

A Systematic Review of Multi-Label Feature Selection and a New Method Based on Label Construction

Newton Spolaôr^{a,b,*}, Maria Carolina Monard^a, Grigorios Tsoumakas^c, Huei
Diana Lee^b

^a*Laboratory of Computational Intelligence,
Institute of Mathematics and Computer Science, University of São Paulo, Brazil.
ZIP code: 13560-970. Tel.: +55 16 3373-9646. Fax: +55 16 3373-9751*

^b*Laboratory of Bioinformatics,
State West Paraná University, Foz do Iguaçu, Brazil.
ZIP code: 85867-900. Tel.: +55 45 3576-8815. Fax: +55 45 3575-2733*

^c*Machine Learning and Knowledge Discovery Group,
Department of Informatics, Aristotle University of Thessaloniki, Greece.
ZIP code: 54124. Tel.: +30 23 1099-8145*

Abstract

Each example in a multi-label dataset is associated with multiple labels, which are often correlated. Learning from this data can be improved when dimensionality reduction tasks, such as feature selection, are applied. The standard approach for multi-label feature selection transforms the multi-label dataset into single-label datasets before using traditional feature selection algorithms. However, this approach often ignores label dependence. In this work, we propose an alternative method, LCFS, that constructs new labels based on relations between the original labels. By doing so, the label set from the data is augmented with second-order information before applying

*Corresponding author

Email addresses: newtonspolaor@gmail.com (Newton Spolaôr),
mcmonard@icmc.usp.br (Maria Carolina Monard), greg@csd.auth.gr (Grigorios
Tsoumakas), hueidianalee@gmail.com (Huei Diana Lee)

the standard approach. To assess LCFS, an experimental evaluation using Information Gain as a measure to estimate the importance of features was carried out on 10 benchmark multi-label datasets. This evaluation compared four LCFS settings with the standard approach, using random feature selection as a reference. For each dataset, the performance of a feature selection method is estimated by the quality of the classifiers built from the data described by the features selected by the method. The results show that a simple LCFS setting gave rise to classifiers similar to, or better than, the ones built using the standard approach. Furthermore, this work also pioneers the use of the systematic review method to survey the related work on multi-label feature selection. The summary of the 99 papers found promotes the idea that exploring label dependence during feature selection can lead to good results.

Keywords: feature ranking, filter feature selection, binary relevance, information gain, systematic review

1. Introduction

In multi-label learning, each example is associated with multiple labels simultaneously. A key difference between multi-label and traditional binary or multi-class single-label learning is that the labels in multi-label learning are not mutually exclusive. Thus, in comparison with traditional single-label learning, multi-label learning is more general and challenging to solve. The issue of learning from multi-label data has attracted significant attention from the community, motivated by an increasing number of new applications in bioinformatics [1, 2], emotion analysis [3], text mining [4, 5] and image

10 analysis [6], among others.

11 As other machine learning tasks, multi-label learning also suffers from
12 the “*curse of dimensionality*”. Dimensionality reduction (feature selection),
13 which aims to find a small subset of features that describes the dataset as
14 well as, or even better than, the original set of features does [7], is an effective
15 way to mitigate the curse of dimensionality.

16 The standard approach for multi-label Feature Selection (FS), which
17 transforms the multi-label dataset into single-label datasets before using tra-
18 ditional FS algorithms, is implementable in the Binary Relevance (*BR*) ap-
19 proach [8]. However, a *BR* drawback is that label dependence is often ignored.
20 Thus, a significant challenge regarding this approach is how to explore the
21 labels structure to improve multi-label learning performance simultaneously
22 with dimensionality reduction.

23 An alternative to overcoming this problem would be to construct labels
24 based on relations among the original labels and include the new labels dur-
25 ing the feature selection phase. The main idea of variable (label or feature)
26 construction is to gather information about the relations among the original
27 variables from data and infer additional variables [9]. Although feature con-
28 struction methods are less usual than feature selection methods [10], they
29 have already been used to support single-label [11, 12] and multi-label learn-
30 ing [13, 14, 15]. Nevertheless, to the best of our knowledge, there is little
31 research on *label construction* for multi-label data.

32 In this work, we present the Label Construction for Feature Selection
33 (*LCFS*) method, originally proposed in [16], to build binary variables (new
34 labels) based on label relations. These variables are then included as new

35 labels in the original dataset and the standard multi-label FS approach based
36 on *BR* is used in the augmented dataset to select features. Afterwards, the
37 dataset described by the selected features and the original labels can be
38 submitted to any multi-label learning algorithm.

39 The *LCFS* method was experimentally compared with the standard multi-
40 label FS approach based on *BR* on 10 benchmark datasets. We also used
41 Random Feature Selection (*RFS*) as a reference. Both *LCFS* and the stan-
42 dard approach consider the frequently used measure Information Gain (*IG*)
43 to evaluate features.

44 The experimental results suggest that setting *LCFS* with simple strate-
45 gies to build binary variables (new labels) from pairs of labels gives rise to
46 classifiers similar to, or better than, the ones built using the standard ap-
47 proach based on *BR*. Good *LCFS* results are also observed when the number
48 of features selected is small.

49 Furthermore, we also applied the Systematic Review (SR) method [17]
50 to survey the literature on multi-label FS. The summary of the 99 papers
51 found shows that good results are obtained from research which takes into
52 account label dependence. Another finding is that *IG* is the most frequently
53 used importance measure, as can be observed in 23 out of the 99 papers.

54 The rest of this paper is organized as follows: Section 2 briefly describes
55 multi-label learning and feature selection. It also presents the systematic
56 review method, the *LCFS* method and the multi-label datasets used in the
57 experimental evaluation. Section 3 describes the experimental setting used
58 to obtain the results discussed in Section 4. Section 5 shows the related
59 work found by applying the SR method. Section 6 concludes the paper and

60 highlights future work.

61 2. Material and methods

62 This section briefly describes multi-label learning and feature selection,
63 as well as the systematic literature review method. It also describes the
64 *LCFS* feature selection method and the characteristics of the 10 benchmark
65 multi-label datasets used in the experimental evaluation.

66 2.1. Multi-label learning

67 Let D be a dataset composed of N examples $E_i = (\mathbf{x}_i, Y_i)$, $i = 1 \dots N$.
68 Each example E_i is associated with a feature vector $\mathbf{x}_i = (x_{i1}, x_{i2}, \dots, x_{iM})$
69 described by M features (attributes) X_j , $j = 1 \dots M$, and its multi-label Y_i ,
70 which consists of a subset of labels $Y_i \subseteq L$, where $L = \{y_1, y_2, \dots, y_q\}$ is the
71 set of q labels. Table 1 shows this representation. In this scenario, the multi-
72 label classification task consists of generating a classifier H which, given an
73 unseen example $E = (\mathbf{x}, ?)$, is capable of accurately predicting its multi-label
74 Y , *i.e.*, $H(E) \rightarrow Y$.

Table 1: Multi-label data

| | X_1 | X_2 | \dots | X_M | Y |
|----------|----------|----------|----------|----------|----------|
| E_1 | x_{11} | x_{12} | \dots | x_{1M} | Y_1 |
| E_2 | x_{21} | x_{22} | \dots | x_{2M} | Y_2 |
| \vdots | \vdots | \vdots | \ddots | \vdots | \vdots |
| E_N | x_{N1} | x_{N2} | \dots | x_{NM} | Y_N |

75 2.1.1. Categorizing multi-label learning algorithms

76 Multi-label learning methods can be organized into two main categories [8]:
77 problem transformation and algorithm adaptation. The key philosophy for

78 the former is to fit data to algorithms, while for the latter is to fit algorithms
79 to data [18]. In particular:

- 80 • Problem transformation methods decompose the multi-label learning
81 problem into a set of single-label (binary or multi-class) learning tasks;
- 82 • Algorithm adaptation methods adapt specific learning algorithms to
83 handle multi-label datasets directly.

84 The multi-label learning algorithm *BRkNN*, which modifies the single-
85 label *lazy k*-Nearest Neighbor (*kNN*) algorithm to classify multi-label ex-
86 amples, belongs to the algorithm adaptation category. To better tackle the
87 multi-label problem, the extensions *BRkNN-a* and *BRkNN-b* are proposed
88 in [19]. Both extensions are based on a label confidence score, which is esti-
89 mated for each label from the percentage of the *k*-Nearest Neighbors having
90 this label. *BRkNN-a* classifies a new example *E* using the labels in the multi-
91 labels of the *k*-Nearest Neighbors which have a confidence score greater than
92 0.5, *i.e.*, labels included in the multi-labels of at least half of the *k*-Nearest
93 Neighbors of *E*. If no label satisfies this condition, it outputs the label with
94 the greatest confidence score. On the other hand, *BRkNN-b* classifies *E* with
95 the $\lceil s \rceil$ (nearest integer of *s*) labels that have the greatest confidence score,
96 where *s* is the average size of the multi-labels of the *k*-Nearest Neighbors of
97 *E*. By conducting an experimental comparison with the state-of-the-art *lazy*
98 algorithm *MLkNN* [20], the authors found that *BRkNN-b* achieved compet-
99 itive results.

100 As *lazy* learning algorithms are sensitive to irrelevant features, they are
101 a good choice to indicate the quality of a feature selection method. Thus, in

102 this work, we use *BRkNN-b* to assess the quality of the classifiers built using
103 the original datasets — All Features (AF) — and the classifiers built using
104 the datasets described by the selected features.

105 As exploring label dependence during learning could improve the classi-
106 fier performance [21], Zhang and Zhou [18] proposed another categorization
107 of multi-label learning methods which takes into account the degree of la-
108 bel dependence exploration. First-order strategies ignore the co-existence of
109 other labels. The Binary Relevance (*BR*) approach, a problem transforma-
110 tion method, exemplifies this category by transforming a multi-label dataset
111 into q single-label binary datasets, learning from each single-label problem
112 separately and combining the results. Second-order strategies can consider
113 pairwise relations between labels, such as interactions between any pair of
114 labels, or the ranking between relevant and irrelevant labels. High-order
115 strategies consider relations among more labels.

116 Although high-order strategies potentially model wider label dependences,
117 they are usually computationally more demanding. This work focuses on
118 finding second-order relations between single labels from the multi-label dataset
119 and representing them as new labels. The idea is that, by labeling examples
120 with the original and the constructed labels, feature selection methods based
121 on the *BR* approach to incorporate label pairwise information are feasible.

122 2.1.2. Evaluation Measures

123 Unlike single-label classification where the classification of a new exam-
124 ple has only two possible outcomes, correct or incorrect, multi-label clas-
125 sification should also take into account *partially* correct classification. As
126 a consequence, multi-label evaluation measures consider the performance of

127 the classifier from diverse aspects and are, thus, of a different nature. A
 128 complete discussion on multi-label evaluation measures is out of the scope
 129 of this work and can be found in [8]. In what follows, we describe the four
 130 evaluation measures used in this work.

131 *F-measure*, *Hamming loss* and *Accuracy*, defined by Equations 1 to 3,
 132 are example-based evaluation measures, where Δ represents the symmetric
 133 difference of two sets, Y_i and Z_i are the true and the predicted multi-labels
 134 respectively.

$$F\text{-measure}(H, D) = \frac{1}{|D|} \sum_{i=1}^{|D|} \frac{2|Y_i \cap Z_i|}{|Z_i| + |Y_i|}. \quad (1)$$

$$Hamming\ loss(H, D) = \frac{1}{|D|} \sum_{i=1}^{|D|} \frac{|Y_i \Delta Z_i|}{|L|}. \quad (2)$$

$$Accuracy(H, D) = \frac{1}{|D|} \sum_{i=1}^{|D|} \frac{|Y_i \cap Z_i|}{|Y_i \cup Z_i|}. \quad (3)$$

135 In addition, *Micro-averaged F-measure* (F_b), defined by Equation 4, is a
 136 label-based measure, where $T_{P_{y_i}}$, $F_{P_{y_i}}$, $T_{N_{y_i}}$ and $F_{N_{y_i}}$ represent, respectively,
 137 the number of true/false positives/negatives for a label $y_j \in L$.

$$F_b(H, D) = \frac{2 \sum_{j=1}^q T_{P_{y_j}}}{2 \sum_{j=1}^q T_{P_{y_j}} + \sum_{j=1}^q F_{P_{y_j}} + \sum_{j=1}^q F_{N_{y_j}}}. \quad (4)$$

138 All these performance measures range in the interval $[0, 1]$. For *Hamming loss*,
 139 the smaller the value, the better the multi-label classifier performance is,
 140 while for the other measures, greater values indicate better performance.

141 *2.2. Feature selection*

142 Regardless of the multi-label learning approach, any FS method addresses
143 a few relevant issues, such as the interaction with the learning algorithm
144 and the feature importance measure. The first issue is taken into account
145 in different ways by the wrapper, embedded and filter approaches. The
146 wrapper and the embedded approaches involve interaction with the learning
147 algorithm, such that the features are selected for a specific learning algorithm.
148 On the other hand, the filter approach uses general properties of the dataset
149 to remove unimportant features from it, regardless of the learning algorithm.
150 Thus, the features chosen may not be the best ones for a specific learning
151 algorithm. The FS algorithms considered in this work use the filter approach.

152 Many measures have been proposed to estimate the importance of features
153 based on characteristics of the dataset. As Section 5 reports, a frequently
154 used single-label FS measure is Information Gain (*IG*), which evaluates each
155 feature according to the dependence between this feature and a single label,
156 as defined by Equation 5 — the higher the *IG* value for a feature X_j , the
157 stronger is the relationship between X_j and the label.

$$IG(D, X_j) = entropy(D) - \sum_v \frac{|D_v| entropy(D_v)}{|D|}. \quad (5)$$

158 In other words, the *IG* of feature X_j , $j = 1 \dots M$, calculates the difference
159 between the entropy of dataset D and the weighted sum of the entropy of
160 each subset $D_v \subseteq D$, where D_v consists of the set of examples where X_j has

161 the value v . Therefore, if X_j has 10 distinct values¹ in D , the sum would be
162 applied to 10 different D_v datasets.

163 Using the *BR* approach, any single-label FS measure can be used to select
164 features from multi-label data, as shown in [22], in which the single label FS
165 measures *IG* and ReliefF are used. The procedure is simple: initially, using
166 the *BR* approach the multi-label dataset is transformed into q single-label
167 datasets, one per label. Afterwards, the single-label FS measure is applied
168 to each feature X_j having as single-label y_i , $i = 1 \dots q$, and the q results are
169 averaged to obtain the final importance value of feature X_j . Finally, the
170 importance value of the M features could be ranked to guide the selection of
171 the better subset of features.

172 In this work, we also use Random Feature Selection (*RFS*) as a reference,
173 in which the features are randomly selected, *i.e.*, no label or multi-label
174 information is considered. Next, the Systematic Review (SR) method used
175 to survey the literature on multi-label feature selection is briefly described.

176 2.3. Systematic review

177 The systematic literature review method provides a rigorous and repli-
178 cable process to review the evidence relevant to a particular research ques-
179 tion [17]. Although this method emerged in areas such as Medicine, currently
180 there are guidelines and applications in other areas. In Computer Science,
181 several applications can be found, including a systematic review that surveys
182 other systematic reviews [23]. Figure 1 summarizes the workflow of the three
183 systematic review steps — planning, conducting and reporting — as well as

¹Discretization is applied to numerical features before using *IG*.

184 its main inputs and outputs.

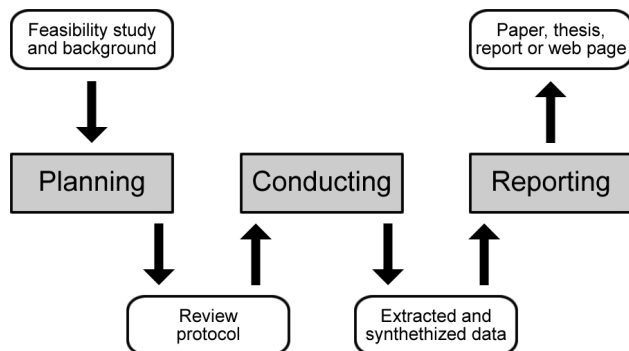


Figure 1: An overview of the systematic literature review method

185 The input of the planning step consists of a feasibility study about the sys-
186 tematic review method and the background related to the research question.
187 In particular, the feasibility study can be carried out by analyzing whether
188 or not a systematic review is needed concerning the topic of interest.

189 Planning yields a review protocol, which describes the systematic review
190 components and improves the method replicability. This protocol should be
191 consulted to conduct the method, which provides the data able to answer
192 specific research question(s), which is (are) the SR core. In the last step,
193 reporting, this data is disseminated in papers or a Ph.D. thesis, for example.

194 Although there are a few surveys on multi-label learning [18, 8], to the
195 best of our knowledge there is no previous SR neither on multi-label learning
196 nor on multi-label feature selection. This motivated us to conduct a pioneer-
197 ing systematic review to rigorously survey the related work on multi-label
198 feature selection.

199 The systematic review research question we aim to answer is “*what are*
200 *the publications of feature selection in multi-labeled data?*”. Further details

201 regarding the review protocol and the instantiation of the SR method for
202 multi-label FS are described in [24]. The reporting step to disseminate the
203 related work summary is described in Section 5.

204 2.4. The LCFS method

205 Given a multi-label dataset D with the set of single labels $L =$
206 $\{y_1, y_2, y_3, \dots, y_q\}$, the main idea of *LCFS* is to construct q' new single labels
207 by combining the original labels within pairs (y_i, y_j) , $i \neq j$, $y_i \in L$ and $y_j \in L$.
208 In each iteration, *LCFS* selects a pair of labels (y_i, y_j) from L and combines
209 the labels within this pair to generate a new label y_{ij} . After repeating this
210 procedure q' times, the q' new labels are included in the label set L , such
211 that information about pairwise relationships between original labels can be
212 used by the *BR* approach for feature selection.

213 The *LCFS* method consists of two steps, each one concerned with an-
214 swering a different question:

- 215 1. Selection: which pairs of labels (y_i, y_j) should be chosen?
- 216 2. Generation: how to combine these labels to generate the new labels y_{ij} ?

217 Figure 2 illustrates these steps for $q' = 1$.

218 Thus, instantiating *LCFS* involves choosing a strategy to select label
219 pairs and a strategy to combine the labels within each pair. An additional
220 parameter is the number of new labels q' that will be constructed. In what
221 follows, the two *LCFS* steps are described.

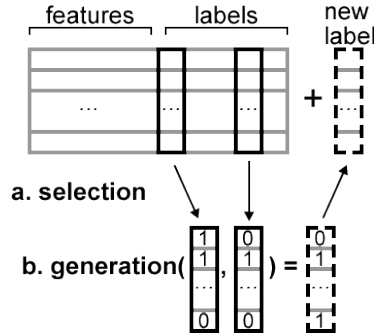


Figure 2: Applying the two steps of *LCFS* to construct $q' = 1$ new labels [16]

222 *2.4.1. Step 1: selection*

223 Given the set of labels $L = \{y_1, y_2, y_3, \dots, y_q\}$ of dataset D , *LCFS* chooses
 224 q' different pairs of labels² (y_i, y_j) , $i \neq j$, according to a selection strategy.
 225 The idea is that these pairs capture some pairwise relationships between the
 226 labels to be considered by feature selection.

227 *LCFS* supports different selection strategies, such as the simple Random
 228 Selection (*RS*), as well as heuristic strategies based on the number of ex-
 229 amples labeled by each original single label (label frequency). In particular,
 230 two strategies considering label frequency are Co-occurrence-based Selection
 231 (*CS*) and related Labels Selection (*LS*). *CS* sorts in descending order label
 232 pairs according to the co-occurrence c_c , *i.e.*, the number of examples labeled
 233 by both labels within a pair — $(1, 1)$ — , and selects the first q' different
 234 pairs. On the other hand, *LS* counts:

- 235 1. The number of examples in which the labels within a pair agree, c_e —
 236 $(1, 1)$ or $(0, 0)$;

²In this work, two label pairs are considered different if they do not have a common label. For example, (y_3, y_5) and (y_1, y_3) are not considered different pairs because they share the label y_3 .

237 2. The number of examples in which the labels within a pair disagree,
238 c_d — (1, 0) or (0, 1).

239 Then, the pairs are sorted, in descending order, into two lists according to
240 the values of c_e and c_d . The pair with the greatest value is selected, removed
241 from the correspondent list and the procedure is repeated until selecting q'
242 *different pairs*.

243 2.4.2. Step 2: generation

244 In this step, *LCFS* combines both labels from all previously selected pairs
245 (y_i, y_j) , $i \neq j$, to construct the new labels y_{ij} . The idea is that the values
246 of y_{ij} represent a pairwise relationship between y_i and y_j . In the end, all
247 examples in D are labeled by the q original labels and the q' new labels. *LCFS*
248 supports different combination strategies between binary variables (labels).
249 In this work, we use three simple logical operators to generate the values of
250 the new labels of each example in D . The logical operators are:

251 **AND** : $y_{ij} = 1$ iff $y_i = y_j = 1$; $y_{ij} = 0$ otherwise.

252 **XOR** : $y_{ij} = 1$ iff $y_i \neq y_j$; $y_{ij} = 0$ otherwise.

253 **XNOR** : $y_{ij} = 1$ iff $y_i = y_j$; $y_{ij} = 0$ otherwise.

254 The AND operator clearly highlights co-occurring labels. XNOR, also
255 known as the coincidence function, assigns the value 1 to y_{ij} *iff* the labels
256 y_i and y_j agree, whereas XOR does the opposite. Although other logical
257 operators, such as OR, could be included, we consider that AND, XOR and
258 XNOR are enough to represent relations between the original labels.

259 Finally, after generating the q' new labels, the traditional *BR* approach
 260 for FS can be applied to the dataset now labeled by the $q + q'$ labels. Note
 261 that, by combining *BR* with *LCFS*, any single-label FS algorithm can be
 262 applied to the augmented dataset with second-order label information [18].

263 The *LCFS* method was implemented in Mulan [25], a multi-label
 264 learning package based on Weka [26], and is available to the community
 265 at [http://www.labic.icmc.usp.br/pub/mcmonard/Implementations/](http://www.labic.icmc.usp.br/pub/mcmonard/Implementations/Multilabel/lcfs.zip)
 266 [Multilabel/lcfs.zip](http://www.labic.icmc.usp.br/pub/mcmonard/Implementations/Multilabel/lcfs.zip).

267 2.4.3. Illustrative example of the selection strategies

268 The application of the simple logical operators AND, XOR and XNOR
 269 to generate labels is straightforward. To illustrate the strategies to select
 270 *different pairs* of labels, consider the multi-label dataset described in Table 2,
 271 with $L = \{y_1, y_2, y_3, y_4\}$, and the number of new labels to be constructed
 272 $q' = \frac{q}{2} = 2$.

Table 2: Illustrative dataset for the *LCFS* method

| | y_1 | y_2 | y_3 | y_4 |
|-------|-------|-------|-------|-------|
| E_1 | 1 | 0 | 1 | 0 |
| E_2 | 1 | 1 | 0 | 0 |
| E_3 | 0 | 0 | 0 | 1 |
| E_4 | 1 | 1 | 1 | 0 |
| E_5 | 1 | 0 | 0 | 1 |
| E_6 | 1 | 1 | 0 | 1 |
| E_7 | 0 | 1 | 0 | 1 |

273 *RS* randomly selects $q' = 2$ *different pairs*. On the other hand, *CS* and
 274 *LS* sort the pairs of labels (y_i, y_j) , $i \neq j$, in descending order according to the
 275 number of examples fulfilling a specific condition — c_c , c_e and c_d . Table 3
 276 shows the c_c , c_e and c_d values calculated by *CS* and *LS* for each label pair

277 involving y_1 , y_2 , y_3 and y_4 .

Table 3: Number of examples fulfilling specific conditions for each label pair

| | (y_1, y_2) | (y_1, y_3) | (y_1, y_4) | (y_2, y_3) | (y_2, y_4) | (y_3, y_4) |
|-------|--------------|--------------|--------------|--------------|--------------|--------------|
| c_c | 3 | 2 | 2 | 1 | 2 | 0 |
| c_e | 4 | 4 | 2 | 3 | 3 | 1 |
| c_d | 3 | 3 | 5 | 4 | 4 | 6 |

278 First, CS selects the pair (y_1, y_2) , which has the highest co-occurrence
 279 (c_c value). Then it considers the next pair in its ordered list. As (y_1, y_3) is
 280 not a pair different from (y_1, y_2) due to the label y_1 , CS goes to the next
 281 list element, (y_1, y_4) . This procedure is performed successively until finding
 282 (y_3, y_4) , which is a pair different from the label pair previously selected. As
 283 $q' = 2$ pairs of labels were selected, the CS selection strategy ends.

284 LS compares the frequencies (numbers of examples) c_e and c_d from the
 285 first label pair in each ordered list: (y_1, y_2) and (y_3, y_4) . As $c_d(y_3, y_4) >$
 286 $c_e(y_1, y_2)$, only (y_3, y_4) is selected. The procedure goes to the next label pair
 287 in the list from which (y_3, y_4) was selected, *i.e.*, the list sorted according
 288 to c_d . However, as the current label pair, (y_1, y_4) , is not different from the
 289 pair previously selected due to the label y_4 , the procedure moves to the
 290 next iteration. As $c_d(y_2, y_3) = c_e(y_1, y_2)$ and (y_1, y_2) is a *different pair*, the
 291 strategy selects (y_1, y_2) before ending.

292 2.5. Multi-label datasets

293 Table 4 summarizes the characteristics of the 10 datasets used in this
 294 work. For each dataset, it shows: dataset name (Dataset); dataset domain
 295 (Domain); number of examples (N); number of features (M); feature type
 296 ($Type$); number of labels ($|L|$); label cardinality (LC), which is the average
 297 number of labels associated with each example; label density (LD), which is

298 the cardinality normalized by $|L|$; and the number of different multi-labels
 299 ($\#Diff$).

Table 4: Dataset description

| Dataset | Domain | N | M | Type | $ L $ | LC | LD | $\#Diff$ |
|-----------------------|---------|-------|------|----------|-------|--------|-------|----------|
| 1- <i>Cal500</i> | music | 502 | 68 | numeric | 174 | 26.044 | 0.150 | 502 |
| 2- <i>Corel5k</i> | image | 5000 | 499 | discrete | 374 | 3.522 | 0.009 | 3175 |
| 3- <i>Corel16k001</i> | image | 13766 | 500 | discrete | 153 | 2.859 | 0.019 | 4803 |
| 4- <i>Emotions</i> | music | 593 | 72 | numeric | 6 | 1.869 | 0.311 | 27 |
| 5- <i>Fapesp</i> | text | 332 | 8669 | discrete | 66 | 1.774 | 0.027 | 206 |
| 6- <i>Genbase*</i> | biology | 662 | 1185 | discrete | 27 | 1.252 | 0.046 | 32 |
| 7- <i>Llog-f*</i> | text | 1253 | 1004 | discrete | 75 | 1.375 | 0.018 | 303 |
| 8- <i>Magtag5k</i> | music | 5260 | 68 | numeric | 136 | 4.839 | 0.036 | 4163 |
| 9- <i>Scene</i> | image | 2407 | 294 | numeric | 6 | 1.074 | 0.179 | 15 |
| 10- <i>Yeast</i> | biology | 2417 | 103 | numeric | 14 | 4.237 | 0.303 | 198 |

300 Except for datasets *5-Fapesp* and *8-Magtag5k*, the other datasets are
 301 available in the Mulan³ and Meka⁴ repositories. In particular, *5-Fapesp* was
 302 built by members of our research laboratory⁵ [27]. Dataset *8-Magtag5k*⁶ is
 303 further described in [28]. Furthermore, *6-Genbase** and *7-Llog-f** are pre-
 304 processed versions of the publicly available datasets in which an identification
 305 feature and unlabeled examples, respectively, were removed.

306 Besides the dataset characteristics shown in Table 4, information related
 307 to label frequency is also important to characterize multi-label datasets. To
 308 this end, we use quartiles⁷ to describe the datasets label frequency distribu-
 309 tion.

310 Figure 3 depicts the single label frequencies for each dataset by boxplots.
 311 Recall that the bottom and the top of the box are the first and third quartiles,

³<http://mulan.sourceforge.net/datasets.html>

⁴<http://meke.sourceforge.net/#datasets>

⁵The dataset can be obtained from the authors.

⁶<http://tl.di.fc.ul.pt/t/magtag5k.zip>

⁷One of the three values that divides a sorted group of data into four equal parts, each one with 25% of the data.

312 and the band inside the box is the second quartile. Thus, the spacing between
 313 the different parts of the box indicates the degree of dispersion, as well as
 314 the skewness in the dataset. Moreover, the minimum and maximum label
 315 frequencies are also shown.

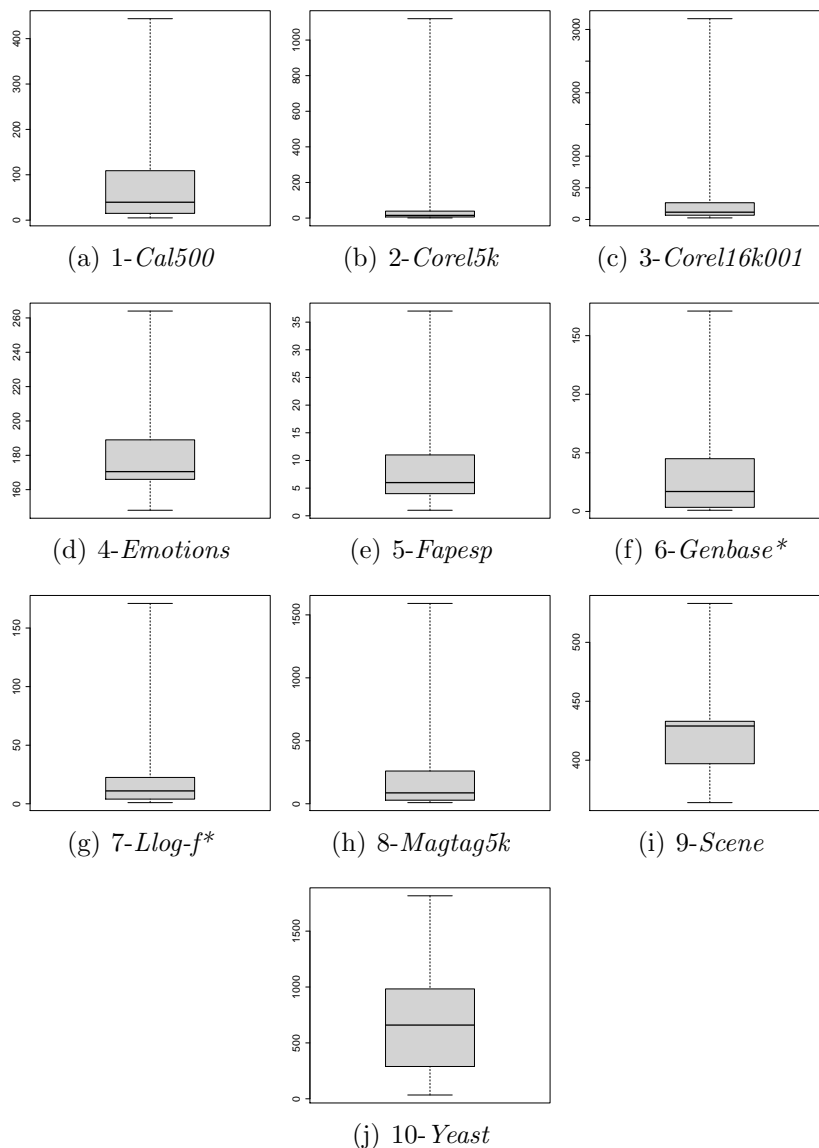


Figure 3: Boxplots of the single-label frequencies for each dataset

316 As can be observed, there is a large variation in the label frequency
317 across the 10 multi-label datasets. It is also worth observing that datasets 3-
318 *Corel16k001* (Figure 3(c), max.=3170, min.=25) and 5-*Fapesp* (Figure 3(e),
319 max.=37, min.=1) show, respectively, the highest (3145) and the smallest
320 (36) absolute difference between the maximum and the minimum label fre-
321 quencies.

322 3. Experimental setting

323 The experiments were carried out on the 10 multi-label datasets described
324 in Table 4. The performance of a FS method was assessed by the *BRkNN-b*
325 classifiers built using the features selected by the method. In particular, the
326 four evaluation measures described in Section 2.1.2 are used to assess the
327 quality of these classifiers, as well as the classifiers built using All Features
328 (AF). The evaluation measures were estimated according to the 10-fold cross-
329 validation strategy.

330 The number of nearest neighbors k was set as 10 for all datasets. This
331 is a commonly-used value that leads lazy learning algorithms to achieve sat-
332 isfactory results [20, 29]. On the other hand, this differs from our previous
333 work [16], in which the goal was to find the k value of *BRkNN-b* that max-
334 imized the *Example-based F-measure* value achieved by the classifier built
335 from each original dataset. All the remaining parameters related to classi-
336 fication and feature selection were executed with the default values used by
337 the Mulan⁸ [25] and Weka⁹ [26] frameworks.

⁸<http://mulan.sourceforge.net>

⁹<http://www.cs.waikato.ac.nz/ml/weka>

338 In this work, a combination between Information Gain and Binary Relevance ($IG-BR$) is performed according to the filter approach:
 339

- 340 1. In the dataset annotated with the original set of labels (standard approach);
- 341
- 342 2. In the dataset annotated with the original labels and the ones constructed by a $LCFS$ setting.
- 343

344 Figure 4 depicts some differences between these approaches. Moreover, by including second-order label information, $LCFS$ is considered as a second-order strategy for multi-label FS, whereas the original $IG-BR$ is considered
 345
 346 as a first-order strategy.
 347

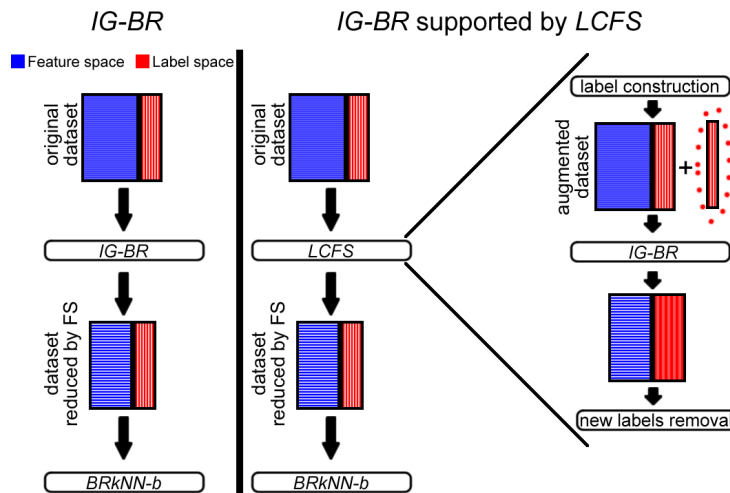


Figure 4: An overview of the FS process conducted by the original $IG-BR$ (standard approach) and by $IG-BR$ supported by $LCFS$

348 Regardless of the label set used, $IG-BR$ transforms a multi-label dataset
 349 into single-label datasets, applies IG to each single-label dataset and averages
 350 the IG score of each feature X_j , $j = 1 \dots M$, across all labels. The resulting
 351 feature ranking sorts the M averaged IG values in descending order. Recall

352 that the labels constructed by *LCFS*, which are only used to select features,
 353 are removed before classification.

354 It should be observed that the averaging strategy used in this work to
 355 aggregate *IG* scores in *BR* was highlighted in an experimental comparison
 356 on 20 multi-label textual datasets [30].

357 Table 5 shows the four settings combining different selection and genera-
 358 tion strategies considered by *LCFS*.

Table 5: *LCFS* settings evaluated in this work

| Setting | Selection | Generation |
|-------------|-----------|--|
| <i>LS-X</i> | <i>LS</i> | XOR or XNOR is chosen based on the lists sorted by the values of c_e and c_d |
| <i>CS-A</i> | <i>CS</i> | AND |
| <i>RS-A</i> | <i>RS</i> | AND |
| <i>RS-X</i> | <i>RS</i> | XOR or XNOR is randomly chosen |

359 Recall that the *LS* selection strategy sorts the label pairs based on the
 360 values of c_e and c_d , *i.e.*, label agreement and disagreement. For a given
 361 label pair (y_i, y_j) , *LS-X* applies the XNOR operator to generate the new
 362 label y_{ij} if the pair was selected from the list sorted by c_e ; otherwise, it
 363 applies the XOR operator. It should be emphasized that *LS-X* and *CS-A*
 364 consider a relationship between the selection and the generation strategies,
 365 as they take into account label agreement/disagreement and co-occurrence
 366 respectively. Finally, *RS-X* randomly selects the XOR or XNOR operator.
 367 See Section 2.4.3 for an illustrative example.

368 We set the number of new labels $q' = \lfloor \frac{q}{2} \rfloor$, *i.e.*, every single label is selected
 369 once if q is even, or one single label is left out if q is odd.

370 In both cases, *IG-BR* and the four *LCFS* settings, the feature subsets
 371 $X' \subset X$, $|X'| = 10\%M, 20\%M, \dots, 90\%M$, ranked by each FS method, are
 372 used to describe a dataset. This dataset is then submitted to the multi-label

373 learning algorithm *BRkNN-b*.

374 Moreover, the *RFS* method, which was not considered in our previous
375 work [16], is included as a reference. In this case 10% M up to 90% M features
376 are randomly chosen from the M features. *RFS* is executed three times per
377 fold, due to its stochasticity, and the three outputs are averaged to yield the
378 result of each fold.

379 As previously mentioned, multi-label evaluation measures consider the
380 performance of a classifier from diverse aspects, as most algorithms learn from
381 training examples by explicitly or implicitly optimizing one specific metric.
382 To this end, in this work we also used *General_B* [31], a simple baseline
383 learning algorithm which learns by only looking at the multi-labels of the
384 dataset. As this algorithm does not necessarily concentrate on optimizing
385 specific loss functions, it can be used as a global baseline for the difficult
386 task of evaluating multi-label predictions.

387 The rationale behind *General_B* is very simple. It consists of ranking the
388 q single labels in L according to their individual relative frequencies in the
389 multi-labels in order to include the σ most frequent labels in the predicted
390 multi-label Z . To obtain a representative Z , *General_B* defines σ as the closest
391 integer value of the label cardinality LC — Section 2.5. In case of ties (single
392 labels with the same frequency), the label co-occurrence measure chooses the
393 label which maximizes its co-occurrence with better ranked labels.

394 4. Results and discussion

395 In this section, we compare the learning performance of the *BRkNN-b*
396 classifiers ($k = 10$) built from the datasets described by the features selected

397 by three groups of FS approaches:

- 398 1. The standard *IG-BR*;
- 399 2. *IG-BR* after applying the four *LCFS* settings to construct the new sets
400 of labels: *LS-X*, *CS-A*, *RS-A* and *RS-X*;
- 401 3. The reference Random Feature Selection (*RFS*).

402 The main difference between groups (1) and (2) consists in the label space
403 submitted for FS (Figure 4). The four *LCFS* settings take into account the
404 different strategies illustrated in Section 2.4.3. Finally, as is the case with
405 group (1), the method in group (3) selects features directly from the original
406 datasets.

407 As a *BRkNN-b* classifier is built from a dataset described by the best 10%
408 up to 90% of the features ranked by each FS method, 54 cases, *i.e.*, 9 feature
409 subsets \times 6 FS methods, are evaluated for each multi-label dataset and eval-
410 uation measure. All the experimental results, as well as tables and graph-
411 ical representations, can be found in the supplementary material available
412 at [http://www.labic.icmc.usp.br/pub/mcmonard/ExperimentalResults/
413 NEUCOM2015.pdf](http://www.labic.icmc.usp.br/pub/mcmonard/ExperimentalResults/NEUCOM2015.pdf). In what follows, some of the experimental results are sum-
414 marized and discussed.

415 4.1. Results overview

416 First of all, it is worth noticing that the datasets and measures in which
417 the *BRkNN-b* classifiers built using all features failed to improve on the
418 baseline classifier *General_B* — Table 6. It should be observed that this
419 situation is not unusual in multi-label learning. In [32] we carried out the
420 SR process to find papers reporting experimental evaluation measure values

421 of classifiers which were constructed using publicly available datasets, and
 422 reported on several statistics. From the 10 datasets most frequently used in
 423 the selected papers, the statistics show that 12.8% of these published results
 424 were worse than or equal to the ones obtained by $General_B$. Moreover,
 425 this percentage is unevenly distributed among the datasets. In the “worst”
 426 dataset, 43.0% of such results were reported, and in the “best” one only
 427 0.6%.

Table 6: Cases where $General_B$ outperforms the $BRkNN$ classifier built using all features

| Dataset | F -measure | Hamming loss | Accuracy | F_b |
|-----------------------|--------------|--------------|----------|-------|
| 1- <i>Cal500</i> | ✓ | ✓ | ✓ | ✓ |
| 2- <i>Corel5K</i> | ✓ | ✓ | ✓ | ✓ |
| 3- <i>Corel16k001</i> | ✓ | ✓ | ✓ | ✓ |
| 5- <i>Fapesp</i> | | ✓ | | |
| 7- <i>Llog-f*</i> | | ✓ | | |

428 Nevertheless, dataset 1-*Cal500* is the only one in which the FS methods
 429 for all feature subsets considered fails to improve on $General_B$ in the four
 430 evaluation measures used in this work. Considering the datasets individually,
 431 it can be observed that, as expected, the degree of improvement of each FS
 432 method using different feature subsets is dependent on the particular dataset.
 433 One of the datasets which obtained good results in three of the four evaluation
 434 measures using small feature subsets, is dataset 5-*Fapesp*. Recall that finding
 435 a small number of good features is an aim of the FS task. Figure 5 shows the
 436 performance of the $BRkNN$ - b classifiers according to each evaluation measure
 437 (y -axis) built using each feature subset (x -axis) in this dataset.

438 Figures 5(a), 5(c) and 5(d) show that very good results were obtained
 439 by using 10% of the features selected. Moreover, there is a considerable
 440 difference with the RFS method used as a reference, as well as the classifier

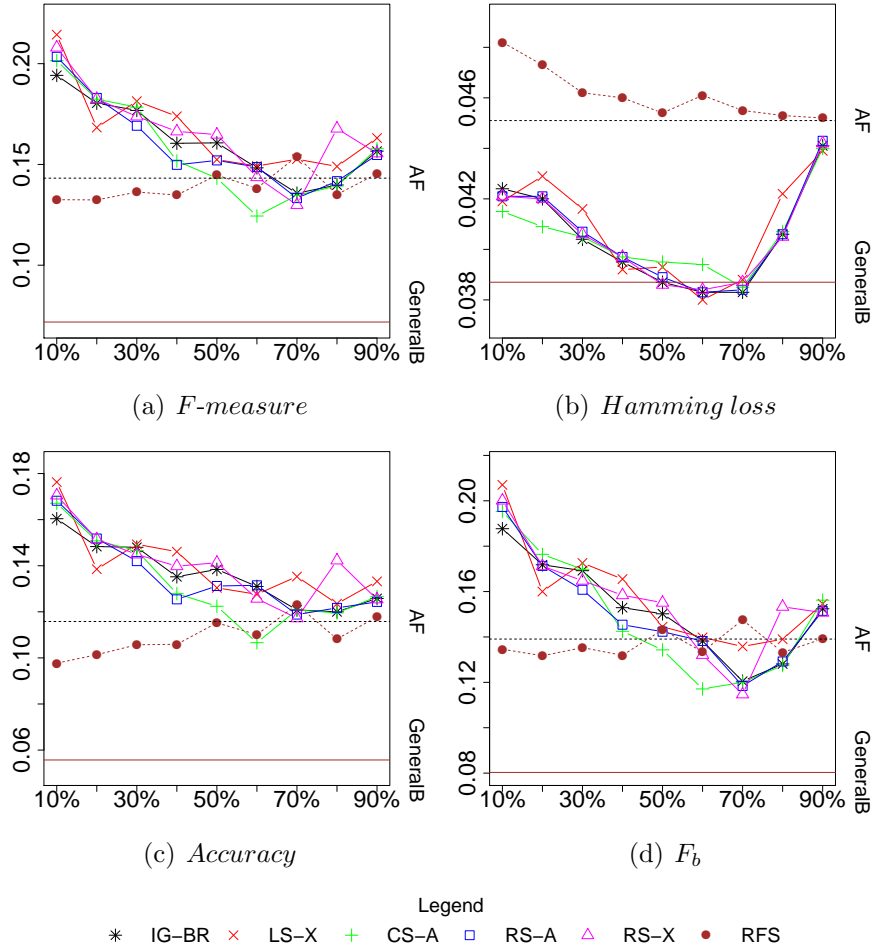


Figure 5: Performance of FS methods assessed by the correspondent $BRkNN-b$ classifier in dataset *5-Fapesp*

441 built using all features (AF). For *F-measure*, *Accuracy* and F_b , the best
442 results were obtained by the FS methods *LS-X*, followed by *RS-X*, *RS-A*,
443 *CS-A* and *IG-BR* in that order.

444 *Hamming loss* results better than the ones obtained by the baseline
445 classifier $General_B$ were achieved only for three feature subsets: $50\%M$,
446 $60\%M$ and $70\%M$ — Figure 5(b). Nevertheless, the heuristic FS methods
447 are still notably better than *RFS* and AF. Thus, regardless of the evaluation
448 measure, the *IG-BR* and *LCFS* methods are highlighted in this dataset.

449 In fact, *Hamming loss* is the evaluation measure which more often
450 showed worse results than the ones obtained by the baseline classifier $General_B$.
451 Figure 6 shows the number of datasets in which a *BRkNN-b* classifier, built
452 using a feature subset chosen by a FS method, achieved *Hamming loss* val-
453 ues worse than the ones obtained by $General_B$. In this figure, the horizontal
454 thick line shows the number of datasets in which a classifier built using All
455 Features (AF), *i.e.*, without feature selection, was worse than $General_B$.

456 Figure 6 shows that for small feature subsets ($|X'| \leq 50\%M$), three *LCFS*
457 variations (*CS-A*, *RS-A* and *RS-X*) were able to reduce by one the number
458 of datasets in which the *BRkNN-b* classifiers were worse than the baseline
459 classifier $General_B$. In particular, *RS-X* is the one which obtained that result
460 with the smallest feature subsets ($|X'| = 10\%M$). On the other hand, the
461 original *IG-BR* only achieved that result when using larger feature subsets
462 ($|X'| = 60\%M$ and $|X'| = 70\%M$).

463 4.2. Statistical comparison among the FS methods

464 To assess whether the overall differences in performance across the multi-
465 label FS methods are statistically significant, we used the Friedman’s test

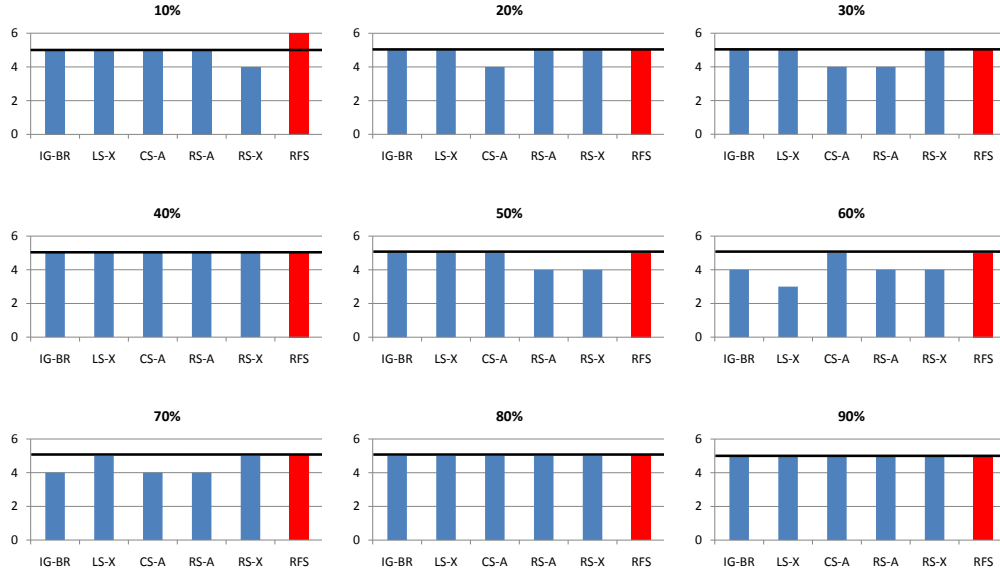


Figure 6: Number of datasets in which a $BRkNN-b$ classifier built using a feature subset chosen by a FS method achieved *Hamming loss* values worse than the ones obtained by the *baseline General_B*

466 and the Nemenyi’s post-hoc test as recommended by Demšar [33]. The Fried-
 467 man’s test is a non-parametric test for multiple hypotheses testing. It ranks
 468 the methods according to their performance for each dataset separately. The
 469 best performing method obtains the rank of 1, the second best the rank of 2,
 470 and so on. In case of ties, it assigns average ranks. If a statistically significant
 471 difference in the performance is detected, the next step is a post-hoc test to
 472 detect between which methods those differences appear.

473 We applied the Friedman’s statistical test under the null hypothesis that
 474 the performances of the classifiers built using the features selected by each
 475 FS method are equivalent. The statistical results can be found in the suple-
 476 mentary material. As the hypothesis was rejected at the significance level
 477 $\alpha = 0.05$ for all measures, we proceed with the Nemenyi’s multiple comparison

478 post-hoc test to detect which differences among the methods are significant.
 479 This post-hoc test points out a significant difference whenever the average
 480 rank of two methods differ by more than a Critical Difference (CD). Fur-
 481 thermore, the results of the post-hoc test can be visually represented with
 482 a simple diagram. Figure 7 shows the correspondent diagrams, on the four
 483 evaluation measures considered in this work, for the smaller feature subsets
 484 evaluated — $|X'| = 10\%M, 20\%M, 30\%M$. The diagrams for other feature
 485 subsets evaluated can be found in the supplementary material. The lines
 486 for the average ranks of the methods that do not differ significantly (at the
 487 significance level of $\alpha = 0.05$) are connected with a line. Observe that for
 488 the *Hamming loss* diagrams, the higher the average ranking, the better the
 489 FS method is, whereas for the remaining diagrams, the lower the average
 490 ranking, the better the method is.

491 In all cases in which a significant difference was found, at least a heuristic
 492 FS method outperformed *RFS*. Moreover, in all cases *RFS* was always ranked
 493 last. Table 7 summarizes all algorithms significantly better than *RFS* for each
 494 evaluation measure and feature subset size. For example, the first entry in
 495 this table shows that *RS-A*, *LS-X* and *RS-X* are significantly better than
 496 *RFS* when *F-measure* is considered.

Table 7: Multi-label FS methods significantly better than *RFS* ($\alpha = 0.05$)

| $ X' $ | <i>F-measure</i> | <i>Hamming loss</i> | <i>Accuracy</i> | F_b |
|--------|---|--|---|--|
| 10% | <i>RS-A</i> , <i>LS-X</i> , <i>RS-X</i> | <i>RS-A</i> , <i>LS-X</i> , <i>CS-A</i> | <i>RS-A</i> , <i>LS-X</i> , <i>RS-X</i> | <i>RS-A</i> , <i>RS-X</i> , <i>LS-X</i> |
| 20% | <i>RS-X</i> | <i>RS-X</i> , <i>RS-A</i> , <i>CS-A</i> , <i>IG-BR</i> | <i>RS-X</i> | <i>RS-X</i> , <i>CS-A</i> , <i>IG-BR</i> , <i>RS-A</i> |
| 30% | <i>RS-X</i> | <i>RS-X</i> , <i>CS-A</i> | <i>RS-X</i> | <i>RS-X</i> , <i>CS-A</i> |
| 70% | | <i>LS-X</i> , <i>RS-X</i> , <i>IG-BR</i> | | |
| 80% | <i>RS-X</i> | <i>RS-X</i> | <i>RS-X</i> | <i>RS-X</i> |

497 Regardless of the evaluation measure, from the 17 out of a total of 36

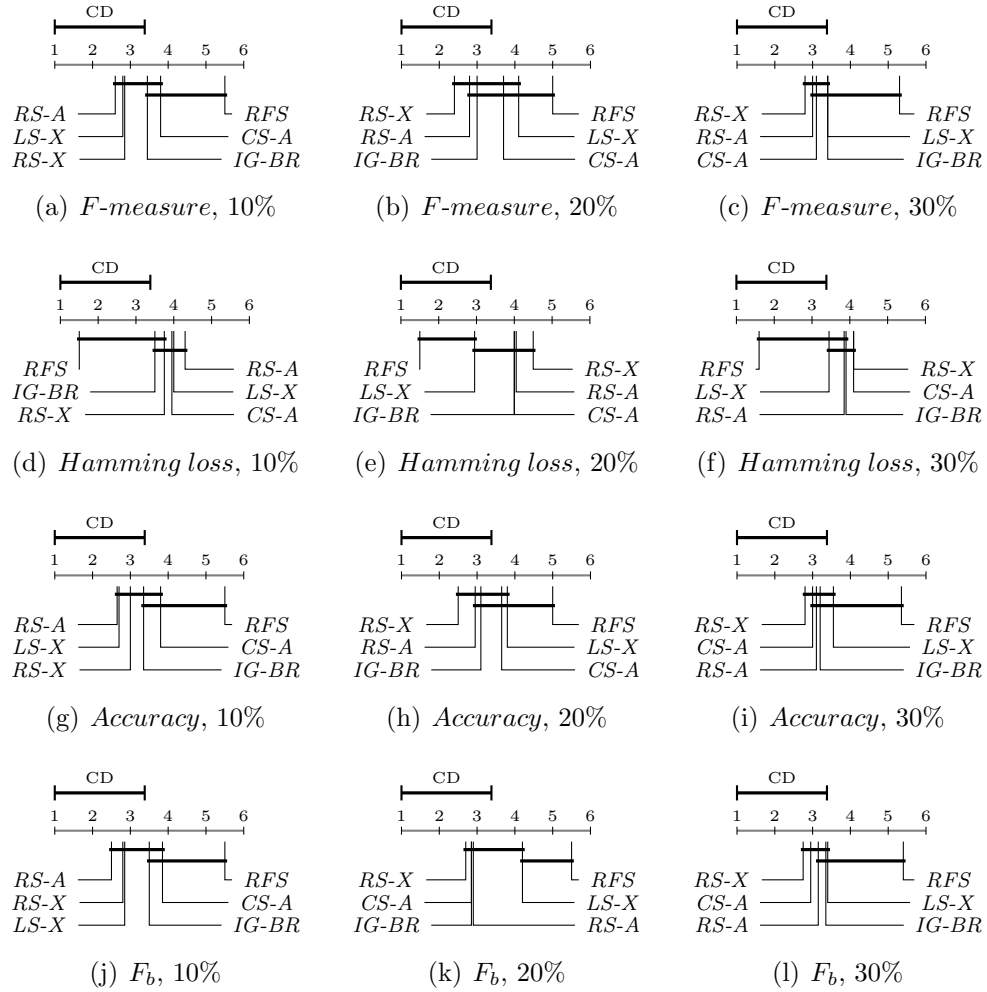


Figure 7: Nemenyi's test comparison of the performance measure values achieved by $BRkNN$ - b classifiers built after selecting the best $|X'| = 10\%M, 20\%M, 30\%M$ of the features ranked by each FS method. Groups of classifiers that are not significantly different according to the Critical Difference (CD) — at $\alpha = 0.05$ — are connected

498 cases analyzed in which significant differences were found, it can be observed
 499 that the *LCFS* setting *RS-X* was the one highlighted more often as being
 500 significantly better than *RFS* (15 times out of 17), followed by *RS-A*, *LS-X*
 501 and *CS-A* (5 times out of 17 each), and finally *IG-BR* (3 times out of 17).
 502 As expected, significant differences tend to diminish when the feature subset
 503 size is large.

504 The Friedman’s test also provides information about the best method
 505 built after FS by the rankings averaged across all datasets — Table 8. In
 506 this table, each symbol identifies a FS method: $-$ (*IG-BR*), $*$ (*LS-X*), o (*CS-*
 507 *A*), \times (*RS-A*), $+$ (*RS-X*) and \bullet (*RFS*). The last rows and columns sum up
 508 the results for each method. Note that there is more than one FS method in
 509 some cells, as the average rankings achieved by the correspondent classifiers
 510 are equal.

Table 8: Best FS method based on the Friedman’s test average rankings calculated for each feature subset size and evaluation measure

| | 10%M | 20%M | 30%M | 40%M | 50%M | 60%M | 70%M | 80%M | 90%M | - | * | o | × | + | • |
|----------------------|------|------|---------------|---------------|---------------|---------------|------|------|------|---|----|---|---|----|---|
| <i>F-measure</i> | × | + | + | * | $\frac{o}{x}$ | * | * | + | + | 0 | 3 | 1 | 2 | 4 | 0 |
| <i>Hamming loss</i> | × | + | $\frac{o}{+}$ | * | × | $\frac{-}{+}$ | * | + | * | 1 | 3 | 1 | 3 | 4 | 0 |
| <i>Accuracy</i> | × | + | + | $\frac{o}{+}$ | $\frac{o}{+}$ | * | * | + | + | 0 | 3 | 2 | 1 | 5 | 0 |
| <i>F_b</i> | × | + | + | * | × | $\frac{-}{*}$ | * | + | * | 1 | 4 | 0 | 2 | 3 | 0 |
| - | 0 | 0 | 0 | 0 | 0 | 2 | 0 | 0 | 0 | 2 | | | | | |
| * | 0 | 0 | 0 | 4 | 0 | 3 | 4 | 0 | 2 | | 13 | | | | |
| o | 0 | 0 | 1 | 1 | 2 | 0 | 0 | 0 | 0 | | | 4 | | | |
| × | 4 | 0 | 0 | 1 | 3 | 0 | 0 | 0 | 0 | | | | 8 | | |
| + | 0 | 4 | 4 | 1 | 0 | 1 | 0 | 4 | 2 | | | | | 16 | |
| • | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | | | | | | 0 |

511 Regardless of the evaluation measures, *RS-X* (+) achieved the best aver-
 512 age rankings more often, mainly when $|X'| < \frac{|X|}{2}$, *i.e.*, less than half of the
 513 features are used. This setting was already highlighted in Figure 7. In addi-
 514 tion, another *LCFS* setting, *RS-A* (×), obtained the best average ranking for

515 the smallest feature subsets ($|X'| = 10\%M$). *IG-BR* (-), in turn, obtained
516 the best average ranking in only two cases, when a large number of features
517 ($|X'| = 60\%M$) was selected. On the other hand, no classifier built using the
518 features chosen by *RFS* achieved the best average ranking.

519 As the *LCFS* setting *RS-X* was prominent in the statistical comparison
520 and *IG-BR* represents the standard approach, we focused on the comparison
521 of both methods.

522 We applied the Wilcoxon signed-ranks test, recommended for compar-
523 isons of two algorithms [33], with the null hypothesis that both methods are
524 equivalent and $\alpha = 0.05$. By applying this test for each evaluation measure,
525 *RS-X* was significantly better than *IG-BR* twice — when the classifiers are
526 built using the feature subset size $|X'| = 80\%M$ and evaluated by *F-measure*
527 and *Accuracy*.

528 4.3. LCFS RS-X versus AF

529 We compared the performance of the *BRkNN-b* classifiers built using
530 the features selected by *IG-BR* in the original datasets and in the datasets
531 augmented by using the four *LCFS* settings, as well as the classifiers built
532 using features randomly chosen by *RFS*. In this comparison, *RS-X* showed
533 good results when fewer features were selected. However, the quality of the
534 classifiers has not been taken into account. To this end, we compare the
535 performance of the classifiers built by *BRkNN-b*, using up to 30% of the
536 features selected by *RS-X*, with the performance achieved by the *BRkNN-b*
537 classifiers using all features, *i.e.*, the original dataset. Table 9 shows, for each
538 dataset, and for each one of the four evaluation measures used in this work,
539 whenever the classifiers built using the features selected by *RS-X* achieved

540 evaluation measure values better than or equal to (indicated by \star), or at
 541 most 5% worse (indicated by \star) than the ones obtained by the classifiers
 542 using all features. The symbol 0 indicates the other cases.

Table 9: Classifiers built using the features selected by *RS-X* vs the classifiers built using all features

| Dataset | $ X' = 10\%M$ | $ X' = 20\%M$ | $ X' = 30\%M$ |
|-----------------------|---------------------------|---------------------------|---------------------------|
| 1- <i>Cal500</i> | $\star/\star/\star/\star$ | $\star/\star/\star/\star$ | $\star/\star/\star/\star$ |
| 2- <i>Corel5k</i> | $\star/\star/\star/\star$ | $\star/\star/\star/\star$ | $\star/\star/\star/\star$ |
| 3- <i>Corel16k001</i> | $\star/\star/\star/\star$ | $\star/\star/\star/\star$ | $\star/\star/\star/\star$ |
| 4- <i>Emotions</i> | 0 / 0 / 0 / 0 | $\star/0/\star/\star$ | $\star/0/\star/\star$ |
| 5- <i>Fapesp</i> | $\star/\star/\star/\star$ | $\star/\star/\star/\star$ | $\star/\star/\star/\star$ |
| 6- <i>Genbase*</i> | $\star/\star/\star/\star$ | $\star/\star/\star/\star$ | $\star/\star/\star/\star$ |
| 7- <i>Llog-f*</i> | $\star/\star/\star/\star$ | $\star/\star/\star/\star$ | $\star/\star/\star/\star$ |
| 8- <i>Magtag5k</i> | 0 / 0 / 0 / 0 | 0 / 0 / 0 / 0 | 0 / 0 / 0 / 0 |
| 9- <i>Scene</i> | 0 / 0 / 0 / 0 | 0 / 0 / 0 / 0 | 0 / 0 / 0 / 0 |
| 10- <i>Yeast</i> | 0 / 0 / 0 / 0 | $\star/0/\star/\star$ | $\star/\star/\star/\star$ |

543 As can be observed, very good results were obtained in 5 out of the 10
 544 datasets, in which the four evaluation measure values of the classifiers based
 545 on the *RS-X* setting were better than or equal to the ones obtained by the
 546 AF classifiers. Good results were obtained in all cases in dataset 1-*Cal500*,
 547 as well as in some cases in datasets 4-*Emotion* and 10-*Yeast* when 20% and
 548 30% of the features are selected, as the results are at most 5% worse than
 549 the AF ones. On the other hand, poor results were obtained in datasets
 550 8-*Magtag5k* and 9-*Scene* even when 30% of the features are considered. In
 551 fact, it is necessary to consider 70% of the features selected in these datasets
 552 in order to obtain ($\star/\star/\star/\star$).

553 However, the good results obtained in dataset 1-*Cal500* should be con-
 554 sidered with care, as the classifiers built using AF are worse than the ones
 555 built by the baseline classifier *General_B*. In fact, this seems to be a difficult
 556 dataset to learn from. Table 4 shows that this dataset has $N = 502$ examples,
 557 as well as the same number of different multi-labels ($\#Diff$), in which the

558 average number of labels associated with each example is $LC = 20.044$ from
 559 a total of $|L| = 174$ labels.

560 5. Related work found by the systematic review method

561 Feature selection has been an active research topic in supervised learning,
 562 and there are many related publications and comprehensive surveys [7]. Al-
 563 though most FS publications are related to single-label learning, a number of
 564 papers have recently reported results to support multi-label learning [18, 8].

565 Aiming at capturing a wide, replicable and rigorous overview of the topic,
 566 we have instantiated the systematic literature review method [17] for multi-
 567 label FS in [24] and updated it in March/2015. Table 10 summarizes the 99
 568 publications found in terms of two categorizations described in Sections 2.1
 569 and 2.2: the order of label dependence and the interaction of the FS method
 570 with the learning algorithm. More information regarding the 99 references
 571 is described in the supplementary material available at [http://www.labic.
 572 icmc.usp.br/pub/mcmonard/ExperimentalResults/NEUCOM2015.pdf](http://www.labic.icmc.usp.br/pub/mcmonard/ExperimentalResults/NEUCOM2015.pdf).

Table 10: Number of papers published per approach found by the systematic literature review process (total = 99 related publications)

| categorization | approach | #publications |
|---|--------------|---------------|
| order of label dependence | first-order | 53 |
| | second-order | 15 |
| | high-order | 14 |
| | hybrid | 8 |
| | unrecognized | 9 |
| interaction with the learning algorithm | filter | 70 |
| | embedded | 12 |
| | wrapper | 8 |
| | hybrid | 3 |
| | unrecognized | 6 |

573 As can be observed, filters and first-order strategies have been the most
 574 usual choices in multi-label FS. This behavior could be partly explained by

575 the relative lower computational cost in comparison with other alternatives.
576 In addition, these strategies can be combined, as exemplified by *IG-BR* and
577 some proposals from the related work [34, 35, 36, 5, 37]. In particular, these
578 filter methods apply the Information Gain importance measure in binary
579 data directly or indirectly transformed by the *BR* approach, a first-order
580 strategy.

581 Regarding importance measures, Information Gain has been the most
582 often used measure (23 out of 99 papers). Mutual information [38], chi-
583 squared [39], ReliefF [40] and correlation-based feature selection [1] come
584 next.

585 The method we present, *LCFS*, pioneers label construction as a second-
586 order strategy. Other methods that take into account label relations have
587 also been proposed, reporting good results [41, 42, 43, 44, 45, 46, 47, 48, 49].
588 By organizing the related work according to the order of label dependence
589 exploration, the SR can be useful for further research on multi-label feature
590 selection.

591 Although the number of papers with unrecognized and hybrid¹⁰ strategies
592 in Table 10 is relatively low, it indicates the need to consider a taxonomy spe-
593 cific for multi-label FS. Well established taxonomies for single-label feature
594 selection would be useful as a starting point [7].

595 The SR also provides information regarding the number of papers per
596 publication year. Figure 8 suggests that interest in multi-label feature selec-
597 tion is increasing as time goes by.

¹⁰In this work, a hybrid strategy is considered whenever the FS method falls into two or more categories.

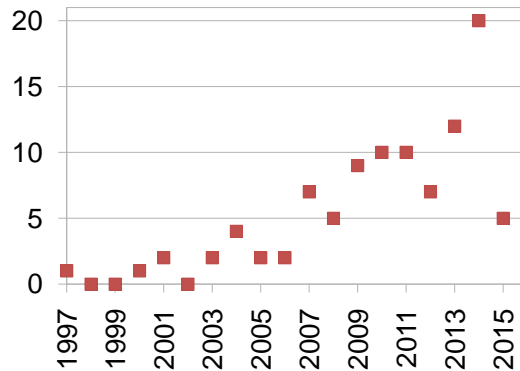


Figure 8: Number of papers published per year found by the systematic literature review process (total = 99 related publications)

598 **6. Conclusion**

599 This work presented and evaluated *LCFS*, a method that constructs new
600 labels based on relations between the original dataset labels. The new labels
601 are included in the dataset before applying the standard multi-label feature
602 selection approach based on binary relevance. By doing so, *LCFS* considers
603 second-order label information during filter feature selection.

604 The experimental evaluation on 10 benchmark multi-label datasets shows
605 that the *LCFS* setting *RS-X* gave rise to classifiers similar to, or better
606 than, the ones built by simply combining *BR* and Information Gain — *IG-*
607 *BR*. *RS-X* is a simple alternative that randomly selects a pair of labels and
608 combines them by the XOR or XNOR operator to yield each new label.
609 Thus, the *LCFS* setting is competitive with *IG-BR* by slightly increasing the
610 computational cost due to the application of a binary operator. Moreover,
611 the evaluated method contributed to outperform classifiers built using all
612 features, *i.e.*, without FS, as well as the baseline classifier *General_B* and
613 random feature selection.

614 As an additional contribution, this work pioneers the use of the system-
615 atic review method to survey the related work on multi-label FS. We or-
616 ganized the 99 papers found in terms of two categorizations proposed for
617 multi-label learning methods and single-label FS algorithms. By doing so,
618 it was observed that most of them consider first-order strategies, *i.e.*, ignore
619 label dependence, and follow the filter approach. In the summary of the
620 99 papers, it was also found evidence that agrees with *LCFS* experimental
621 achievements, as some related papers reported good results when exploring
622 label dependence.

623 As future work, we plan to use other multi-label learning algorithms to
624 evaluate FS, as exemplified in previous work [47], aiming to reduce the po-
625 tential influence of a specific algorithm. Furthermore, we plan to evaluate
626 *LCFS* strategies based on label weighting [50] in benchmark and synthetic
627 datasets [51]. By applying Exploratory Data Analysis [52] in these cases, we
628 expect to find relations among the quality of filter FS methods and multi-
629 label datasets properties.

630 **Acknowledgements**

631 This research was supported by the São Paulo Research Foundation
632 (FAPESP), grant 2011/02393-4. This agency did not have any further in-
633 volvement in this paper.

634 **Vitae**



Newton Spolaôr holds a Ph.D. degree in Computer Science at the Institute of Mathematics and Computer Science from the University of São Paulo, Brazil (2014). He also holds an M.Sc. degree in Information Engineering from the Federal University of ABC (2010) and a B.Sc. degree in Computer Science from the State West Paraná University (2008). His main research interests are feature selection, multi-label learning, systematic review, time series clustering and analysis of biomedical images.



Maria Carolina Monard is Emeritus Professor in Computer Science at the Institute of Mathematics and Computer Science at the University of São Paulo, Brazil. She holds a Ph.D. degree from the Pontifical Catholic University of Rio de Janeiro, Brazil (1980), and an M.Sc. degree from Southampton University, UK (1968). Her research interests are in the field of Artificial Intelligence, more specifically in Machine Learning, Data and Text Mining.



Grigorios Tsoumakias is Assistant Professor at the Department of Informatics at the Aristotle University of Thessaloniki (AUTH), Greece. He holds a Ph.D. degree in Informatics from AUTH (2005), an M.Sc. in Artificial Intelligence from the University of Edinburgh (2000) and a degree in Informatics from AUTH (1999). His research interests include various aspects of machine learning, knowledge discovery and data mining, including ensemble methods, distributed data mining, text classification and multi-label learning.



Huei Diana Lee is Assistant Professor at the Engineering Institute and Exact Sciences at the State West Paraná University, Brazil. She holds a Ph.D. and an M.Sc. degree in Computer Science from the Institute of Mathematics and Computer Science at the University of São Paulo (2005 and 2000). She also holds a B.Sc. degree in Computer Science from the São Paulo State University, Brazil (1994). Her research interests are in bioinformatics, intelligent data analysis and telemedicine.

635 **References**

- 636 [1] S. Jungjit, A. A. Freitas, M. Michaelis, J. Cinatl, Extending multi-
637 label feature selection with kegg pathway information for microarray
638 data analysis, in: IEEE Conference on Computational Intelligence in
639 Bioinformatics and Computational Biology, 2014, pp. 1–8.
- 640 [2] D. Heider, R. Senge, W. Cheng, E. Hüllermeier, Multilabel classifica-
641 tion for exploiting cross-resistance information in HIV-1 drug resistance
642 prediction, *Bioinformatics* 29 (2013) 1946–1952.
- 643 [3] P. K. Bhowmick, A. Basu, P. Mitra, A. Prasad, Multi-label text clas-
644 sification approach for sentence level news emotion analysis, in: Inter-

- 645 national Conference on Pattern Recognition and Machine Intelligence,
646 2009, pp. 261–266.
- 647 [4] A. Esuli, T. Fagni, F. Sebastiani, Boosting multi-label hierarchical text
648 categorization, *Information Retrieval* 11 (2008) 287–313.
- 649 [5] W. Chen, J. Yan, B. Zhang, Z. Chen, Q. Yang, Document transforma-
650 tion for multi-label feature selection in text categorization, in: *IEEE*
651 *International Conference on Data Mining*, 2007, pp. 451–456.
- 652 [6] M. R. Boutell, J. Luo, X. Shen, C. M. Brown, Learning multi-label
653 scene classification, *Pattern Recognition* 37 (2004) 1757–1771.
- 654 [7] H. Liu, H. Motoda, *Computational Methods of Feature Selection*, Chap-
655 man & Hall/CRC, 2008.
- 656 [8] G. Tsoumakas, I. Katakis, I. P. Vlahavas, Mining multi-label data, in:
657 O. Maimon, L. Rokach (Eds.), *Data Mining and Knowledge Discovery*
658 *Handbook*, Springer, 2010, pp. 667–685.
- 659 [9] H. Motoda, H. Liu, Feature selection, extraction and construction,
660 *Communication of Institute of Information and Computing Machinery*
661 5 (2002) 67–72.
- 662 [10] K. Lillywhite, D.-J. Lee, B. Tippetts, J. Archibald, A feature con-
663 struction method for general object recognition, *Pattern Recognition* 46
664 (2013) 3300–3314.
- 665 [11] D. García, A. González, R. Pérez, A feature construction approach

- 666 for genetic iterative rule learning algorithm, *Journal of Computer and*
667 *System Sciences* 80 (2014) 101–117.
- 668 [12] S. Piramuthu, R. Sikora, Iterative feature construction for improving
669 inductive learning algorithms, *Expert Systems with Applications* 36
670 (2009) 3401–3406.
- 671 [13] R. Prati, F. Olivetti de Franca, Extending features for multilabel clas-
672 sification with swarm biclustering, in: *IEEE Congress on Evolutionary*
673 *Computation*, 2013, pp. 2964–2971.
- 674 [14] W. Duivesteijn, E. Loza Mencía, J. Fürnkranz, A. Knobbe, Multi-
675 label LeGo - enhancing multi-label classifiers with local patterns, in:
676 J. Hollmén, F. Klawonn, A. Tucker (Eds.), *Advances in Intelligent*
677 *Data Analysis XI*, volume 7619 of *Lecture Notes in Computer Science*,
678 Springer, 2012, pp. 114–125.
- 679 [15] Y. Yang, S. Gopal, Multilabel classification with meta-level features in
680 a learning-to-rank framework, *Machine Learning* 88 (2012) 47–68.
- 681 [16] N. Spolaôr, M. C. Monard, G. Tsoumakas, H. D. Lee, Label construction
682 for multi-label feature selection, in: *Brazilian Conference on Intelligent*
683 *Systems*, IEEE, 2014, pp. 1–6.
- 684 [17] B. A. Kitchenham, S. Charters, Guidelines for performing systematic
685 literature reviews in software engineering, EBSE-2007-01 Technical Re-
686 port. 65 pg., 2007. Evidence-based Software Engineering.
- 687 [18] M.-L. Zhang, Z.-H. Zhou, A review on multi-label learning algorithms,

- 688 IEEE Transactions on Knowledge and Data Engineering 26 (2014) 1819–
689 1837.
- 690 [19] E. Spyromitros, G. Tsoumakas, I. Vlahavas, An empirical study of lazy
691 multilabel classification algorithms, in: Hellenic conference on Artificial
692 Intelligence, Springer-Verlag, 2008, pp. 401–406.
- 693 [20] M.-L. Zhang, Z.-H. Zhou, A k-nearest neighbor based algorithm for
694 multi-label classification, IEEE International Conference on Granular
695 Computing 2 (2005) 718–721.
- 696 [21] K. Dembczynski, W. Waegeman, W. Cheng, E. Hüllermeier, On label
697 dependence and loss minimization in multi-label classification, Machine
698 Learning 88 (2012) 5–45.
- 699 [22] N. Spolaôr, E. A. Cherman, M. C. Monard, H. D. Lee, A comparison of
700 multi-label feature selection methods using the problem transformation
701 approach, Electronic Notes in Theoretical Computer Science 292 (2013)
702 135–151.
- 703 [23] B. Kitchenham, R. Pretorius, D. Budgen, O. P. Brereton, M. Turner,
704 M. Niazi, S. Linkman, Systematic literature reviews in software en-
705 gineering - a tertiary study, Information and Software Technology 52
706 (2010) 792–805.
- 707 [24] N. Spolaôr, M. C. Monard, H. D. Lee, A systematic review to identify
708 feature selection publications in multi-labeled data, ICMC Technical
709 Report No 374. 31 pg., 2012. University of São Paulo.

- 710 [25] G. Tsoumakas, E. Spyromitros-Xioufis, J. Vilcek, I. Vlahavas, Mulan:
711 A java library for multi-label learning, *Journal of Machine Learning*
712 *Research* 12 (2011) 2411–2414.
- 713 [26] I. H. Witten, E. Frank, *Data Mining: Practical Machine Learning Tools*
714 *and Techniques*, Morgan Kaufmann, 2011.
- 715 [27] R. G. Rossi, S. O. Rezende, Building a topic hierarchy using the bag-of-
716 related-words representation, in: *Symposium on Document Engineering*,
717 2011, pp. 195–204.
- 718 [28] G. Marques, M. A. Domingues, T. Langlois, F. Gouyon, Three current
719 issues in music autotagging, in: *Conference of the International Society*
720 *for Music Information Retrieval*, 2011, pp. 795–800.
- 721 [29] M. Robnik-Šikonja, I. Kononenko, Theoretical and empirical analysis of
722 ReliefF and RReliefF, *Machine Learning* 53 (2003) 23–69.
- 723 [30] N. Spolaôr, G. Tsoumakas, Evaluating feature selection methods for
724 multi-label text classification, in: *BioASQ workshop*, 2013, pp. 1–12.
- 725 [31] J. Metz, L. F. Abreu, E. A. Cherman, M. C. Monard, On the esti-
726 mation of predictive evaluation measure baselines for multi-label learn-
727 ing, in: J. Pavón, N. D. Duque-Méndez, R. Fuentes-Fernández (Eds.),
728 *Advances in Artificial Intelligence - IBERAMIA 2012*, volume 7637 of
729 *Lecture Notes in Computer Science*, Springer, 2012, pp. 189–198.
- 730 [32] J. Metz, N. Spolaôr, E. A. Cherman, M. C. Monard, Comparing pub-
731 lished multi-label classifier performance measures to the ones obtained

- 732 by a simple multi-label baseline classifier, in: Arxiv.org, Unpublished
733 results, pp. 1–19. URL: <http://arxiv.org/abs/1503.06952>.
- 734 [33] J. Demšar, Statistical comparison of classifiers over multiple data sets,
735 *Journal of Machine Learning Research* 7 (2006) 1–30.
- 736 [34] M. Du, M. Pierce, L. Pivovarova, R. Yangarber, Supervised classifica-
737 tion using balanced training, in: L. Besacier, A.-H. Dediu, C. Martín-
738 Vide (Eds.), *Statistical Language and Speech Processing*, volume 8791 of
739 *Lecture Notes in Computer Science*, Springer International Publishing,
740 2014, pp. 147–158.
- 741 [35] M. Karabulut, Fuzzy unordered rule induction algorithm in text catego-
742 rization on top of geometric particle swarm optimization term selection,
743 *Knowledge-Based Systems* 54 (2013) 288–297.
- 744 [36] G.-P. Liu, J.-J. Yan, Y.-Q. Wang, J.-J. Fu, Z.-X. Xu, R. Guo, P. Qian,
745 Application of multilabel learning using the relevant feature for each
746 label in chronic syndrome diagnosis, *Evidence-Based Complementary
747 and Alternative Medicine* 2012 (2012) 1–9.
- 748 [37] Z. Zheng, X. Wu, R. Srihari, Feature selection for text categorization on
749 imbalanced data, *Special Interest Group on Knowledge Discovery and
750 Data Mining Explorations Newsletter* 6 (2004) 80–89.
- 751 [38] Y. Yang, J. O. Pedersen, A comparative study on feature selection in
752 text categorization, in: *Proceedings of the Fourteenth International
753 Conference on Machine Learning*, 1997, pp. 412–420.

- 754 [39] E. Spyromitros-Xioufis, K. Sechidis, G. Tsoumakas, I. Vlahavas, Mlkd’s
755 participation at the clef 2011 photo annotation and concept-based re-
756 trieval tasks, in: ImageCLEF Lab of CLEF 2011 Conference on Multi-
757 lingual and Multimodal Information Access Evaluation, 2011, pp. 1–15.
- 758 [40] D. Kong, C. Ding, H. Huang, H. Zhao, Multi-label ReliefF and F-
759 statistic feature selections for image annotation, in: IEEE Conference
760 on Computer Vision and Pattern Recognition, 2012, pp. 2352–2359.
- 761 [41] J. Lee, D.-W. Kim, Memetic feature selection algorithm for multi-label
762 classification, *Information Sciences* 293 (2015) 80–96.
- 763 [42] J. Lee, D.-W. Kim, Mutual information-based multi-label feature selec-
764 tion using interaction information, *Expert Systems with Applications*
765 42 (2015) 2013–2025.
- 766 [43] O. G. R. Pupo, C. Morell, S. V. Soto, Scalable extensions of the relieff
767 algorithm for weighting and selecting features on the multi-label learning
768 context (in press), *Neurocomputing* (2015).
- 769 [44] F. M. Rodrigues, C. J. Camara, A. M. P. Canuto, A. M. Santos, Confi-
770 dence factor and feature selection for semi-supervised multi-label classifi-
771 cation methods, in: International Joint Conference on Neural Networks,
772 2014, pp. 864–871.
- 773 [45] K. Sechidis, N. Nikolaou, G. Brown, Information theoretic feature se-
774 lection in multi-label data through composite likelihood, in: P. Fränti,
775 G. Brown, M. Loog, F. Escolano, M. Pelillo (Eds.), *Structural, Syntac-*

- 776 tic, and Statistical Pattern Recognition, volume 8621 of *Lecture Notes*
777 *in Computer Science*, Springer Berlin Heidelberg, 2014, pp. 143–152.
- 778 [46] I. Slavkov, J. Karcheska, D. Kocev, S. Kalajdziski, S. Džeroski, Reli-
779 eff for hierarchical multi-label classification, in: A. Appice, M. Ceci,
780 C. Loglisci, G. Manco, E. Masciari, Z. W. Ras (Eds.), *New Frontiers in*
781 *Mining Complex Patterns*, volume 8399 of *Lecture Notes in Computer*
782 *Science*, Springer International Publishing, 2014, pp. 148–161.
- 783 [47] N. Spolaôr, M. C. Monard, Evaluating ReliefF-based multi-label feature
784 selection algorithm, in: A. L. C. Bazzan, K. Pichara (Eds.), *Advances in*
785 *Artificial Intelligence – IBERAMIA 2014*, volume 8864 of *Lecture Notes*
786 *in Computer Science*, Springer International Publishing, 2014, pp. 194–
787 205.
- 788 [48] Y. Yu, Y. Wang, Feature selection for multi-label learning using mutual
789 information and ga, in: D. Miao, W. Pedrycz, D. Ślęzak, G. Peters,
790 Q. Hu, R. Wang (Eds.), *Rough Sets and Knowledge Technology*, vol-
791 ume 8818 of *Lecture Notes in Computer Science*, Springer International
792 Publishing, 2014, pp. 454–463.
- 793 [49] N. Spolaôr, E. A. Cherman, M. C. Monard, H. D. Lee, ReliefF for multi-
794 label feature selection, in: *Brazilian Conference on Intelligent Systems*,
795 2013, pp. 6–11.
- 796 [50] S. Jungjit, A. A. Freitas, M. Michaelis, J. Cinatl, Two extensions to
797 multi-label correlation-based feature selection: A case study in bioin-

798 formatics, in: IEEE International Conference on Systems, Man, and
799 Cybernetics, 2013, pp. 1519–1524.

800 [51] J. T. Tomás, N. Spolaôr, E. A. Cherman, M. C. Monard, A framework to
801 generate synthetic multi-label datasets, *Electronic Notes in Theoretical*
802 *Computer Science* 302 (2014) 155–176.

803 [52] G. J. Myatt, W. P. Johnson, *Making Sense of Data I - A Practical*
804 *Guide to Exploratory Data Analysis and Data Mining*, John Wiley &
805 Sons, 2014.