i - \beta^j - y \geq 0, \quad \forall i \neq j; i, j \in \mathcal{C}, u \in I_i$$

$$\|A\| \leq 1$$

$$-z_k \leq \sum_{k:h \leq k} \sum_{c=1}^{\mathcal{C}} \alpha_k^c \leq z_h \quad \forall h = 1, 2, \dots, N \quad (28)$$

$$\alpha_k^c \text{ unrestricted} \quad \forall k = 1, 2, \dots, N; c \in \mathcal{C}$$

$$y \text{ unrestricted}$$

$$\beta^c \geq 0 \quad \forall c \in \mathcal{C}$$

$$z_k \in \{0, 1\} \quad \forall k = 1, 2, \dots, N$$

346

347 In this formulation if $\|\cdot\|$ is the L_∞ , then the normalization constraint is

348 redundant.

349 Due to the presence of a nonlinear constraint ($\|A\| \leq 1$), Problem (28) is a

350 biobjective mixed integer nonlinear program.

351 Many classical SVM implementations have used the Euclidean norm [13],

352 yielding a quadratic program. Mangasarian [22] proposes the use of other

353 norms. In particular Linear Programming approaches have been implemented

369 L_∞ norm, instead of the L_1 norm, in the definition of geometrical margin.

370 Problem (29) is a biobjective mixed integer linear problem, which can be
371 tackled for instance, by adapting the two-phase method of [34] designed for
372 solving biobjective knapsack problems.

373 In the first phase, one obtains the so-called supported solutions, namely, those
374 which are found as solution of the scalarized problem

$$\begin{aligned} \max \quad & \lambda_1 \theta(A, b) - \lambda_2 \pi(A, b) \\ \text{s.t. :} \quad & \|A\| \leq 1 \end{aligned} \tag{30}$$

$$(A, b) \in \mathbb{R}^{NC} \times \mathbb{R}_+^C$$

375

376 for some weights $\lambda_1, \lambda_2 \in [0, 1]$, with $\lambda_1 + \lambda_2 = 1$. These points describe, in the
377 outcome space, the frontier of the convex hull of the Pareto-optimal outcomes.

Since we face a bi-objective problem, the set of possible weights

$$\Lambda = \{(\lambda_1, \lambda_2) \in \mathbb{R}_+^2 : \lambda_1 + \lambda_2 = 1\}$$

378 that describe the supported efficient outcomes is unidimensional, and only a
379 finite number of weights describe different outcomes. This fact can be exploited
380 to find all supported outcomes in a sequential way.

381 A solution with minimal (zero) cost is the trivial solution $(A, b) = (0, 0)$. Note
382 that with this solution, points are classified arbitrarily by the tie-break rules,

383 since all the score functions will be zero.

When we are optimizing only the first objective, namely maximizing the margin, the optimal value can be obtained by solving Problem (20), which can be easily reformulated as a linear program. Denote by θ^* its optimal value. Given an optimal solution (A^*, b^*) of (20), a feasible solution (A^*, b^*, z^*) of the biobjective problem (26) can be built by setting

$$z_i^* = \begin{cases} 1, & \text{if } \alpha_i^{*c} \neq 0 \text{ for some } c \in C, \\ 0, & \text{otherwise.} \end{cases}$$

If (A^*, b^*) is the unique optimal solution, then (A^*, b^*, z^*) will be a Pareto-optimal point. Otherwise, a Pareto-optimal point of (26) can be found by maximizing the margin, i.e., by solving,

$$\begin{aligned} \min \quad & \pi(A, b) \\ \text{s.t. :} \quad & \|A\| \leq 1 \\ & \theta(A, b) \geq \theta^* \\ & (A, b) \in \mathbb{R}^{NC} \times \mathbb{R}_+^C \end{aligned}$$

Once we have both a Pareto-optimal solution with minimal cost, i.e. $(0, 0)$, and a Pareto-optimal solution with maximal margin, namely (A_0, b_0) , we construct an ordered list (sorted by either margin or by cost) whose elements can be

built from any two consecutive already known elements (A_1, b_1) and (A_2, b_2) by the scalarized Problem (30) for certain λ_1 and λ_2 . Denote θ_1 and θ_2 the margin of solution (A_1, b_1) and (A_2, b_2) respectively and costs Π^1 and Π^2 . The scalarization needed in the problem is

$$\lambda_1 = \frac{\theta^2 - \theta^1}{\theta^2 - \theta^1 + \Pi^2 - \Pi^1}$$

$$\lambda_2 = \frac{\Pi^2 - \Pi^1}{\theta^2 - \theta^1 + \Pi^2 - \Pi^1}.$$

384 All optimal solutions of such scalarized problem are Pareto-optimal points. If
 385 both (or any of) (A_1, b_1) and (A_2, b_2) are solutions of the scalarized problem,
 386 the set of its optimal solutions yield the only supported Pareto outcomes be-
 387 tween those of (A_1, b_1) and (A_2, b_2) , so we do not need to seek more supported
 388 Pareto points between them. Since the number of Pareto outcomes is finite,
 389 the process ends in finite time.

When all the supported Pareto outcomes are found, the non-supported ones may be obtained in the following way. Let (A_1, b_1) be any Pareto-optimal point with cost $\Pi^1 > 0$. Let $\hat{\pi}$ be the minimal feature cost that is positive,

$$\hat{\pi} = \min_{k=1,2,\dots,N} \{\Pi_k : \Pi_k > 0\}.$$

390 Then a Pareto-optimal point, with cost strictly lower than Π^1 , is obtained by
 391 solving the problem

$$\begin{aligned} \max \quad & \theta(A, b) \\ \text{s.t. :} \quad & \|A\| \leq 1 \end{aligned} \tag{31}$$

$$\pi(A, b) \leq \Pi^1 - \hat{\pi}$$

$$(A, b) \in \mathbb{R}^{NC} \times \mathbb{R}_+^C.$$

392

393 Then, the next Pareto-optimal point can be found in the same way. Thus,
 394 starting from any supported Pareto-optimal point with cost greater than zero,
 395 the non-supported Pareto-optimal outcomes between it and the next sup-
 396 ported one can be found.

397 **6 Soft-margin biobjective optimization**

398 In classification problems, when the number of parameters to be fitted is large,
 399 model may incur a phenomenon called overfitting. It is said to happen when
 400 a classification rule achieves very good performance in the training sample I ,
 401 but does not generalize well, thus yielding a bad performance in future objects.

402 Moreover, it may happen that I is not separable in the feature space. Then,
 403 the models proposed in the previous section do not apply, since they look for
 404 rules which correctly classify all the objects in I . As stated in Property 4,

405 other feature space could be used, but usually they are more complicated and
406 thus the model would incur overfitting.

407 In order to both avoid overfitting and deal with the non-separability of I ,
408 the typical SVM approach, called soft-margin maximization [13], is based on
409 allowing some objects in I to be misclassified. This is done by adding to the
410 model some slack variables $\xi \in \mathbb{R}_+^n$, where n is the cardinal of the train-
411 ing sample. Using this idea, the biobjective Problem (28) is replaced by the

412 following problem:

$$\begin{aligned}
& \max y \\
& \min \sum_{k=1}^N \Pi_k z_k \\
& s.t. : \sum_{k=1}^N \phi_k(x^u) (\alpha_k^i - \alpha_k^j) + \beta^i - \beta^j - y + \xi^u \geq 0, \quad \forall i \neq j; i, j \in \mathcal{C}, u \in I_i \\
& \|A\| + \gamma \sum_{u \in I} \xi^u \leq N \\
& -N z_k \leq \sum_{k:h \preceq k} \sum_{c=1}^C \alpha_k^c \leq N z_h \quad \forall h = 1, 2, \dots, N \\
& \alpha_k^c \text{ unrestricted} \quad \forall k = 1, 2, \dots, N; c \in \mathcal{C} \\
& y \text{ unrestricted} \\
& \beta^c \geq 0 \quad \forall c \in \mathcal{C} \\
& z_k \in \{0, 1\} \quad \forall k = 1, 2, \dots, N \\
& \xi^u \geq 0 \quad \forall u \in I,
\end{aligned} \tag{32}$$

413

414 for some user-defined value γ , which trades off the perturbations ξ^u and the
415 margin.

416 In the same way as for the hard-margin approach, when $\|\cdot\|$ is a polyhedral
417 norm, this problem can be formulated as a Biobjective Mixed Integer Problem.

418 For instance, if $\|\cdot\|$ is a scaled L_1 -norm, then Problem (32) can be formulated

419 as follows:

$$\max y$$

$$\min \sum_{k=1}^N \Pi_k z_k$$

$$s.t. : \sum_{k=1}^N \phi_k(x^u) \left(\alpha_{+k}^i - \alpha_{-k}^i - \alpha_{+k}^j + \alpha_{-k}^j \right) + \beta^i - \beta^j - y + \xi^u \geq 0,$$

$$\forall i \neq j; i, j \in \mathcal{C}, u \in I_i$$

$$\sum_{c=1}^{\mathcal{C}} \sum_{k=1}^N \left(\alpha_{+k}^c + \alpha_{-k}^c \right) + \gamma \sum_{u \in I} \xi^u \leq N$$

$$\sum_{k:h \preceq k} \sum_{c=1}^{\mathcal{C}} \left(\alpha_{+k}^c + \alpha_{-k}^c \right) \leq N z_h \quad \forall h = 1, 2, \dots, N \quad (33)$$

y unrestricted

$$\alpha_{+k}^c \geq 0 \quad \forall k = 1, 2, \dots, N; c \in \mathcal{C}$$

$$\alpha_{-k}^c \geq 0 \quad \forall k = 1, 2, \dots, N; c \in \mathcal{C}$$

$$\beta^c \geq 0 \quad \forall c \in \mathcal{C}$$

$$z_k \in \{0, 1\} \quad \forall k = 1, 2, \dots, N$$

$$\xi^u \geq 0 \quad \forall u \in I$$

420

421 The Two-Phase Method proposed in Section 5 to find the Pareto-optimal

422 classifiers can also be used for solving (33). Note that in this case, the solution

423 with minimal (zero) cost is not the trivial solution $(A, b) = (0, 0)$, but any

424 optimal solution (33) with A set equal to the null matrix. The following steps
425 of the method remain analogous to the hard-margin approach, and will not
426 be repeated here.

427 **7 Numerical results**

428 In order to explore both, costs and quality, of the Pareto score functions
429 obtained, we have performed a series of numerical tests on four standard
430 databases, publicly available from the UCI Machine Learning Repository [6],
431 namely, the BUPA Liver-disorders Database, called here `bupa`; the Pima In-
432 dians Diabetes Database, called here `pima`; the New Diagnostic Database,
433 contained in the Wisconsin Breast Cancer Databases, called here `wdbc`, and
434 the Credit Screening Databases, called here `credit`.

435 For each database, the name of the file (as called in the database), the total
436 number of objects $|\Omega|$, the number of groups C and the number of variables
437 (all quantitative) p are given in Table 2.

438 For the sake of simplicity, the features are chosen as the original variables in
439 the database x_1, x_2, \dots, x_p and their products, yielding monomials of degree
440 up to g . However, other feature spaces, as those proposed by [7], might give
441 better classification rates.

442 Two types of costs are considered for the original variables. For the four
443 databases, costs are independently chosen, randomly in the interval $(0, 1)$.

444 Moreover, for the databases `bupa` and `pima` there exists a file, donated by
445 Turney [31] and publicly available in the UCI repository [6], which contains
446 an example for possible costs for the measurement of the variables. The cost
447 information comes from the Ontario Health Insurance Program’s fee schedule.
448 For these databases we have also considered such given costs. The remaining
449 features have zero cost. The partial order is given as follows: feature $\phi = x_k$
450 precedes all features of the form $\phi(x) = x_k q(x)$ for some monomial $q(x)$ of
451 degree up to $g - 1$.

452 Data were standardized by subtracting its mean and dividing by its standard
453 deviation. Then, from each database, a random sample with two thirds of
454 the objects is drawn and used as training sample I . The supported Pareto-
455 optimal solutions of Problem (32) were computed by the first phase of the
456 Two-Phase Method [34], described in Section 5. The non-supported Pareto-
457 optimal solutions can also be computed using formulation (31). The trade-off
458 parameter γ is chosen to be equal to the number of objects in I .

459 The results are plotted in Figures 2-9. In the right side of such figures, mea-
460 surement costs of the Pareto-optimal rules (except for zero-cost solutions) are
461 plotted against the margin. Since only Pareto-optimal solutions are consid-
462 ered, we see that, the higher the cost, the higher the margin.

463 This is the plot the final user will obtain in real-world applications, and chose,
464 with this information, one classification rule.

465 However, margin maximization is only a surrogate for the minimization of the
466 misclassification rate, which will remain unknown. In the right side of Figures
467 2-9 we have plotted, for the Pareto-optimal classifiers obtained, costs against
468 the percentage of correctly classified objects in the testing sample. Figures
469 show clearly that high correct classification rates correspond to high costs.
470 Moreover, the trade-off between measurement costs and margin translates
471 into a similar trade-off between measurement costs and percentage of correctly
472 classified objects.

473 For comparative purposes, in Table 3, the percentage of correctly classified
474 objects is shown for different classification methods, such as classification
475 trees [10], k -nearest neighbor classifier [12] and the classical SVM approach
476 as implemented in SVMlight [21]. It can be observed that the classification
477 behavior of the Pareto-optimal classifiers are among the best ones, even for
478 low classification costs.

479 The method proposed in this paper, can thus be seen as a procedure that
480 generates a series of classification rules with different costs, and expected good
481 classification behavior supported by the theoretical generalization properties
482 of the margin maximizer (e.g. Vapnik [33]). Choosing one classification rule
483 among them can be done by the user after plotting the measurement costs
484 against margins, as illustrated in the examples.

485 **List of Figures**

486	1	separable, but not OAR-separable	36
487	2	Database ‘bupa’, $g = 1$, random costs.	37
488	3	Database ‘bupa’, $g = 1$, Turney’s costs.	37
489	4	Database ‘pima’, $g = 1$, random costs.	37
490	5	Database ‘pima’, $g = 1$, Turney’s costs.	38
491	6	Database ‘credit’, $g = 1$, random costs.	38
492	7	Database ‘credit’, $g = 2$, random costs.	38
493	8	Database ‘wdbc’, $g = 1$, random costs.	39
494	9	Database ‘wdbc’, $g = 2$, random costs.	39

495 **List of Tables**

496	1	Example of feature cost.	39
497	2	Parameters of the databases. *only the numerical variables	
498		were used.	39
499	3	Behavior of other methods.	40

500 **References**

- 501 [1] E.L. Allwein, R.E. Schapire, and Y. Singer. Reducing multiclass to binary: A
502 unifying approach for margin classifiers. *Journal of Machine Learning Research*,
503 1:113–141, 2000.
- 504 [2] C. Apte. The big (data) dig. *OR/MS Today*, February 2003.
- 505 [3] V. Bayer-Zubek. *Learning Cost-Sensitive Diagnostic Policies from Data*. PhD
506 thesis, Oregon State University, July 2003.
507 <http://eecs.oregonstate.edu/library/?call=2003-13>.
- 508 [4] K.P. Bennet and O.L. Mangasarian. Robust linear programming discrimination
509 of two linearly inseparable sets. *Optimization Methods and Software*, 1:23–24,
510 1992.
- 511 [5] K.P. Bennet. Combining support vector and mathematical programming
512 methods for induction. In B. Scholkopf, C. Burges, and A. Smola, editors,
513 *Advances in Kernel Methods - Support Vector Learning*, pages 307–326. MIT
514 Press, Cambridge, MA, 1999.
- 515 [6] C.L. Blake and C.J. Merz. UCI Repository of Machine Learning Databases.
516 <http://www.ics.uci.edu/~mllearn/MLRepository.html>, 1998. University of
517 California, Irvine, Dept. of Information and Computer Sciences.
- 518 [7] E. Boros, P.L. Hammer, T. Ibaraki, and A. Kogan. A logical analysis of
519 numerical data. *Mathematical Programming*, 79:163–190, 1997.
- 520 [8] P.S. Bradley, U.M. Fayyad, and O.L. Mangasarian. Mathematical programming
521 for data mining: formulations and challenges. *INFORMS Journal on*
522 *Computing*, 11(3):217–238, 1999.
- 523 [9] P.S. Bradley, O. Mangasarian, and D. Musicant. Optimization methods in
524 massive datasets. In J. Abello, P.M. Pardalos, and M.G.C. Resende, editors,
525 *Handbook of Massive Datasets*, pages 439–472. Kluwer Academic Pub., 2002.
- 526 [10] L. Breiman, J.H. Friedmann, R.A. Olshen, and C.J. Stone. *Classification and*
527 *Regression Trees*. Wadsworth, Belmont, CA, 1984.
- 528 [11] C. Cortes and V. Vapnik. Support-vector network. *Machine Learning*, 1(1):113–
529 141, 1995.
- 530 [12] T.M. Cover and P.E. Hart. Nearest neighbor pattern classification. *IEEE*
531 *Transactions on Information Theory*, 13:21–27, 1967.
- 532 [13] N. Cristianini and J. Shawe-Taylor. *An introduction to support vector machines*.
533 Cambridge University Press, 2000.
- 534 [14] A. Demiriz, K.P. Bennett, and J. Shawe-Taylor. Linear programming boosting
535 via column generation. *Machine Learning*, 46(1):225–254, 2002.

- 536 [15] A.P. Duarte Silva and A. Stam. Second order mathematical programming
537 formulations for discriminant analysis. *European Journal of Operational*
538 *Research*, 72:4–22, 1994.
- 539 [16] J.E. Falk and V.E. Karlov. Robust separation of finite sets via quadratics.
540 *Computers and Operations Research*, 28:537–561, 2001.
- 541 [17] Y. Freund and R.E. Schapire. A decision-theoretic generalization of on-line
542 learning and an application to boosting. *Journal of Computer and System*
543 *Sciences*, 55(1):119–139, 1997.
- 544 [18] I. Guyon and A. Elisseeff. An introduction to variable and feature selection.
545 *Journal of Machine Learning Research*, 3(1157-1182), 2003.
- 546 [19] T. Hastie and R. Tibshirani. Classification by pairwise coupling. *The Annals*
547 *of Statistics*, 26(2):451–471, 1998.
- 548 [20] R. Herbrich. *Learning Theory Classifiers. Theory and Algorithms*. MIT Press,
549 2002.
- 550 [21] T. Joachims. *Learning to Classify Text Using Support Vector Machines*. Kluwer,
551 2002.
- 552 [22] O.L. Mangasarian. Linear and nonlinear separation of patterns by linear
553 programming. *Operations Research*, 13:444–452, 1965.
- 554 [23] O.L. Mangasarian. Mathematical programming in data mining. *Data Mining*
555 *and Knowledge Discovery*, 42(1):183–201, 1997.
- 556 [24] S.W. Norton. Generating better decision trees. In *Proceedings of the Eleventh*
557 *International Joint Conference on Artificial Intelligence, IJCAI-89*, pages 800–
558 805, Detroit, Michigan, 1989.
- 559 [25] M. Núñez. The use of background knowledge in decision tree induction. *Machine*
560 *Learning*, 6(3):231–250, 1991.
- 561 [26] P. Paclik, R.P.W. Duin, G.M.P. van Kempen, and R. Kohlus. On feature
562 selection with measurement cost and grouped features. *Lecture Notes in*
563 *Computer Science*, 2396:461–469, 2002.
- 564 [27] J.P. Pedroso and N. Murata. Support vector machines with different norms:
565 motivation, formulations and results. *Pattern Recognition Letters*, 22(12):1263–
566 1272, 2001.
- 567 [28] A.M. Rubinov, A.M. Bagirovand, N.V. Soukhoroukova, and J. Yearwood.
568 Unsupervised and supervised data classification via nonsmooth and global
569 optimization. *TOP*, 11(1):1–93, 2003.
- 570 [29] A. Smola, T.T. Friess, and B. Schlkopf. Semiparametric support vector and
571 linear programming machines. In M.J. Kearns, S.A. Solla, and D.A. Cohn,
572 editors, *Advances in Neural Information Processing Systems 10*, pages 585–591.
573 MIT Press, 1999.

- 574 [30] M. Tan. Cost-sensitive learning of classification knowledge and its applications
575 in robotics. *Machine Learning*, 13(1):7–33, 1993.
- 576 [31] P.D. Turney. Cost-sensitive classification: Empirical evaluation of a hybrid
577 genetic decision tree induction algorithm. *Journal of Artificial Intelligence*
578 *Research*, 2:369–409, 1995.
- 579 [32] V. Vapnik. *The nature of statistical learning theory*. Springer-Verlag, 1995.
- 580 [33] V. Vapnik. *Statistical learning theory*. Wiley, 1998.
- 581 [34] M. Visée, J. Teghem, M. Pirlot, and E.L. Ulungu. Two-phases method and
582 branch and bound procedures to solve the bi-objective knapsack problem.
583 *Journal of Global Optimization*, 12:139–155, 1998.
- 584 [35] J. Weston, A. Gammerman, M.O. Stitson, V. Vapnik, V. Vovk, and C. Watkins.
585 Support vector density estimation. In B. Schölkopf, C. Burges, and A. Smola,
586 editors, *Advances in Kernel Methods - Support Vector Learning*, pages 293 –
587 305. MIT Press, 1999.
- 588 [36] J. Weston, S. Mukherjee, O. Chapelle, M. Pontil, T. Poggio, and V. Vapnik.
589 Feature selection for SVMs. In T.K. Leen, T.G. Dietterich, and V. Tresp,
590 editors, *Advances in Neural Information Processing Systems 13*, pages 668–674.
591 MIT Press, 2001.

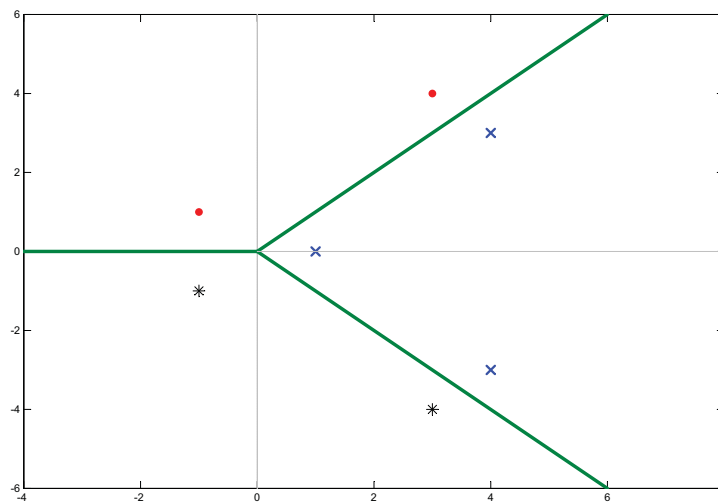


Fig. 1. separable, but not OAR-separable

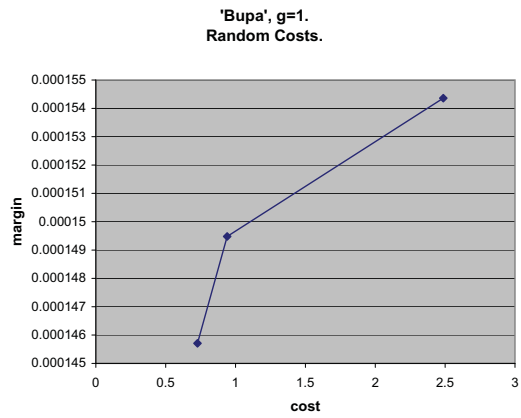
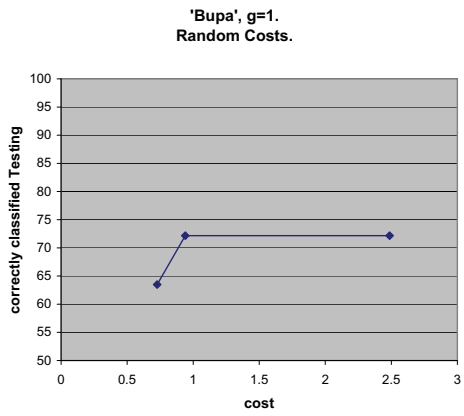


Fig. 2. Database 'bupa', $g = 1$, random costs.

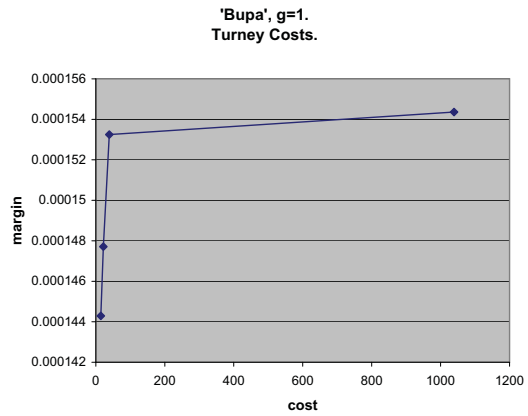
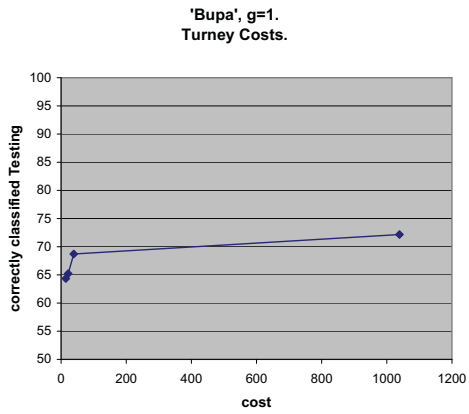


Fig. 3. Database 'bupa', $g = 1$, Turney's costs.

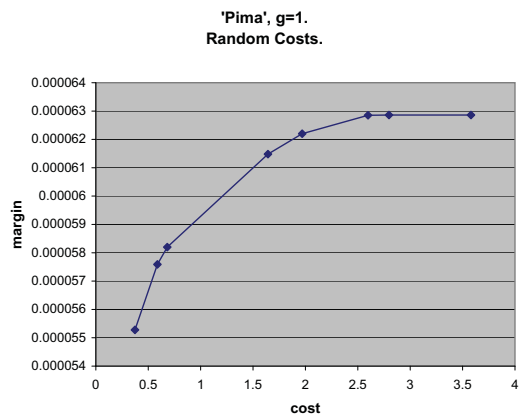
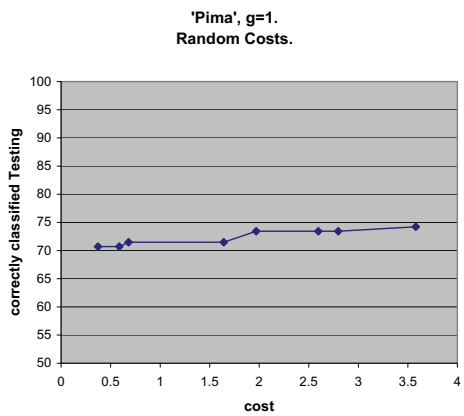


Fig. 4. Database 'pima', $g = 1$, random costs.

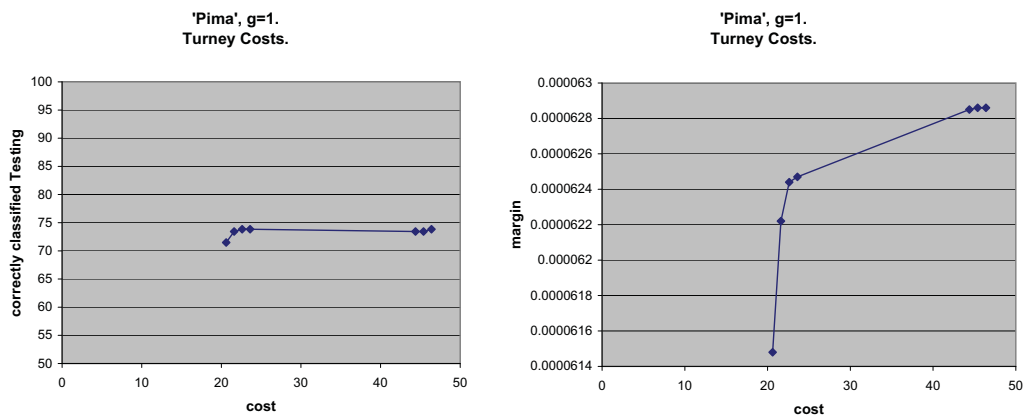


Fig. 5. Database 'pima', $g = 1$, Turney's costs.

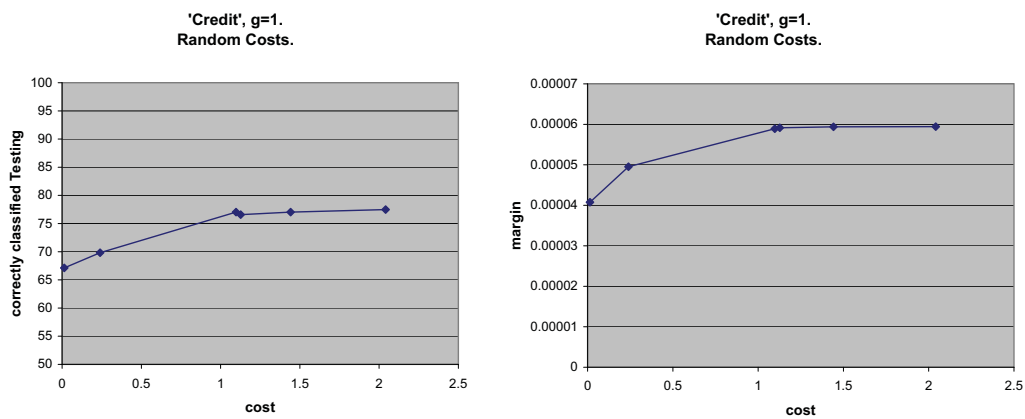


Fig. 6. Database 'credit', $g = 1$, random costs.

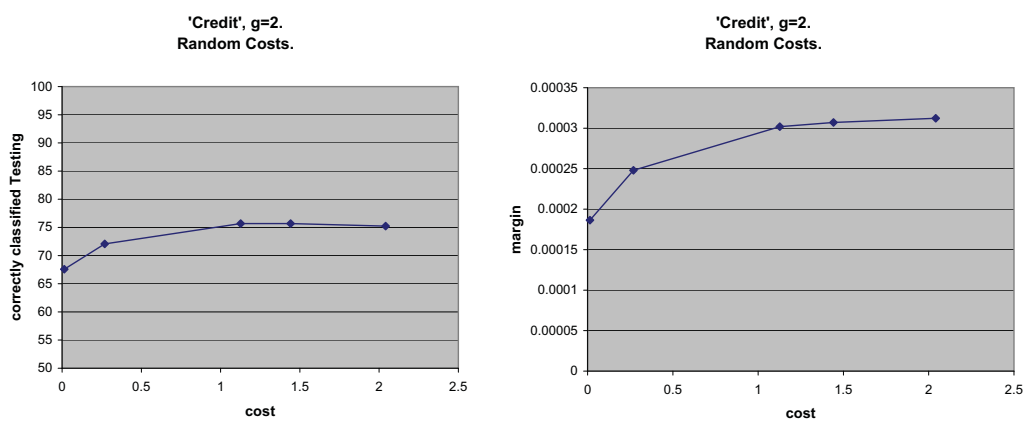


Fig. 7. Database 'credit', $g = 2$, random costs.

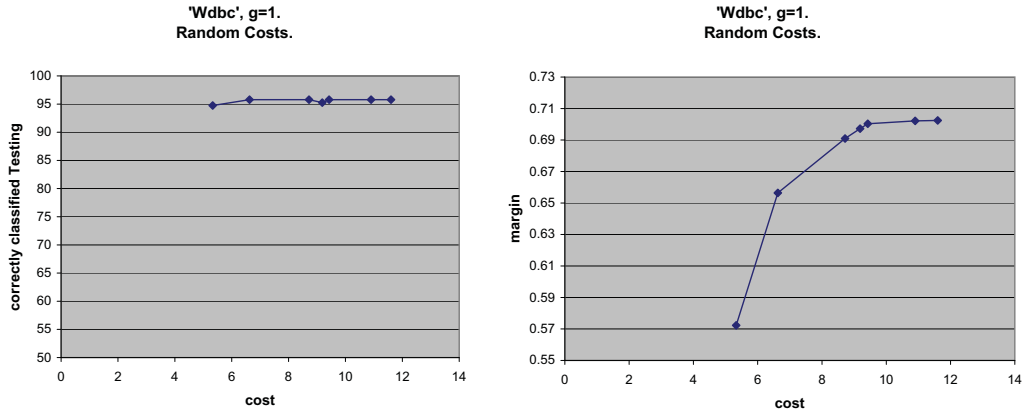


Fig. 8. Database 'wdbc', $g = 1$, random costs.

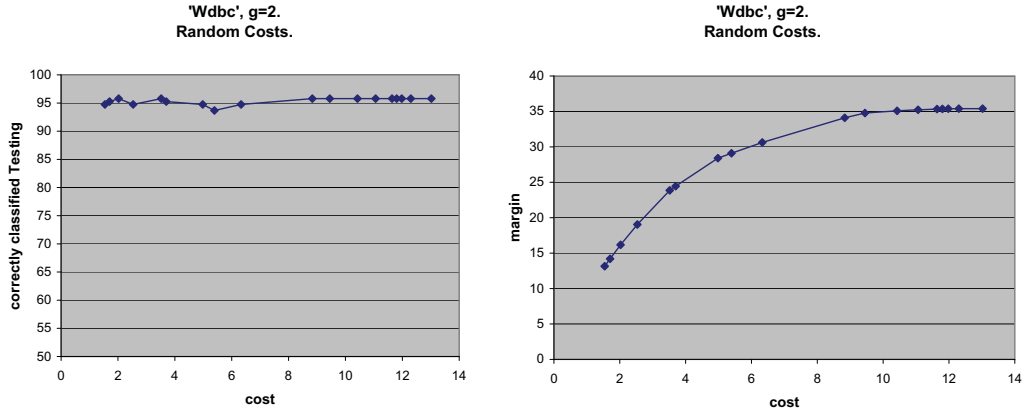


Fig. 9. Database 'wdbc', $g = 2$, random costs.

features	ϕ_1	ϕ_2	ϕ_3	ϕ_4	ϕ_5
costs	2	5	3	0	2

Table 1
Example of feature cost.

Database	filename	$ \Omega $	C	p
bupa	bupa.data	345	2	6
pima	pima-indians-diabetes.data	768	2	8
wdbc	wdbc.data	569	2	30
credit	crx.data*	768	2	8

Table 2
Parameters of the databases. *only the numerical variables were used.

method	'bupa'	'pima'	'wdbc'	'credit'
1-Nearest Neighbor	60.87	64.84	94.74	72.07
2-Nearest Neighbor	57.39	69.14	94.21	70.72
3-Nearest Neighbor	60.00	72.27	95.26	73.87
4-Nearest Neighbor	60.87	72.27	95.26	72.52
5-Nearest Neighbor	62.61	71.48	95.79	72.07
Classification Tree	67.83	70.31	90.53	72.97
SVM with linear kernel	72.17	74.22	95.79	77.48
SVM with polynomial kernel, grade =2	66.96	38.28	94.21	65.32
SVM with polynomial kernel, grade =3	59.13	66.41	93.68	69.37
SVM with polynomial kernel, grade =4	58.26	62.89	89.47	59.01
SVM with polynomial kernel, grade =5	57.39	67.19	91.58	75.23
SVM with radial basis function kernel	68.70	64.84	63.16	77.48

Table 3
Behavior of other methods.