

Article

# How Is a Data-Driven Approach Better than Random Choice in Label Space Division for Multi-Label Classification?

Piotr Szymański <sup>1,2,\*</sup>, Tomasz Kajdanowicz <sup>1</sup> and Kristian Kersting <sup>3</sup>

<sup>1</sup> Department of Computational Intelligence, Wrocław University of Technology, Wybrzeże Stanisława Wyspiańskiego 27, 50-370 Wrocław, Poland; tomasz.kajdanowicz@pwr.edu.pl

<sup>2</sup> Illimites Foundation, Gajowicka 64 lok. 1, 53-422 Wrocław, Poland

<sup>3</sup> Department of Computer Science, TU Dortmund University, August-Schmidt-Straße 4, 44221 Dortmund, Germany; kristian.kersting@cs.tu-dortmund.de

\* Correspondence: piotr.szymanski@pwr.edu.pl; Tel.: +48-71-320-3609

Academic Editor: Andreas Holzinger

Received: 1 February 2016 ; Accepted: 19 July 2016; Published: 30 July 2016

**Abstract:** We propose using five data-driven community detection approaches from social networks to partition the label space in the task of multi-label classification as an alternative to random partitioning into equal subsets as performed by *RAkELd*. We evaluate modularity-maximizing using fast greedy and leading eigenvector approximations, infomap, walktrap and label propagation algorithms. For this purpose, we propose to construct a label co-occurrence graph (both weighted and unweighted versions) based on training data and perform community detection to partition the label set. Then, each partition constitutes a label space for separate multi-label classification sub-problems. As a result, we obtain an ensemble of multi-label classifiers that jointly covers the whole label space. Based on the binary relevance and label powerset classification methods, we compare community detection methods to label space divisions against random baselines on 12 benchmark datasets over five evaluation measures. We discover that data-driven approaches are more efficient and more likely to outperform *RAkELd* than binary relevance or label powerset is, in every evaluated measure. For all measures, apart from Hamming loss, data-driven approaches are significantly better than *RAkELd* ( $\alpha = 0.05$ ), and at least one data-driven approach is more likely to outperform *RAkELd* than a priori methods in the case of *RAkELd*'s best performance. This is the largest *RAkELd* evaluation published to date with 250 samplings per value for 10 values of *RAkELd* parameter  $k$  on 12 datasets published to date.

**Keywords:** label space clustering; label co-occurrence; label grouping; multi-label classification; clustering; machine learning; random  $k$ -label sets; ensemble classification

## 1. Introduction

Shannon's work on the unpredictability of information content inspired a search for the area of multi-label classification that requires more insight: where has the field still been using random approaches to handling data uncertainty when non-random methods could shed light and provide the ability to make better predictions?

Interestingly enough, random methods are prevalent in well-cited and multi-label classification approaches, especially in the problem of label space partitioning, which is a core issue in the problem-transformation approach to multi-label classification.

A great family of multi-label classification methods, called problem transformation approaches, depends on converting an instance of a multi-label classification problem into one or more single-label

single-class or multi-class classification problems, performs such classification and then converts the results back to multi-label classification results.

Such a situation stems from the fact that historically, the field of classification started out with solving single-label classification problems; in general, a classification problem of understanding the relationship (function) between a set of objects and a set of categories that should be assigned to it. If the object is allocated to at most one category, the problem is called a single-label classification. When multiple assignments per instance are allowed, we are dealing with a multi-label classification scenario.

In the single-label scenario, in one variant, we deal with a case when there is only one category, i.e., the problem is a binary choice: whether to assign a category or not, such a scenario is called single-class classification, e.g., the case of classifying whether there is a car in the picture or not. The other case is when we have to choose at most one from many possible classes; such a case is called multi-class classification, i.e., classifying a picture with the dominant brand of cars present in it. The multi-label variant of this example would concern classifying a picture with all car brands present in it.

As both single- and multi-class classification problems have been considerably researched during the last few decades, one can naturally see that it is reasonable to transform the multi-label classification case, by dividing the label space, into a single- or a multi-class scenario. A great introduction to the field was written by [1].

The two basic approaches to such a transformation are binary relevance and label powerset. Binary Relevance (BR) takes an a priori assumption that the label space is entirely separable, thus converting the multi-label problem into a family of single-class problems, one for every label, and making a decision whether to apply it or not. Converting the results back to multi-label classification is based just on taking a union of all assigned labels. Regarding our example, binary relevance assumes that correlations between car brands are not important enough and discards them a priori by classifying with each brand separately.

Label Powerset (LP) makes an opposite a priori assumption: the label space is non-divisible label-wise, and each possible label combination becomes a unique class. Such an approach yields a multi-class classification problem on a set of classes equal to the powerset of the label set, i.e., growing exponentially if one treats all label combinations as possible. In practice, this would be intractable. Thus, as [2] note, label powerset is most commonly implemented to handle only these combinations of classes that occur in the training set and, as such, is prone to overfitting. It is also, per [3], prone to differences in label combination distributions between the training set and the test set, as well as to an imbalance in how label combinations are distributed in general.

To remedy the overfitting of label powerset, Tsoumakas et al. [3] propose to divide the label space into smaller subspaces and use label powerset in these subspaces. The source of proposed improvements come from the fact that it should be easier for label powerset to handle a large number of label space combinations in a smaller space. Two proposed approaches are called random  $k$ -label sets (RA $k$ EL). RA $k$ EL comes in two variants: a label space partitioning RA $k$ EL $_d$ , which divides the label set into  $k$  disjoint subsets, and RA $k$ EL $_o$ , which is a sampling approach that allows overlapping of label subspaces. In our example, RA $k$ EL would randomly select a subset of brands and use the label powerset approach for brand combinations in all of the subspaces.

While these methods were developed, we saw advances in other fields that brought us more and more tools to explore relations between entities in data. Social and complex networks have been flourishing after most of the well-established methods were published. In this paper, we propose a data-driven approach for label space partitioning in multi-label classification. While we tackle the problem of classification, our goal is to spark a reflection on how data-driven approaches to machine learning using new methods from complex/social networks can improve established procedures. We show that this direction is worth pursuing, by comparing method-driven and data-driven approaches towards partitioning the label space.

Why should one rely on label space division at random? Should not a data-driven approach be better than random choice? Are methods that perform simplistic a priori assumptions truly worse than the random approach? What are the variances of result quality upon label space partitioning? Instead of selecting random subspaces of brands, we could consider that some city brands occur more often with each other and less so with other suburban brands. Based on such a premise, we could build a weighted graph depicting the frequency of how often two brands occur together in photos. Then, using well-established community detection methods on this graph, we could provide a data-driven partition for the label space.

In this paper, we wish to follow Shannon's ambition to search for a data-driven solution, an approach of finding structure instead of accepting uncertainty. We run RAKELd on 12 benchmark datasets, with different values of parameter  $k$ , taking  $k$  to be equal to 10%, 20%, ..., 90% of the label set size. We draw 250 distinct partitions of the label set into subsets of  $k$  labels, per every value of  $k$ . In case there are less than 250 possible partitions of the label space (e.g., because there are less than 10 labels), we consider all possible partitions. We then compare these results against the performance of methods based both on a priori assumptions—binary relevance and label powerset—and also well-established community detection methods employed in social and complex network analysis to detect patterns in the label space. For each of the measures, we state three hypotheses:

- RH1: a data-driven approach performs statistically better than the average random baseline;
- RH2: a data-driven approach is more likely to outperform RAKELd than methods based on a priori assumptions;
- RH3: a data-driven approach has a higher likelihood to outperform RAKELd in the worst case than methods based on a priori assumptions;
- RH4: a data-driven approach is more likely to perform better than RAKELd, than otherwise, i.e., the worst-case likelihood is greater than 0.5;
- RH5: the data-driven approach is more time efficient than RAKELd.

We describe the multi-label classification problem and existing methods in Section 3, our new proposed data-driven approach in Section 4 and compare the results of the likelihood of a data-driven approach being better than randomness in Section 6. We provide the technical detail of the experimental scenario in Section 5. We conclude with the main findings of the paper and future work in Section 7.

## 2. Related Work

Our study builds on two kinds of approaches to multi-label classification: problem transformation and ensemble. We extend the label powerset problem transformation method by employing an ensemble of classifiers to classify obtained partitions of the label space separately. We show that partitioning the label space randomly, as is done in the RAKEL approach, can be improved by using a variety of methods to infer the structure of partitions from the training. We extend the original evaluation of random  $k$ -label sets' performance using a larger sampling of the label space and providing deeper insight into how RAKELd performs. We also provide some insights into the nature of random label space partitioning in RAKELd. Finally, we provide alternatives to random label space partitioning that are very likely to yield better results than RAKELd depending on the selected measure, and we show which methods to use depending on the generalization strategy.

The classifier chains [4] approach to label space partitioning is based on a Bayesian conditioning scheme, in which labels are ordered in a chain, and then, the  $n$ -th classification is performed taking into account the output of the last  $n - 1$  classifications. These methods suffer from a variety of challenges: the results are not stable when ordering of labels in the chain changes, and finding the optimal chain is NP-hard. Existing methods that optimize towards best quality cannot handle more than 15 labels in a dataset (e.g., Bayes-optimal Probabilistic Classifier Chains (PCC) [5]). Furthermore,

in every classifier chain approach, one always needs to train at least the same number of classifiers as there are labels and if ensemble approaches are applied, much more. In our approach, we use community detection methods to divide the label space into a fewer number of cases to classify as multi-class problems, instead of transforming to a large number of single-class problems that are interdependent. We also do not strive to find the directly optimal solution to community detection problems on the label co-occurrence graph, to avoid overfitting; instead, we perform approximations of the optimal solutions. This approach provides a large overhead over random approaches. We note that it would be an interesting question whether the random orderings in classifier chains are as suboptimal of a solution, as random partitioning turns out to be in label space partitioning. Yet, it is not the subject of the study and is open to further research.

Tsoumakas et al.'s [6] Hierarchy Of Multilabel classifierS (HOMER) is a method of two-step hierarchical multi-label classification in which the label space is divided based on label assignment vectors; then, observations are classified first with cluster meta-labels; and finally, for each cluster they were labeled with, they are classified with labels of that cluster. We do not compare to HOMER directly in this article, due to the different nature of the classification scheme, as the subject of this study is to evaluate how data-driven label space partitioning using complex/social network methods, which can be seen as weak classifiers (as all objects are assigned automatically to all subsets), can improve the random partitioning multi-label classification. HOMER uses a strong classifier to decide which object should be classified in which subspace. Although we do not compare directly to HOMER, due to the difference in the classification scheme and base classifier, our research shows similarities to Tsoumakas's result that abandoning random label space partitioning for the  $k$ -means-based data-driven approach improves classification results. Thus, our results are in accord, yet we provide a much wider study, as we have performed a much larger sampling of the random space than the authors of HOMER in their method-describing paper.

Madjarov et al. [7] compare the performance of 12 multi-label classifiers on 11 benchmark sets evaluated by 16 measures. To provide statistical support, they use a Friedman multiple comparison procedure with a post hoc Nemenyi test. They include the RAKELo procedure in their study, i.e., the random label subsetting instead of partitioning. They do not evaluate the partitioning strategy RAKELd, which is the main subject of this study. Our main contribution, the study of how RAKELd performs against more informed approaches, therefore fills the unexplored space of Madjarov et al.'s extensive comparison. Note that, due to computational limits, we use Classification and Regression Trees (CART) instead of Support Vector Machines (SVM) as the single-label base classifier, as explained in Section 5.2.

Zhang et al. [8] review theoretical aspects and reported the experimental performance of eight multi-label algorithms and categorize them by the order of correlations taken into account and the evaluation measure that they try to optimize.

### 3. Multi-Label Classification

In this section, we aim to provide a more rigid description of the methods we use in the experimental scenario. We start by formalizing the notion of classification. Classification aims to understand a phenomenon, a relationship (function  $f : X \rightarrow Y$ ) between objects and categories, by generalizing from empirically-collected data  $D$ :

- objects are represented as feature vectors  $\bar{x}$  from the input space  $X$ ;
- categories, i.e., labels or classes come from a set  $L$ , and it spans the output space  $Y$ :
  - in the case of single-label single-class classification,  $|L| = 1$  and  $Y = \{0, 1\}$
  - in the case of single-label multi-class classification,  $|L| > 1$  and  $Y = \{0, 1, \dots, |L|\}$
  - in the case of multi-label classification,  $Y = 2^L$
- the empirical evidence collected:  $D = (D_x, D_y) \subset X \times Y$ ;
- a quality criterion function  $q$ .

In practice the empirical evidence  $D$  is split into two groups: the training set for learning the classifier and the test set to use for evaluating the quality of classifier performance. For the purpose of this section, we will denote  $D_{train}$  as the training set.

The goal of classification is to learn a classifier  $h : X \rightarrow Y$ , such that  $h$  generalizes  $D_{train}$  in a way that maximizes  $q$ .

We are focusing on problem-transformation approaches that perform multi-label classification by transforming it to a single-label classification and convert the results back to the multi-label case. In this paper, we use CART decision trees as the single-label base classifier. CART decision trees are a single-label classification method capable of both single- and multi-class classification. A decision tree constructs a binary tree in which every node performs a split based on the value of a chosen feature  $X_i$  from the feature space  $X$ . For every feature, a threshold is found that minimizes an impurity function calculated for a threshold on the data available in the current node's subtree. For the selected pair of a feature  $X_i$  and a threshold  $\theta$ , a new split is performed in the current node. Objects with values of the feature lesser than  $\theta$  are evaluated in the left subtree of the node, while the rest in the right. The process is repeated recursively for every new node in the binary tree until the depth limit selected for the method is reached or there is just one observation left to evaluate. In all of our scenarios, we use CART as the base single-label single- and multi-class base classifier.

Binary relevance learns  $|L|$  single-label single-class base classifiers  $b_j : X \rightarrow \{0, 1\}$  for each  $L_j \in L$  and outputs the multi-label classification using a classifier  $h(\bar{x}) = \{L_j \in L : b_j(\bar{x}) = 1\}$ .

Label powerset constructs a bijection  $lp : 2^L \rightarrow C$  between each subset of labels  $L_i$  and a class  $C_i$ . Label powerset then learns a single-label multi-class base classifier  $b : X \rightarrow C$  and transforms its output into a multi-label classification result  $lp^{-1}(b(\bar{x}))$ .

RAkELd performs a random partition of the label set  $L$  into  $k$  subsets  $L_j|_1^k$ . For each  $L_j$ , a label powerset classifier  $b_j : X \rightarrow L_j$  is learned. For a given input vector  $\bar{x}$ , the RAkELd classifier  $h$  performs multi-label classification with each  $b_j$  classifier and then sums the results, which can be formally described as  $h(\bar{x}) = \bigcup_{j=0}^k b_j(\bar{x})$ . Following the RAkELd scenario from [3], all partitions of the set  $L$  into  $k$  subsets are equally probable.

#### 4. The Data-Driven Approach

Having described the baseline random scenario of RAkELd, we now turn to explaining how complex/social network community detection methods fit into a data-driven perspective for label space division. In this scenario, we are transforming the problem exactly like RAkELd, but instead of performing random space partitioning, we construct a label co-occurrence graph from the training data and perform community detection on this graph to obtain a label space division.

##### 4.1. Label Co-Occurrence Graph

We construct the label co-occurrence graph as follows. We start with an undirected co-occurrence graph  $G$  with the label set  $L$  as the vertex set and the set of edges constructed from all pairs of labels that were assigned together at least once to an input object  $\bar{x}$  in the training set (here,  $l_{i,j}, \dots$  denote labels, i.e., elements of the set  $L$ ):

$$E = \{ \{ \lambda_i, \lambda_j \} : (\exists (\bar{x}, \Lambda) \in D_{train}) (\lambda_i \in \Lambda \wedge \lambda_j \in \Lambda) \}$$

One can also extend this unweighted graph  $G$  to a weighted graph by defining a function  $w : L \rightarrow \mathbb{N}$ :

$$\begin{aligned} w(\lambda_i, \lambda_j) &= \text{number of input objects } \bar{x} \text{ that have both labels assigned} = \\ &= \left| \{ \bar{x} : (\bar{x}, \Lambda) \in D_{train} \wedge \lambda_i \in \Lambda \wedge \lambda_j \in \Lambda \} \right| \end{aligned}$$

Using such a graph  $G$ , weighted or unweighted, we find a label space division by using one of the following community detection methods to partition the graph's vertex set, which is equal to the label set.

#### 4.2. Dividing the Label Space

Community detection methods are based on different principles as different fields defined communities differently. We are employing a variety of methods.

Modularity-based approaches, such as the fast greedy [9] and the spectral leading eigenvector algorithms, are based on detecting a partition of label sets that maximizes the modularity measure by [10]. Behind this measure lies an assumption that true community structure in a network corresponds to a statistically-surprising arrangement of edges [10], i.e., that a community structure in a real phenomenon should exhibit a structure different from an average case of a random graph, which is generated under a given null model. A well-established null model is the configuration model, which joins vertices together at random, but maintains the degree distribution of vertices in the graph.

For a given partition of the label set, the modularity measure is the difference between how many edges of the empirically-observed graph have both ends inside of a given community, i.e.,  $e(C) = \{(u, v) \in E : u \in C \wedge v \in C\}$  versus how many edges starting in this community would end in a different one in the random case:  $r(C) = \frac{\sum_{v \in C} \text{deg}(v)}{|E|}$ . More formally, this is  $Q(C) = \sum_{c \in C} e(c) - r(c)$ . In the case of weights, instead of counting the number of edges, the total weight of edges is used, and instead of taking vertex degrees in  $r$ , the vertex strengths are used; a precise description of weighted modularity can be found in Newman's paper [11].

Finding  $\bar{C} = \text{argmax}_C Q(C)$  is NP-hard, as shown by Brandes et al. [12]. We thus employ three different approximation-based techniques instead: a greedy, a multi-level hierarchical and a spectral recursively-dividing algorithm.

The fast greedy approach works based on greedy aggregation of communities, starting with singletons and merging the communities iteratively. In each iteration, the algorithm merges two communities based on which merge achieves the highest contribution to modularity. The algorithm stops when there is no possible merge that would increase the value of the current partition's modularity. Its complexity is  $O(N \log^2(N))$ .

The leading eigenvector approximation method depends on calculating a modularity matrix for a split of the graph into two communities. Such a matrix allows one to rewrite the two-community definition of modularity in a matrix form that can be then maximized using the largest positive eigenvalue and signs of the corresponding elements in the eigenvector of the modularity matrix: negative ones assigned to one community, positive ones to another. The algorithm starts with all labels in one community and performs consecutive splits recursively until all elements of the eigenvector have the same sign or the community is a singleton. The method is based on the simplest variant of spectral modularity approximation as proposed by Newman [13]. Its complexity is  $O(M + N^2)$ .

The infomap algorithm concentrates on finding the community structure of the network with respect to flow and to exploit the inference-compression duality to do so [14]. It relies on finding a partition of the vertex set that minimizes the map equation. The map equation divides flows through the graph into intra-community ones and the between-community ones and takes into consideration an entropy-based frequency-weighted average length of codewords used to label nodes in communities and inter-communities. Its complexity is  $O(M)$ .

The label propagation algorithm [15] assigns a unique tag to every vertex in the graph. Next, it iteratively updates the tag of every vertex with the tag assigned to the majority of the elements' neighbors. The updating order is randomly selected at each iteration. The algorithm stops when all vertices have tags that are in accord with the dominant tag in their neighborhood. Its complexity is  $O(N + M)$ .

The walktrap algorithm is based on the intuition that random walks on a graph tend to get “trapped” into densely-connected parts corresponding to communities [16]. It starts with a set of singleton communities and agglomerates obtained communities in a greedy iterative approach based on how close two vertices are in terms of random-walk distance. More precisely, each step merges two communities to maximize the decrease of the mean (averaged over vertices) of squared distances between a vertex and all of the vertices that are in the vertex’s community. The random walk distance between two nodes is measured as the  $L^2$  distance between the random walk probability distribution starting in each of the nodes. The distances are of the same maximum length provided as a parameter to the method. Its expected complexity is  $O(N^2 * \log(N))$ .

In complexity notation,  $N$  is the number of nodes and  $M$  the number of edges.

#### 4.3. Classification Scheme

In our data-driven scheme, the training phase is performed as follows:

1. the label co-occurrence graph is constructed based on the training dataset;
2. the selected community detection algorithm is executed on the label co-occurrence graph;
3. for every community  $L_i$ , a new training dataset  $D_i$  is created by taking the original input space with only the label columns that are present in  $L_i$ ;
4. for every community, a classifier  $h_i$  is learned on training set  $D_i$ .

The classification phase is performed by performing classification on all subspaces detected in the training phase and taking the union of assigned labels:  $h(\bar{x}) = \bigcup_{j=0}^k b_j(\bar{x})$ .

### 5. Experiments and Materials

To prepare the ground for results, in this section, we describe which datasets we have selected for evaluation and why. Then, we present and justify model selection decisions for our experimental scheme. Next, we describe the configuration of our experimental environment. Finally, we describe the measures used for evaluation.

#### 5.1. Datasets

Following Madjarov’s study [7], we have selected 12 different well-cited multi-label classification benchmark datasets. The basic statistics of datasets used in experiments, such as the number of data instances, the number of attributes, the number of labels, the labels’ cardinality, density and the distinct number of label combinations, are available online [17]. We selected the datasets to obtain a balanced representation of problems in terms of the number of objects, the number of labels and domains. At the moment of publishing, this is one of the largest studies of RAKELd, both in terms of datasets examined and in terms of random label partitioning sample count. This study also exhibits a higher ratio of the number of datasets to the number of methods than other studies.

The text domain is represented by 5 datasets: *bibtex*, *delicious*, *enron*, *medical* and *tmc2007-500*. *Bibtex* [18] comes from the ECML/PKDD 2008 Discovery Challenge and is based on data from the Bibsonomy.org publication sharing and bookmarking website. It exemplifies the problem of assigning tags to publications represented as an input space of bibliographic metadata, such as: authors, paper/book title, journal volume, etc. *Delicious* [6] is another user-tagged dataset. It spans over 983 labels obtained by scraping the 140 most popular tags from the del.icio.us bookmarking website, retrieving the 1000 most recent bookmarks, selecting the 200 most popular, deduplication and filtering tags that were used to tag less than 10 websites. For those labels, websites tagged with them were scraped, and from their contents, the top 500 words ranked by the  $\chi^2$  method were selected as input features. *Tmc2007-500* [6] contains an input space consisting of the similarly selected top 500 words appearing in flight safety reports. The labels represent the problems being described in these reports. *Enron* [19] contains emails from senior Enron Corporation employees categorized into topics by the UC Berkeley Enron E-mail Analysis Project [20] with the input space being a bag of

word representation of the e-mails. The *Medical* [4] dataset is from the Medical Natural Language Processing Challenge [21]. The input space is a bag-of-words representation of patient symptom history, and labels represent diseases following the International Classification of Diseases.

The multimedia domain consists of five datasets: *scene*, *corel5k*, *mediamill*, *emotions* and *birds*. The image dataset *scene* [22] semantically indexes still scenes annotated with any of the following categories: beach, fall-foliage, field, mountain, sunset and urban. The *birds* dataset [23] represents a problem of matching bird voice recordings' extracted features with a subset of 19 bird species that are present in the recording; each label represents one species. This dataset was introduced [24] during the The 9th Annual MLSP competition. A larger image set *corel5k* [25] contains normalized-cut segmented images clustered into 499 bins. The bins were labeled with 374 subset labels. The *mediamill* dataset of annotated video segments [26] was introduced during the 2005 NIST TRECVID challenge [27]. It is annotated with 101 labels referring to elements observable in the video. The *emotions* dataset [28] represents the problem of the automated detection of emotion in music, assigning a subset of music emotions based on the Tellegen–Watson–Clark model to each of the songs.

The biological domain is represented with two datasets: *yeast* and *genbase*. The *yeast* [29] dataset concerns the problem of assigning functional classes to genes of the *Saccharomyces cerevisiae* genome. The *genbase* [30] dataset represents the problem of assigning classes to proteins based on detected motifs that serve as input features.

## 5.2. Experiment Design

Using 12 benchmark datasets evaluated with five performance measures, we compare eight approaches to label space partitioning for multi-label classification:

- five methods that divide the label space based on structure inferred from the training data via label co-occurrence graphs, in both unweighted and weighted versions of the graphs;
- two methods that take an a priori assumption about the nature of the label space: binary relevance and label powerset;
- one random label space partitioning approach that draws partitions with equal probability: *RAkELd*.

In the random baseline (*RAkELd*), we perform 250 samplings of random label space partitions into  $k$ -label subsets for each evaluated value of  $k$ . If a dataset had more than 10 labels, we took values of  $k$  ranging from 10% to 90% with a step of 10%, rounding to the closest integer number if necessary. In case of two datasets with a smaller number of labels, i.e., *scene* and *emotions*, we evaluated *RAkELd* for all possible label space partitions due to their low number. The number of label space division samples per dataset can be found in Appendix A (Table A1).

As no one knows the true distribution of classification quality over label space partitions, we have decided to use a large number of samples, 250 per each of the groups, 2500 altogether, to get as close to a representative sample of the population as was possible with our infrastructure limitations.

As the base classifier, we use CART decision trees. While we recognize that the majority of studies prefer to use SVMs, we note that it is intractable to evaluate nearly 32,500 samples of the random label space partitions using SVMs. We have thus decided to use a classifier that presents a reasonable trade-off between quality and computational speed.

We perform statistical evaluation of our approaches by comparing them to average performance of the random baseline of *RAkELd*. We average *RAkELd* results per dataset, which is justified by the fact that this is the expected result one would get without performing extensive parameter optimization. Following Derrac et al.'s [31] de facto standard modus operandi, we use the Friedman test with Iman–Davenport modifications to detect differences between methods, and we check whether a given method is statistically better than the average random baseline using Rom's post-hoc pairwise test. We use these tests' results to confirm or reject RH1.



We do not perform statistical evaluation per group (i.e., isolating each value of  $k$  from 10% to 90%) due to the lack of non-parametric repeated measure tests, as noted by Demsar in the classic paper [32].

Instead, to account for variation, we consider the probability that a given data-driven approach to label space division is better than random partitioning. These probabilities were calculated per dataset, as the fraction of random outputs that yielded worse results than a given method. Thus, for example, if infomap has a 96.5% probability of having higher better subset accuracy (SA) than the random approach in Corel5k, this means that on this dataset, infomap's SA score was better than the scores achieved by 96.5% of all RAKELd experiments. We check the median, the mean and the minimal (i.e., worst case) likelihoods. We use these results to confirm or reject RH2, RH3 and RH4.

### 5.3. Environment

We used `scikit-multilearn` (Version 0.0.1) [33], a `scikit-learn` API compatible library for multi-label classification in python that provides its own implementation of several classifiers and uses `scikit-learn` [34] multi-class classification methods. All of the datasets come from the MULAN [35] dataset library [17] and follow MULAN's division into the train and test subsets.

We use CART decision trees from the `scikit-learn` package (Version 0.15), with the Gini index as the impurity function. We employ community detection methods from the Python version of the `igraph` library [36] for both weighted and unweighted graphs. The performance measures' implementation comes from the `scikit-learn` `metrics` package.

### 5.4. Evaluation Methods

Following Madjarov et al.'s [7] taxonomy of multi-label classification evaluation measures, we use three example-based measures: Hamming loss, subset accuracy and Jaccard similarity, as well as a label-based measure, F1, as evaluated by two averaging schemes: micro and macro. The following definitions are used:

- $X$  is the set of objects used in the testing scenario for evaluation
- $L$  is the set of labels that spans the output space  $Y$ ;
- $\bar{x}$  denotes an example object undergoing classification;
- $h(\bar{x})$  denotes the label set assigned to object  $\bar{x}$  by the evaluated classifier  $h$ ;
- $y$  denotes the set of true labels for the observation  $\bar{x}$ ;
- $tp_j, fp_j, fn_j, tn_j$  are respectively true positives, false positives, false negatives and true negatives of the of label  $L_j$ , counted per label over the output of classifier  $h$  on the set of testing objects  $\bar{x} \in X$ , i.e.,  $h(X)$ ;
- the operator  $[[p]]$  converts the logical value to a number, i.e., it yields 1 if  $p$  is true and 0 if  $p$  is false.

#### 5.4.1. Example-Based Evaluation Methods

Hamming loss is a label-wise decomposable function counting the fraction of labels that were misclassified.  $\otimes$  is the logical exclusive or.

$$\text{HammingLoss}(h) = \frac{1}{|X|} \sum_{\bar{x} \in X} \frac{1}{|L|} \sum_{L_j \in L} [[(L_j \in h(\bar{x})) \otimes (L_j \in y)]]$$

The accuracy score and subset 0/1 loss are instance-wise measures that count the fraction of input observations that have been classified exactly the same as in the golden truth.

$$\text{SubsetAccuracy}(h) = \frac{1}{|X|} \sum_{\bar{x} \in X} [[h(\bar{x}) = y]]$$

Jaccard similarity is a measure of the size of similarity between the prediction and the ground truth comparing what is the cardinality of an intersection of the two, compared to the union of the two. In other words, what fraction of all labels taken into account by any of the prediction or ground truth were assigned to the observation in both of the cases.

$$\text{Jaccard}(h) = \frac{1}{|X|} \sum_{\tilde{x} \in X} \frac{h(\tilde{x}) \cap y}{h(\tilde{x}) \cup y}$$

#### 5.4.2. Label-Based Evaluation Methods

The F1 measure is a harmonic mean of precision and recall where none of the two are more preferred than the other. Precision is the measure of how much the method is immune to Type I error, i.e., falsely classifying negative cases as positives: false positives or FP. It is the fraction of correctly positively-classified cases (i.e., true positives) to all positively-classified cases. It can be interpreted as the probability that an object without a given label will not be labeled as having it. Recall is the measure of how much the method is immune to the Type II error, i.e., falsely classifying positive cases as negatives: false negatives or FN. It is the fraction of correctly positively-classified cases (i.e., true positives) to all positively-classified label. It can be interpreted as the probability that an object with a given label will be labeled as such.

These measures can be averaged from two perspectives that are not equivalent in practice due to a natural non-uniformity of the distribution of labels among input objects in any testing set. Two averaging techniques are well-established, as noted by [37].

Micro-averaging gives equal weight to every input object and performs a global aggregation of true/false positives/negatives, averaging over all objects first. Thus:

$$\begin{aligned} \text{precision}_{\text{micro}}(h) &= \frac{\sum_{j=1}^{|L|} tp_j}{\sum_{j=1}^{|L|} tp_j + fp_j} \\ \text{recall}_{\text{micro}}(h) &= \frac{\sum_{j=1}^{|L|} tp_j}{\sum_{j=1}^{|L|} tp_j + fn_j} \\ \text{F1}_{\text{micro}}(h) &= 2 \cdot \frac{\text{precision}_{\text{micro}}(h) \cdot \text{recall}_{\text{micro}}(h)}{\text{precision}_{\text{micro}}(h) + \text{recall}_{\text{micro}}(h)} \end{aligned}$$

In macro-averaging, the measure is first calculated per label, then averaged over the number of labels. Macro averaging thus gives equal weight to each label, regardless of how often the label appears.

$$\begin{aligned} \text{precision}_{\text{macro}}(h, j) &= \frac{tp_j}{tp_j + fp_j} \\ \text{recall}_{\text{macro}}(h, j) &= \frac{tp_j}{tp_j + fn_j} \\ \text{F1}_{\text{macro}}(h, j) &= 2 \cdot \frac{\text{precision}_{\text{macro}}(h, j) \cdot \text{recall}_{\text{macro}}(h, j)}{\text{precision}_{\text{macro}}(h, j) + \text{recall}_{\text{macro}}(h, j)} \\ \text{F1}_{\text{macro}}(h) &= \frac{1}{|L|} \sum_{j=1}^{|L|} \text{F1}_{\text{macro}}(h, j) \end{aligned}$$

## 6. Results and Discussion

We describe the performance per measure first and then look at how methods behave across measures. We evaluate each of the research hypotheses, RH1 to RH4, for each of the measures. We then look at how these methods performed across datasets. We compare the median and the mean of the achieved probabilities to assess the average advantage over randomness; the higher the better. We compare the median and the means, as in some cases, the methods admit a single worst-performing outlier, while in general providing a great advantage over random approaches. We also check how each method performs in the worst case, i.e., what is the minimum probability of it being better than randomness in label space division?

### 6.1. Micro-Averaged F1 Score

When it comes to ranking of how well methods performed in micro-averaged F1, fast greedy and walktrap approaches used on a weighted label co-occurrence graph performed best, followed by BR, leading eigenvector and unweighted walktrap/modularity-maximizations methods. Furthermore, weighted label propagation and infomap were statistically significantly better than the average random performance. We confirm RH1, with evidence presented in Figures 1 and 2 and Table 1.

fastgreedy.weighted	0
walktrap.weighted	0
BR	5e−05
leading.eigenvector.weighted	7e−05
fastgreedy	7e−05
leading.eigenvector	0.00024
walktrap	0.00045
label.propagation.weighted	0.00241
infomap.weighted	0.00241
label.propagation	
infomap	
LP	
random	

**Figure 1.** Statistical evaluation of the method's performance in terms of micro-averaged F1 score. Gray, baseline; white, statistically identical to the baseline; otherwise, the  $p$ -value of the hypothesis that a method performs better than the baseline.

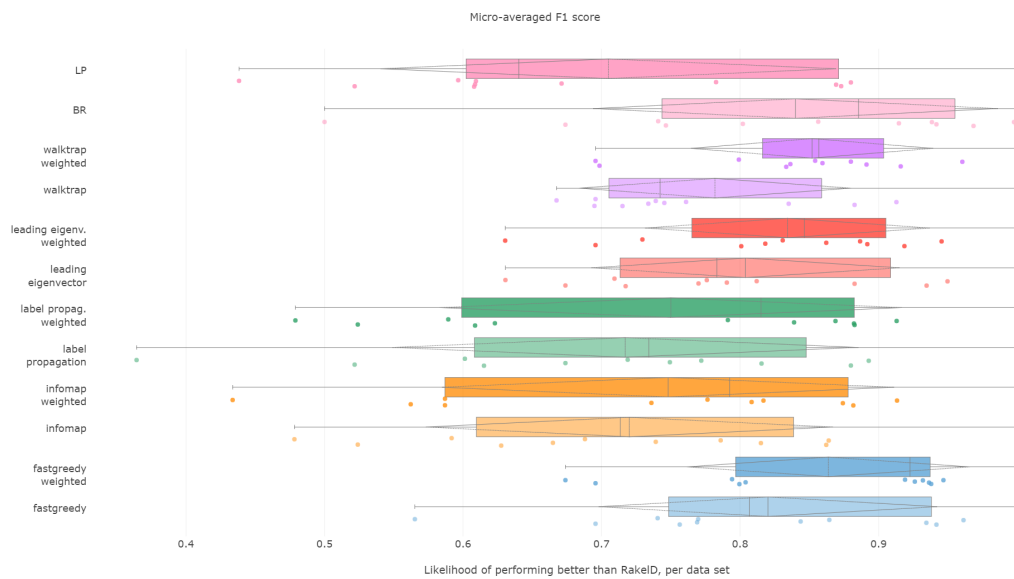
In terms of the micro-averaged F1, the weighted fast greedy approach has both the highest mean (86%) and median (92%) likelihood of scoring better than random baseline. Binary relevance and weighted variants of walktrap, leading eigenvector also performed well with a mean likelihood of 83% to 85% and a lower, but still satisfactory median of 85% to 88%. We confirm RH2.

Modularity-based approaches also turn out to be most resilient. The weighted variant of walktrap was the most resilient with a 69.5% likelihood in the worst case, followed closely by a weighted fast greedy approach with 67% and unweighted walktrap with 66.7%. We note that, apart from a single outlying datasets, all methods (apart from label powerset) had better than 50% likelihood of performing better than RAKELd. Binary relevance's worst case likelihood was exactly 0.5. We thus confirm both RH3 and RH4.

Fast greedy and walktrap weighted approaches yielded the best advantage over RAKELd, both in the average and worst cases. Binary relevance also provided a strong overhead against random label space division, while achieving just 0.5 in the worst case scenario. Thus, when it comes to micro-averaged F1 scores, RAKELd random approaches to label space partitions should be dropped in favor of weighted fast greedy and walktrap methods or binary relevance. All of these methods are

also statistically significantly better than the average random baseline. We therefore confirm RH1, RH2, RH3 and RH4 for micro-averaged F1 scores.

We also note that RAKELd was better than label powerset on micro-averaged F1 in 57% of the cases in the worst case, while Tsoumakas et al.’s original paper [3] provides argumentation of micro-F1 improvements over LP yielded by RAKELd, using SVMs. We note that our observation is not contrary: LP failed to produce significantly different results than the average random baseline in our setting. Instead, our results are complementary, as we use a different classifier, but the intuition can be used to comment on Tsoumakas et al.’s results. While in some cases, RAKELd provides an improvement over LP in F1 score, on average, the probability of drawing a random subspace better is only 30%. We still note that it is much better to use one of the recommended community detection-based approaches instead of a method based on a priori assumptions.



**Figure 2.** Histogram of the methods’ likelihood of performing better than RAKELd in the micro-averaged F1 score aggregated over datasets.

**Table 1.** Likelihood of performing better than RAKELd in the micro-averaged F1 score of every method for each dataset.

	BR	LP	Fast Greedy	Fast Greedy-Weighted	Infomap	Infomap-Weighted	Label-Propagation	Label-Propagation-Weighted	Leading-Eigenvector	Leading-Eigenvector-Weighted	Walktrap	Walktrap-Weighted
Corel5k	0.856444	0.608000	0.961333	0.804000	0.524000	0.881778	0.601333	0.524000	0.949778	0.818222	0.745333	0.799111
bibtex	0.997778	0.782667	0.756444	0.794222	0.664889	0.816889	0.749333	0.882222	0.812000	0.800889	0.835111	0.833333
birds	0.968562	0.438280	0.843736	0.946833	0.591771	0.433657	0.364309	0.478964	0.630606	0.830791	0.694868	0.836338
delicious	0.914667	0.869333	0.941778	0.936444	0.864000	0.874222	0.892889	0.868889	0.934667	0.918667	0.912889	0.916000
emotions	0.500000	0.521739	0.565217	0.673913	0.739130	0.586957	0.673913	0.913043	0.717391	0.630435	0.739130	0.891304
enron	0.802000	0.873000	0.934500	0.938000	0.786000	0.776500	0.815500	0.839000	0.776000	0.945500	0.761000	0.859500
genbase	0.941778	0.880000	0.864444	0.919111	0.862222	0.913333	0.880000	0.882667	0.882667	0.862222	0.882667	0.880000
mediamill	0.740889	0.609333	0.769778	0.932000	0.627556	0.562222	0.615111	0.589333	0.709333	0.886667	0.715111	0.854222
medical	0.938500	0.596500	0.769000	0.799500	0.688000	0.736000	0.772000	0.623000	0.770000	0.729500	0.667500	0.698500
scene	0.673913	0.608696	0.695652	0.695652	0.478261	0.586957	0.521739	0.608696	0.673913	0.695652	0.695652	0.695652
tmc2007-500	0.999343	1.000000	1.000000	1.000000	1.000000	1.000000	1.000000	1.000000	1.000000	1.000000	1.000000	1.000000
yeast	0.746458	0.671141	0.740492	0.926174	0.815063	0.808352	0.718867	0.791201	0.790455	0.891872	0.733781	0.960477

### 6.2. Macro-Averaged F1 Score

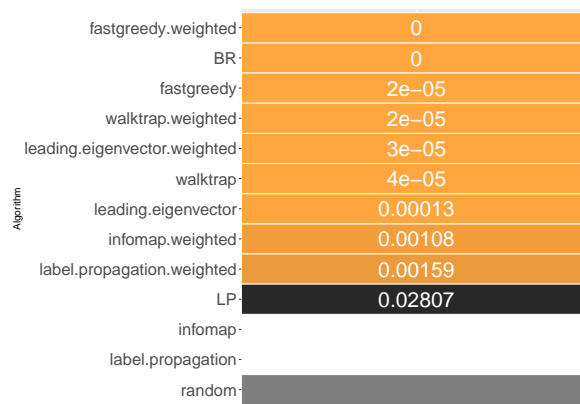
All methods, apart from unweighted label propagation and infomap, performed significantly better than the average random baseline. The highest ranks were achieved by weighted fast greedy, binary relevance and unweighted fast greedy. We confirm RH1, with evidence presented in Figures 3 and 4 and Table 2.

Fast greedy and walktrap approaches used on the weighted label co-occurrence graph were most likely to perform better than RAKELd samples, followed by BR, leading eigenvector and unweighted

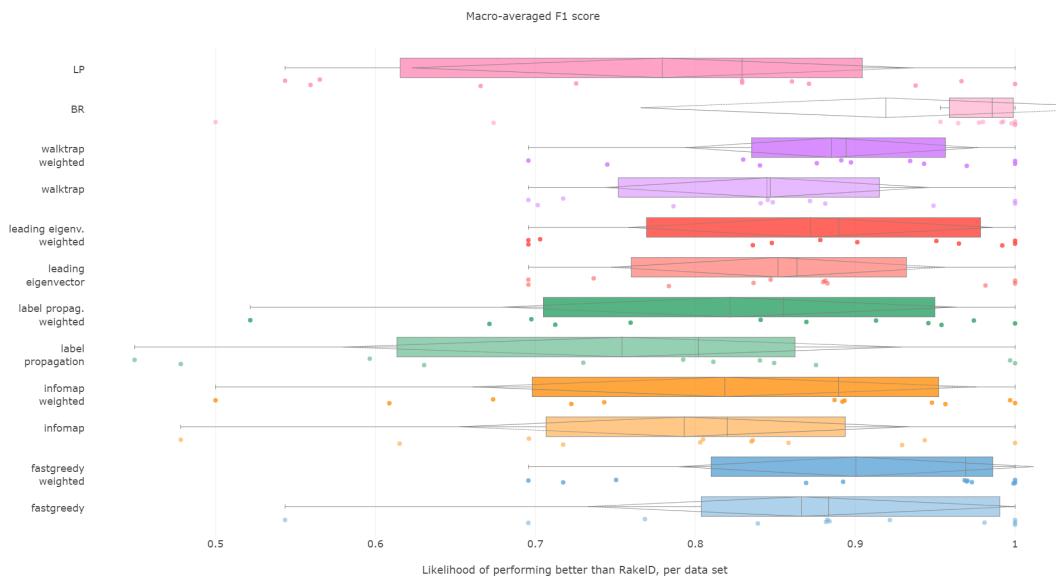
walktrap/modularity-maximizations methods. Furthermore, weighted label propagation and infomap were statistically significantly better than the average random performance.

Binary relevance and weighted fast greedy were the two approaches that surpassed the 90% likelihood of being better than random label space divisions in both the median (98.5% and 97%, respectively) and mean (92% and 90%) cases. Weighted walktrap and leading eigenvector followed closely with both the median and the mean likelihood of 87% to 89%. We thus reject RH2, as binary relevance achieved greater likelihoods than the best data-driven approach.

When it comes to resilience, all modularity (apart from unweighted fast greedy) methods achieve the same high worst-case 70% probability of performing better than RAKELd. Binary relevance underperformed in the worst case, being better exactly in 50% of the cases. All methods on all datasets, apart from the outlier case of infomap’s and label propagation’s performance on the scene dataset, are likely to yield a better macro-averaged F1 score than the random approaches. We confirm RH3 and RH4.



**Figure 3.** Statistical evaluation of the method’s performance in terms of macro-averaged F1 score. Gray, baseline; white, statistically identical to baseline; otherwise, the *p*-value of the hypothesis that a method performs better than the baseline.



**Figure 4.** Histogram of the methods’ likelihood of performing better than RAKELd in the macro-averaged F1 score aggregated over datasets.

**Table 2.** Likelihood of performing better than RAKELd in the macro-averaged F1 score of every method for each dataset.

	BR	LP	Fast Greedy	Fast Greedy-Weighted	Infomap	Infomap-Weighted	Label_Propagation	Label_Propagation-Weighted	Leading_Eigenvector	Leading_Eigenvector-Weighted	Walktrap	Walktrap-Weighted
Corel5k	1.000000	0.665778	0.980889	0.968444	0.615111	0.996889	0.596444	0.712444	0.836444	0.901333	0.845333	0.969778
bibtex	1.000000	0.871111	0.839111	0.869333	0.803111	0.887111	0.849333	0.945778	0.880000	0.878222	0.881333	0.876000
birds	0.992603	0.559408	0.883957	0.999075	0.804901	0.673601	0.449376	0.671290	0.736477	0.847896	0.786408	0.897365
delicious	0.997778	0.937778	1.000000	1.000000	0.929333	0.956444	0.996889	0.974222	1.000000	1.000000	1.000000	1.000000
emotions	0.500000	0.543478	0.543478	0.717391	0.717391	0.608696	0.630435	0.913043	0.695652	0.695652	0.717391	0.891304
enron	0.991500	0.966500	1.000000	0.973000	0.943500	0.948000	0.875500	0.954000	0.981500	0.992000	0.949000	0.830000
genbase	0.953333	0.829333	0.881778	0.892444	0.836000	0.892000	0.840444	0.840889	0.882667	0.836000	0.848444	0.840444
mediamill	0.964444	0.860444	0.882667	0.969778	0.835111	0.743111	0.792444	0.759556	0.881333	0.964889	0.840889	0.943111
medical	0.977500	0.725500	0.768500	0.750500	0.696000	0.722500	0.730000	0.697500	0.783500	0.703000	0.701500	0.745000
scene	0.673913	0.565217	0.695652	0.695652	0.478261	0.500000	0.478261	0.521739	0.695652	0.695652	0.695652	0.695652
tmc2007-500	1.000000	1.000000	1.000000	1.000000	1.000000	1.000000	1.000000	1.000000	1.000000	1.000000	1.000000	1.000000
yeast	0.979866	0.829232	0.921700	0.970172	0.858315	0.893363	0.811335	0.869500	0.847129	0.950783	0.871738	0.934377

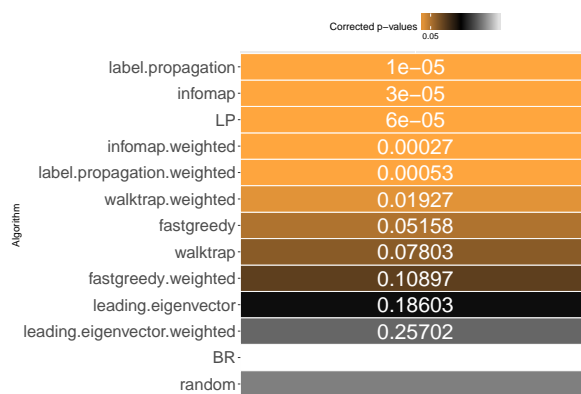
We recommend using binary relevance or weighted fast greedy approaches when generalizing to achieve the best macro-averaged F1 score, as they are both significantly better than average random performance and more likely to perform better than RAKELd samplings, and this likelihood is high even in the worst case. Thus, for macro-averaged F1, we confirm hypotheses RH1, RH3 and RH4. Binary relevance had a slightly better likelihood of beating RAKELd than data-driven approaches, and thus, we reject RH2.

### 6.3. Subset Accuracy

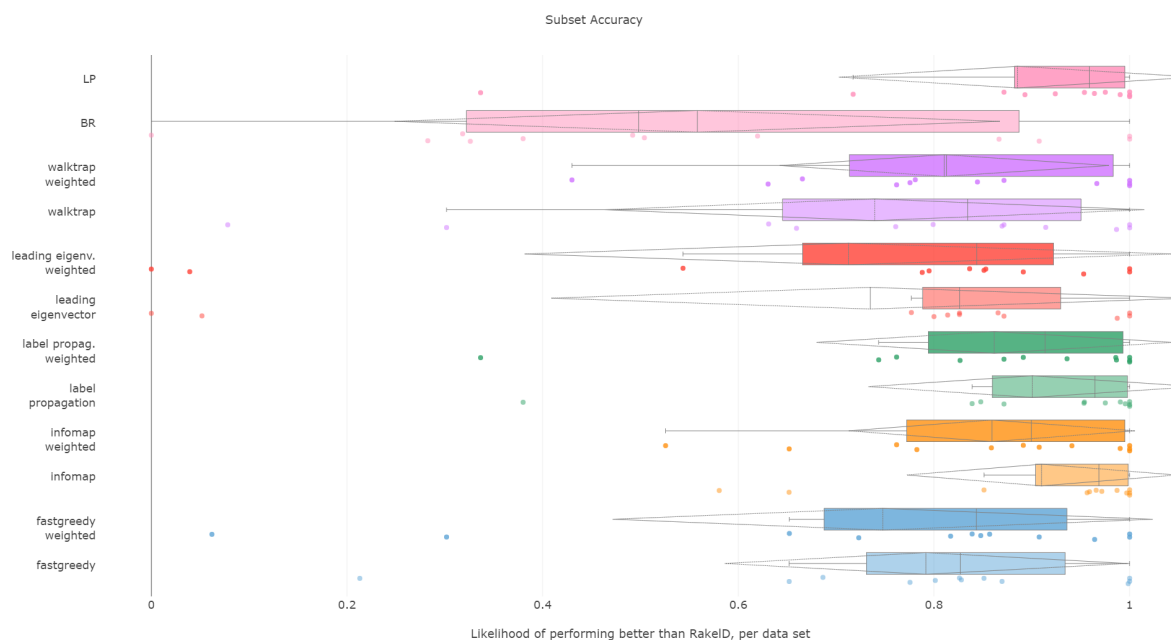
All methods apart from the weighted leading eigenvector modularity maximization approach were statistically significantly better than the average random baseline.

Label propagation, infomap, label powerset, weighted infomap, label propagation and walktrap are the methods that performed statistically significantly better than average random baseline, ordered by ranks. We confirm RH1, with evidence provided in Figures 5 and 6 and Table 3.

Furthermore, unweighted infomap and label propagation are the most likely to yield results of higher subset accuracy than random label space divisions, both regarding the median (96%) and the mean (90% to 91%) likelihood. Label powerset follows with a 95.8% median and 89% mean. Weighted versions of infomap and label powerset are fourth and fifth with five to six percentage points less. We confirm RH2.



**Figure 5.** Statistical evaluation of the method’s performance in terms of Jaccard similarity score. Gray, baseline; white, statistically identical to baseline; otherwise, the p-value of hypothesis that a method performs better than the baseline.



**Figure 6.** Histogram of the methods' likelihood of performing better than RAKELd in subset accuracy aggregated over datasets.

**Table 3.** Likelihood of performing better than RAKELd in the subset accuracy of every method for each dataset.

	BR	LP	Fast Greedy	Fast Greedy-Weighted	Infomap	Infomap-Weighted	Label-Propagation	Label-Propagation-Weighted	Leading-Eigenvector	Leading-Eigenvector-Weighted	Walktrap	Walktrap-Weighted
Corel5k	0.00000	0.953778	0.652000	0.301778	0.965778	0.652000	0.953778	0.826667	0.000000	0.000000	0.301778	0.780889
bibtex	0.492000	0.975111	0.828000	0.723111	0.971556	0.761778	0.975111	0.761778	0.800000	0.788000	0.799111	0.761778
birds	0.380028	0.336570	0.213130	0.061951	0.651872	0.525659	0.380028	0.336570	0.051780	0.039297	0.078132	0.429958
delicious	1.000000	1.000000	1.000000	1.000000	1.000000	1.000000	1.000000	1.000000	1.000000	1.000000	1.000000	1.000000
emotions	0.326087	0.717391	0.826087	0.847826	0.956522	0.891304	0.847826	0.891304	0.826087	0.891304	0.760870	1.000000
enron	0.504000	0.964000	0.775500	0.817000	0.959000	0.941000	0.990500	0.986500	0.865500	0.795000	0.659500	0.775500
genbase	0.907556	0.871556	0.851111	0.907556	0.851111	0.907556	0.871556	0.871556	0.851111	0.851111	0.871556	0.871556
mediamill	0.318222	0.924000	0.801333	0.856889	0.987111	0.858667	0.953333	0.936000	0.776889	0.836444	0.914222	0.844444
medical	0.866500	0.893000	0.686500	0.839000	0.580500	0.782500	0.839000	0.743500	0.814000	0.853000	0.631000	0.665500
scene	0.282609	1.000000	0.869565	0.652174	1.000000	1.000000	1.000000	1.000000	0.826087	0.543478	0.869565	0.630435
tmc2007-500	1.000000	1.000000	1.000000	1.000000	1.000000	1.000000	1.000000	1.000000	1.000000	1.000000	1.000000	1.000000
yeast	0.619687	0.990306	0.998509	0.964206	0.997017	0.990306	0.995526	0.985831	0.987323	0.953020	0.986577	0.966443

Concerning the resiliency of the advantage, only infomap versions proved to be better than RAKELd for more than half of the times: the unweighted version in 58% of cases, the weighted one in 52%. All other methods were below the 50% threshold in the worst case, with label powerset and label propagation likelihood of 33% for both variants. If one or two most wrong outliers were to be discarded, all methods are more than 50% likely to be better than random label space partitioning. We confirm RH3 and RH4.

We thus recommend using unweighted infomap as the data-driven alternative to RAKELd, as it is both significantly better than the random baseline, very likely to perform better than RAKELd and most resilient among the evaluated methods in the worst case. We confirm RH1, RH2, RH3 and RH4 for subset accuracy.

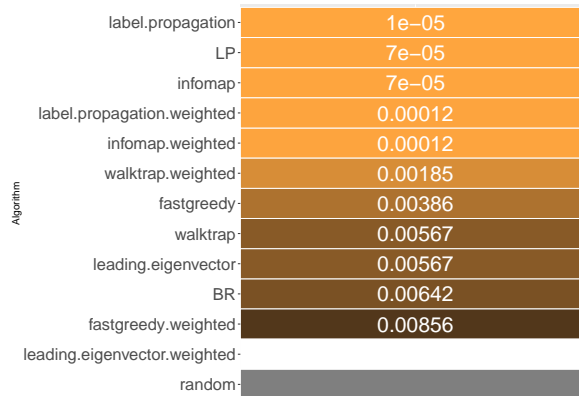
### 6.4. Jaccard Score

All methods apart from the weighted leading eigenvector modularity maximization approach were statistically significantly better than the average random baseline. Unweighted label propagation, label powerset and infomap were the highest ranked methods. We confirm RH1, with evidence provided in Figures 7 and 8 and Table 4.

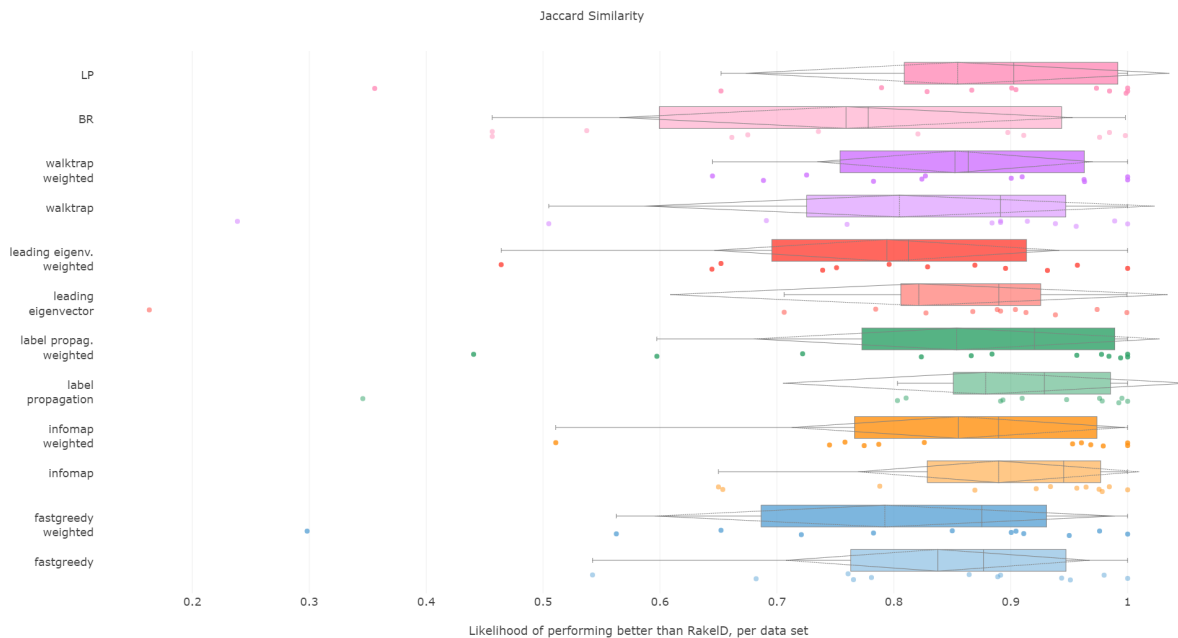
Jaccard score is similar to subset accuracy in rewarding exact label set matches. In effect, it is not surprising to see that unweighted infomap and label propagation are the most likely compared

to RAKELd to yield a result of higher Jaccard score, both in terms of median (94.5% and 92.9%, respectively) and mean likelihoods (88.9% and 87.9%, resp.). Out of the two, infomap provides the most resilient advantage with a 65% probability of performing better than random approaches in the worst case. Label propagation is in the worst cases only 34% to 35% likely to be better than random space partitions.

We recommend using unweighted infomap approach over RAKELd when Jaccard similarity is of importance and confirm RH1, RH2, RH3 and RH4 for this measure.



**Figure 7.** Statistical evaluation of the method’s performance in terms of micro-averaged F1 score. Gray, baseline; white, statistically identical to baseline, otherwise; the *p*-value of the hypothesis that a method performs better than the baseline.



**Figure 8.** Histogram of the methods’ likelihood of performing better than RAKELd in Jaccard similarity aggregated over datasets.



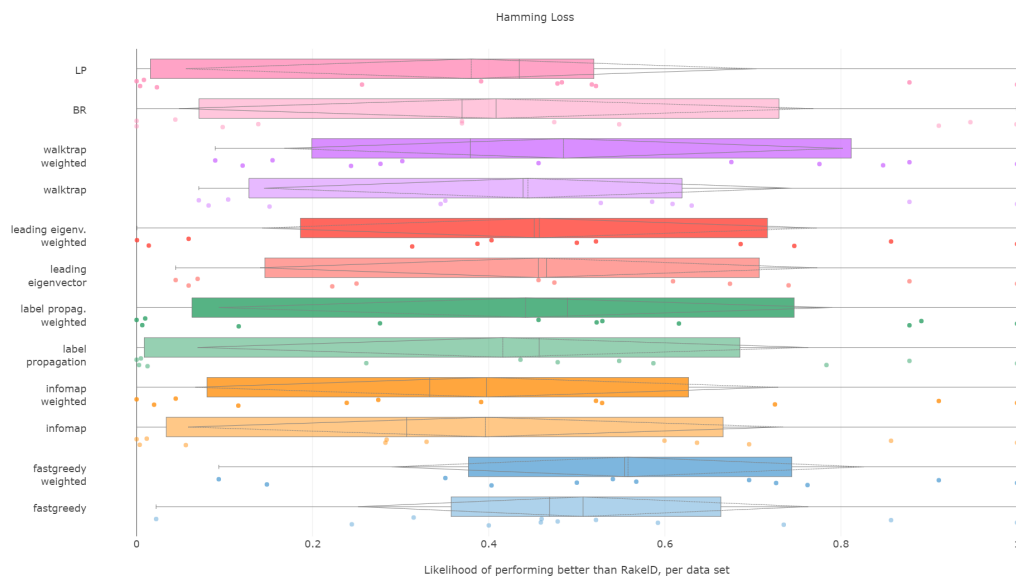
**Table 4.** Likelihood of performing better than RAKELd in Jaccard similarity of every method for each dataset.

	BR	LP	Fast Greedy	Fast Greedy-Weighted	Infomap	Infomap-Weighted	Label Propagation	Label Propagation-Weighted	Leading Eigenvector	Leading Eigenvector-Weighted	Walktrap	Walktrap-Weighted
Corel5k	0.675111	0.828444	0.888889	0.562667	0.788000	0.787111	0.803111	0.597333	0.888444	0.644444	0.504889	0.644889
bibtex	0.984444	0.998667	0.780889	0.720889	0.975356	0.758222	0.995111	0.866222	0.867556	0.796000	0.938222	0.824000
birds	0.820620	0.355987	0.542302	0.298197	0.653722	0.510865	0.345816	0.440592	0.163199	0.464170	0.238558	0.725381
delicious	0.976000	0.973333	0.943556	0.900444	0.964444	0.968444	0.992444	0.984000	0.938222	0.828889	0.988889	0.963111
emotions	0.456522	0.652174	0.760870	0.652174	0.956522	0.826087	0.891304	0.956522	0.891304	0.652174	0.891304	1.000000
enron	0.735500	0.984500	0.951000	0.904500	0.934000	0.960500	0.976000	0.994000	0.784500	0.957000	0.760000	0.827000
genbase	0.911111	0.904444	0.864444	0.911111	0.869333	0.952889	0.909778	0.884000	0.904000	0.869333	0.884000	0.909778
mediamill	0.537333	0.866667	0.682222	0.976000	0.921778	0.774667	0.893333	0.823556	0.706222	0.895556	0.914222	0.900444
medical	0.897500	0.789500	0.765500	0.850000	0.650000	0.745000	0.810500	0.722000	0.827500	0.751000	0.691000	0.688500
scene	0.456522	1.000000	0.891304	0.782609	0.978261	1.000000	0.978261	1.000000	0.913043	0.739130	0.891304	0.782609
tmc2007-500	0.998029	1.000000	1.000000	1.000000	1.000000	1.000000	1.000000	1.000000	1.000000	0.999343	1.000000	1.000000
yeast	0.661447	0.900820	0.979866	0.950037	0.984340	0.979120	0.947800	0.977629	0.973900	0.931394	0.956003	0.962714

6.5. Hamming Loss

Hamming loss is certainly a fascinating case in our experiments. As the measure is evaluated per each label separately, we can expect it to be the most stable over different label space partitions.

The first surprise comes with the Friedman–Iman–Davenport test result, where the test practically fails to find a difference in performance between random approaches and data-driven methods, yielding a *p*-value of 0.049. While the *p*-value is lower than  $\alpha = 0.05$ , the difference cannot be taken as significant given the characteristics of the test. Lack of significance is confirmed by pairwise tests against random baseline (all hypotheses of difference are strongly rejected). We reject RH1, with evidence provided in Figure 9 and Table 5.



**Figure 9.** Histogram of the methods' likelihood of performing better than RAKELd in Hamming loss similarity aggregated over datasets.

**Table 5.** Likelihood of performing better than RAKELd in Hamming loss of every method for each dataset.

	BR	LP	Fast Greedy	Fast Greedy-Weighted	Infomap	Infomap-Weighted	Label Propagation	Label Propagation-Weighted	Leading Eigenvector	Leading Eigenvector-Weighted	Walktrap	Walktrap-Weighted
Corel5k	0.000000	0.004000	0.400000	0.148000	0.003556	0.115556	0.004889	0.009778	0.069333	0.000444	0.350667	0.243556
bibtex	0.097778	0.023111	0.022222	0.093333	0.011556	0.044444	0.012444	0.116000	0.059111	0.070667	0.089333	0.089333
birds	0.474341	0.008322	0.459085	0.540915	0.055941	0.019880	0.002774	0.006472	0.044383	0.013870	0.104022	0.154415
delicious	0.000000	0.000000	0.244444	0.350667	0.000000	0.000000	0.000000	0.000000	0.249778	0.387111	0.081778	0.120444
emotions	0.369565	0.391304	0.478261	0.695652	0.695652	0.521739	0.586957	0.891304	0.673913	0.521739	0.608696	0.847826
enron	0.044000	0.483000	0.460000	0.567500	0.284000	0.274500	0.436000	0.522500	0.474500	0.313000	0.151000	0.277000
genbase	0.911111	0.877778	0.856889	0.911111	0.856889	0.911111	0.877778	0.877778	0.877778	0.856889	0.877778	0.877778
mediamill	0.138222	0.256000	0.314667	0.403111	0.329333	0.238667	0.260889	0.276444	0.222222	0.403111	0.345333	0.301778
medical	0.947000	0.517000	0.735000	0.762000	0.636500	0.725000	0.783500	0.529000	0.740500	0.747000	0.585500	0.675500
scene	0.369565	0.521739	0.521739	0.500000	0.282609	0.391304	0.478261	0.456522	0.456522	0.500000	0.630435	0.456522
tmc2007-500	0.999343	1.000000	1.000000	1.000000	1.000000	1.000000	1.000000	1.000000	1.000000	1.000000	1.000000	1.000000
yeast	0.548098	0.478001	0.592095	0.726324	0.599553	0.528710	0.548098	0.615958	0.609247	0.686055	0.527218	0.775541

Weighted fast greedy was the only approach to be more likely to yield a lower Hamming loss than RAKELd on average, both in median and mean (55%) likelihoods. The unweighted version was better than slightly over 50% more of the cases than RAKELd, with a median likelihood of 46%. Binary relevance and label powerset achieved likelihoods lower by close to 10 percentage points. We thus confirm RH2.

When it comes to worst-case observations binary relevance, label powerset and infomap in both variants were never better than RAKELd. The methods with the most resilient advantage in likelihood (9%) in the worst case were weighted versions of fast greedy and walktrap. We confirm RH3 and reject RH4.

We conclude that the fast greedy approach can be recommended over RAKELd, as even given such a large standard deviation (std) of likelihoods, it still yields lower Hamming loss than random label space divisions on more than half of the datasets. Yet, we reject RH1, RH2 and RH4 for Hamming loss. We confirm RH3, as a priori methods do not provide better performance than RAKELd in the worst case.

### 6.6. Efficiency

The efficiency comparison between the data-driven and random method for subspace division is of a special type. It is due to the fact that the data-driven approach requires a single run, whereas the random ones, a number of repetitions to get stable enough results. Thus, comparing the efficiency of both approaches, it must be based on the methods' overall efficiency and not just the computation time of a single run/sampling.

The efficiency of label space partitioning approaches depends on the number of classifiers trained and the complexity of the partitioning procedure. In the case of the flat classification scheme, which we evaluate in this paper, the number of classifiers is equal to the number of subspaces into which the label space is partitioned. As all objects are classified by all sub-classifiers, the classification time is dependent on the single varying parameter: the label subspace size. RAKELd partitioning is performed in  $O(1)$ , as it is a procedure of randomly drawing a partition. The number of classifiers in the data-driven approach depends on the number of detected communities. Table 6 provides information about the number of detected communities in our experiments.

**Table 6.** Number of communities (number of label subspaces) detected in training sets by each method per dataset.

	Fast Greedy	Fast Greedy-Weighted	Infomap	Infomap-Weighted	Label_Propagation	Label_Propagation-Weighted	Leading_Eigenvector	Leading_Eigenvector-Weighted	Walktrap	Walktrap-Weighted
Corel5k	10	13	8	25	4	7	7	15	22	25
bibtex	3	5	2	8	1	7	4	5	9	8
birds	3	3	1	1	1	1	2	3	2	7
delicious	3	6	1	3	1	2	2	6	5	4
emotions	1	2	1	1	1	2	1	2	1	2
enron	4	4	2	2	2	2	3	3	13	7
genbase	10	11	10	11	10	11	10	11	12	11
mediamill	3	3	1	2	1	1	3	2	3	4
medical	20	19	20	20	18	18	20	18	20	20
scene	3	3	2	2	2	2	3	3	3	3
tmc2007-500	2	3	1	1	1	1	2	3	1	2
yeast	1	4	1	1	1	1	1	4	1	4

That is why we evaluate the mean of the proportion between the running time of the data-driven approach to the time of running a single RAKELd, averaged per value of  $k$ . We compare this proportion to the average number of RAKELd runs required to match the average likelihood of data-driven approaches yielding better results of RAKELd per value of  $k$ , as only that number of RAKELd runs brings a high enough probability that enough samples were taken to yield comparable results. We call that point the efficiency threshold. If for example a data-driven method was better than 60% RAKELd samplings for a given  $k$  in a given measure and the total number of samplings for that value of  $k$  was 250, then the efficiency threshold for that value of  $k$  and that measure would be  $60\% \times 250 = 150$ .

Efficiency figures presented in Appendix B (Figures B1–B5), are presented per method as log plots on the  $y$  axis of proportions averaged per  $k$  and thresholds per  $k$  per measure. We had to use log plots on the  $y$  axis due to the fact that the number of RAKELd runs equivalent to a single data-driven run was so greatly below the efficiency thresholds that the charts would not have been readable. Efficiency results do not vary much across methods; thus, we describe them collectively for all methods.

Points below baseline  $y = 1$  are cases when even a single iteration of RAKELd was slower than a single data-driven run. This happens when the number of classifiers in RAKELd is equal to 10, i.e., with  $k = 10\%$  of labels. The smaller the value of  $k$  becomes, the greater the efficiency of the data-driven methods becomes. As  $k$  increases, the number of classifiers trained by RAKELd decreases, and the single random run becomes faster. For most datasets, the data-driven approach running times become equivalent to six to eight RAKELd runs. The the worst case of data-driven running time reaching 15 RAKELd runs happens for large  $k$  on datasets *genbase* and *medical*. For these datasets, with a low level of label co-occurrence, the obtained co-occurrence graph is mostly disconnected, which yields many singletons, and that causes a large increase in the advantage of RAKELd over data-driven methods, as we never allow  $k = 1$  in the RAKELd scenario. What happens is that RAKELd trains few classifiers, and data-driven methods train more single-class classifiers representing singletons in the graph. We plan to improve this in the future by joining singletons into one subspace.

The slowest performance of data-driven methods was equivalent to 15 RAKELd runs. The average efficiency threshold of Hamming loss spanned between 50 and 70 runs of RAKELd, while other measures' thresholds ranged between 80 and 110 runs. We thus note that RAKELd is far less efficient as a method than data-driven approaches, for every evaluated value of  $k$ .

We also note that before measuring efficiency per  $k$ , one needs to perform parameter estimation of the number of runs for RAKELd' in our case, it would take at least ten runs of RAKELd, one for each value of  $k = 10\%, \dots, 90\%$  of labels. One would usually want to repeat the procedure for each parameter value, at least a few times before selecting the value, as the variance in RAKELd performance is large. With data-driven methods, we gain a high likelihood of being better than RAKELd while running only one iteration.

We thus conclude that data-driven approaches are more efficient than RAKELd for every measure evaluated and for every value of  $k$ . We therefore accept RH5.

## 7. Conclusions

We have compared seven approaches as an alternative to random label space partition. RAKELd served as the random baseline for which we have drawn at most 250 distinct label space partitions for at most ten different values of the parameter  $k$  of label subset sizes. Out of the seven methods, five inferred the label space partitioning from training data in the datasets, while the two others were based on an a priori assumption on how to divide the label space. We evaluated these methods on 12 well-established benchmark datasets.

We conclude that in four of five measures, micro-/macro-averaged F1 score, subset accuracy and Jaccard similarity, all of our proposed methods were more likely to yield better scores than RAKELd apart from single outlying datasets; a data-driven approach was better than average random baseline with statistical significance at  $\alpha = 0.05$ . The data-driven approach was also better than RAKELd in worst-case scenarios. Thus, hypotheses RH1, RH3 and RH4 have been successfully confirmed with these measures.

When it comes to Research Hypothesis 2 (RH2), we have confirmed that with micro-averaged F1, subset accuracy, Hamming loss and Jaccard similarity, the data-driven approaches have a higher likelihood of outperforming RAKELd than a priori methods do. The only exception to this is the case of macro-averaged F1, where binary relevance was most likely to beat random approaches, followed closely by a data-driven approach: weighted fast greedy.

Hamming loss forms a separate case for discussion, as this measure is most unrelated to label groups: it is calculated per label only. With this measure, most data-driven methods performed much worse than in other methods. Our study failed to observe a statistical difference between data-driven methods and the random baseline; thus, we reject hypothesis RH1. For best performing data-driven methods, the worst-case likelihood of yielding a lower Hamming loss than RAKELd was close to 10%, which is far from a resilient score; thus, we also reject RH4. We confirm RH2 and RH3, as there existed a data-driven approach that performed better than a priori approaches.

All in all, the statistical significance of a data-driven approach performing better than the averaged random baseline (RH1) has been confirmed for all measures, except Hamming loss. We have confirmed that the data-driven approach was more likely than binary relevance/label powerset to perform better than RAKELd (RH2) in all measures apart from the macro-averaged F1 score, where it followed the best binary relevance closely. Data-driven approaches were always more likely to outperform RAKELd in the worst case than binary relevance/label powerset, confirming RH3 for all measures. Finally, for all measures apart from Hamming loss, data-driven approaches were more likely to outperform RAKELd in the worst case, than otherwise. RH4 is thus confirmed for all measures, except Hamming loss.

In case of measures that are label-decomposable, the fast greedy community detection approach computed on a weighted label co-occurrence graph yielded the best results among data-driven perspectives and is the recommended choice for F1 measures and Hamming loss. When the measure is instance-decomposable and not label-decomposable, such as subset accuracy or Jaccard similarity, the infomap algorithm should be used on an unweighted label co-occurrence graph.

Data-driven methods also prove to be more time efficient than RAKELd. In case of small values of  $k$ , i.e., a large number of models trained by RAKELd, data-driven methods are more efficient than one iteration of RAKELd on most datasets, whereas the number of models decreases; the advantage of speed gained by RAKELd is not reflected in the likelihood of obtaining a result better than data-driven approaches provide in one run.

We conclude that community detection methods offer a viable alternative to both random and a priori assumption-based label space partitioning approaches. We summarize our findings in Table 7, answering the question in the title of how the data-driven approach to label space partitioning is likely to perform better than random choice.

**Table 7.** The summary of the evaluated hypotheses and proposed recommendations of this paper.

Measure	Micro-Averaged F1	Micro Averaged F1	Subset Accuracy	Jaccard Similarity	Hamming Loss
RH1: The data-driven approach is significantly better than random ( $\alpha = 0.05$ )	Yes	Yes	Yes	Yes	No
RH2: The data-driven approach is more likely to outperform RAKELd than a priori methods	Yes	No	Yes	Yes	Yes
RH3: The data-driven approach is more likely to outperform RAKELd than a priori methods in the worst case	Yes	Yes	Yes	Yes	Yes
RH4: The data-driven approach is more likely to perform better than RAKELd in the worst case, than otherwise	Yes	Yes	Yes	Yes	No
RH5: The data-driven approach is more time efficient than RAKELd	Yes	Yes	Yes	Yes	Yes
Recommended data-driven approach	Weighted fast greedy and weighted walktrap	Weighted fast greedy	Unweighted infomap	Unweighted infomap	Weighted fast greedy

**Acknowledgments:** The work was partially supported by the fellowship co-financed by the European Union within the European Social Fund; The European Commission under the 7th Framework Programme, Coordination and Support Action, Grant Agreement Number 316097 (ENGINE); the European Union's Horizon 2020 research and innovation program under the Marie Skłodowska-Curie Grant Agreement No. 691152 (RENOIR); The National Science Centre research project 2014–2017 Decision No. DEC-2013/09/B/ST6/02317.

This work was partially supported by the Faculty of Computer Science and Management, Wrocław University of Science and Technology statutory funds. Kristian Kersting acknowledges the support by the DFG Collaborative Research Center SFB 876 Project A6 “Resource Efficient Analysis of Graphs”. The authors wish to thank Łukasz Augustyniak for helping with proofreading the article.

**Author Contributions:** Piotr Szymański came up with the concept, conducted the experimental code and wrote most of the paper. Tomasz Kajdanowicz helped design the experiments, supported interpreting the results, supervised the work, wrote parts of the paper and corrected and proofread the paper. Kristian Kersting helped reorient the study towards the data-driven perspective, reorganizing and proofreading the paper. All authors have read and approved the final manuscript.

**Conflicts of Interest:** The authors declare no conflict of interest.

## Appendix A. Result Tables

**Table A1.** Number of random samplings from the universum of RAKELd label space partitions for cases different from 250 samples.

Set Name	$k$	Number of Samplings
birds	17	163
emotions	2	15
	3	10
	4	15
	5	6
scene	2	15
	3	10
	4	15
	5	6
tmc2007-500	21	22
yeast	12	91

**Table A2.** Likelihood of performing better than RAKELd in micro-averaged F1 score aggregated over datasets. Bold numbers signify the best likelihoods of a method performing better than RAKELd in the worst-case (Minimum) and average cases (Mean and Median).

	Minimum	Median	Mean	Std
BR	0.500000	0.885556	0.840028	0.152530
LP	0.438280	0.640237	0.704891	0.171726
fast greedy	0.565217	0.806757	0.820198	0.127463
fast greedy-weighted	0.673913	<b>0.922643</b>	<b>0.863821</b>	0.106477
infomap	0.478261	0.713565	0.720074	0.153453
infomap-weighted	0.433657	0.792426	0.748072	0.170349
label_propagation	0.364309	0.734100	0.717083	0.175682
label_propagation-weighted	0.478964	0.815100	0.750085	0.174347
leading_eigenvector	0.630606	0.783227	0.803901	0.116216
leading_eigenvector-weighted	0.630435	0.846506	0.834201	0.107420
walktrap	0.667500	0.742232	0.781920	0.102776
walktrap-weighted	<b>0.695652</b>	0.856861	0.852037	0.091232

**Table A3.** Likelihood of performing better than RAKELd in macro-averaged F1 score aggregated over datasets. Bold numbers signify the best likelihoods of a method performing better than RAKELd in the worst-case (Minimum) and average cases (Mean and Median).

	Minimum	Median	Mean	Std
BR	0.500000	<b>0.985683</b>	<b>0.919245</b>	0.160273
LP	0.543478	0.829283	0.779482	0.163611
fast greedy	0.543478	0.883312	0.866478	0.139572
fast greedy-weighted	<b>0.695652</b>	0.969111	0.900483	0.116022
infomap	0.478261	0.820006	0.793086	0.147108
infomap-weighted	0.500000	0.889556	0.818476	0.164492
label_propagation	0.449376	0.801890	0.754205	0.182183
label_propagation-weighted	0.521739	0.855195	0.821663	0.148561
leading_eigenvector	<b>0.695652</b>	0.863565	0.851696	0.108860
leading_eigenvector-weighted	<b>0.695652</b>	0.889778	0.872119	0.118911
walktrap	<b>0.695652</b>	0.846889	0.844807	0.105977
walktrap-weighted	<b>0.695652</b>	0.894335	0.885253	0.095436

**Table A4.** Likelihood of performing better than RAKELd in subset accuracy aggregated over datasets. Bold numbers signify the best likelihoods of a method performing better than RAKELd in the worst-case (Minimum) and average cases (Mean and Median).

	Minimum	Median	Mean	Std
BR	0.000000	0.498000	0.558057	0.323240
LP	0.336570	0.958889	0.885476	0.190809
fast greedy	0.213130	0.827043	0.791811	0.214756
fast greedy-weighted	0.061951	0.843413	0.747624	0.288200
infomap	<b>0.580500</b>	<b>0.968667</b>	<b>0.910039</b>	0.143954
infomap-weighted	0.525659	0.899430	0.859231	0.152637
label_propagation	0.380028	0.964444	0.900555	0.174854
label_propagation-weighted	0.336570	0.913652	0.861642	0.189690
leading_eigenvector	0.000000	0.826087	0.734935	0.340537
leading_eigenvector-weighted	0.000000	0.843778	0.712555	0.345277
walktrap	0.078132	0.834338	0.739359	0.288072
walktrap-weighted	0.429958	0.812667	0.810542	0.175723

**Table A5.** Likelihood of performing better than RAKELd in Hamming loss aggregated over datasets. Bold numbers signify the best likelihoods of a method performing better than RAKELd in the worst-case (Minimum) and average cases (Mean and Median).

	Minimum	Median	Mean	Std
BR	0.000000	0.369565	0.408252	0.375954
LP	0.000000	0.434653	0.380021	0.338184
fast greedy	0.022222	0.469130	0.507034	0.266736
fast greedy-weighted	<b>0.093333</b>	<b>0.554208</b>	<b>0.558218</b>	0.280209
infomap	0.000000	0.306667	0.396299	0.352963
infomap-weighted	0.000000	0.332902	0.397576	0.345339
label_propagation	0.000000	0.457130	0.415966	0.361819
label_propagation-weighted	0.000000	0.489511	0.441813	0.363469
leading_eigenvector	0.044383	0.465511	0.456441	0.330251
leading_eigenvector-weighted	0.000444	0.451556	0.457361	0.328774
walktrap	0.070667	0.438943	0.444424	0.312845
walktrap-weighted	0.089333	0.379150	0.484974	0.331186

**Table A6.** Likelihood of performing better than RAKELd in Jaccard similarity aggregated over datasets. Bold numbers signify the best likelihoods of a method performing better than RAKELd in the worst-case (Minimum) and average cases (Mean and Median).

	Minimum	Median	Mean	Std
BR	0.456522	0.778060	0.759178	0.202349
LP	0.355987	0.902632	0.854545	0.189086
fast greedy	0.542302	0.876667	0.837570	0.135707
fast greedy-weighted	0.298197	0.875222	0.792386	0.205646
infomap	<b>0.650000</b>	<b>0.945261</b>	<b>0.889663</b>	0.125510
infomap-weighted	0.510865	0.889488	0.855242	0.148578
label_propagation	0.345816	0.928789	0.878622	0.181165
label_propagation-weighted	0.440592	0.920261	0.853821	0.181186
leading_eigenvector	0.163199	0.889874	0.821436	0.222219
leading_eigenvector-weighted	0.464170	0.812444	0.794091	0.154102
walktrap	0.238558	0.891304	0.804866	0.227916
walktrap-weighted	0.644889	0.863722	0.852369	0.122837

**Table A7.** The  $p$ -values of the assessment of the performance of the multi-label learning approaches compared against random baseline by the Iman–Davenport–Friedman multiple comparison, per measure.

	Iman–Davenport $p$ -Value
Subset Accuracy	<b>0.000000004</b>
F1-macro	<b>0.000000000</b>
F1-micro	<b>0.000000177</b>
Hamming Loss	<b>0.0491215784</b>
Jaccard	<b>0.0000124790</b>

**Table A8.** The post-hoc pairwise comparison  $p$ -values of the assessment of the performance of the multi-label learning approaches compared against random baseline by the Iman–Davenport–Friedman test with Rom post hoc procedure, per measure. Bold numbers signify methods that performed statistically significantly better than RAKELd with a  $p$ -value  $< 0.05$ .

	Accuracy	F1-macro	F1-micro	Hamming Loss	Jaccard
BR	0.3590121	<b>0.0000003</b>	<b>0.0000500</b>	1.0000000	<b>0.0064229</b>
LP	<b>0.0000641</b>	<b>0.0280673</b>	<b>0.0205862</b>	1.0000000	<b>0.0000705</b>
fast greedy	0.0515844	<b>0.0000234</b>	<b>0.0000656</b>	1.0000000	<b>0.0038623</b>
fast greedy-weighted	0.1089704	<b>0.0000001</b>	<b>0.0000001</b>	0.3472374	<b>0.0085647</b>
infomap	<b>0.0000257</b>	<b>0.0484159</b>	<b>0.0205862</b>	1.0000000	<b>0.0000705</b>
infomap-weighted	<b>0.0002717</b>	<b>0.0010778</b>	<b>0.0024092</b>	1.0000000	<b>0.0001198</b>
label_propagation	<b>0.0000112</b>	<b>0.0484159</b>	<b>0.0205862</b>	1.0000000	<b>0.0000098</b>
label_propagation weighted	<b>0.0005315</b>	<b>0.0015858</b>	<b>0.0024092</b>	1.0000000	<b>0.0001196</b>
leading_eigenvector	0.1860319	<b>0.0001282</b>	<b>0.0002372</b>	1.0000000	<b>0.0056673</b>
leading_eigenvector-weighted	0.2570154	<b>0.0000274</b>	<b>0.0000653</b>	1.0000000	<b>0.0259068</b>
walktrap	0.0780264	<b>0.0000397</b>	<b>0.0004457</b>	1.0000000	<b>0.0056673</b>
walktrap-weighted	<b>0.0192676</b>	<b>0.0000239</b>	<b>0.0000040</b>	1.0000000	<b>0.0018482</b>

Appendix B. Efficiency Figures

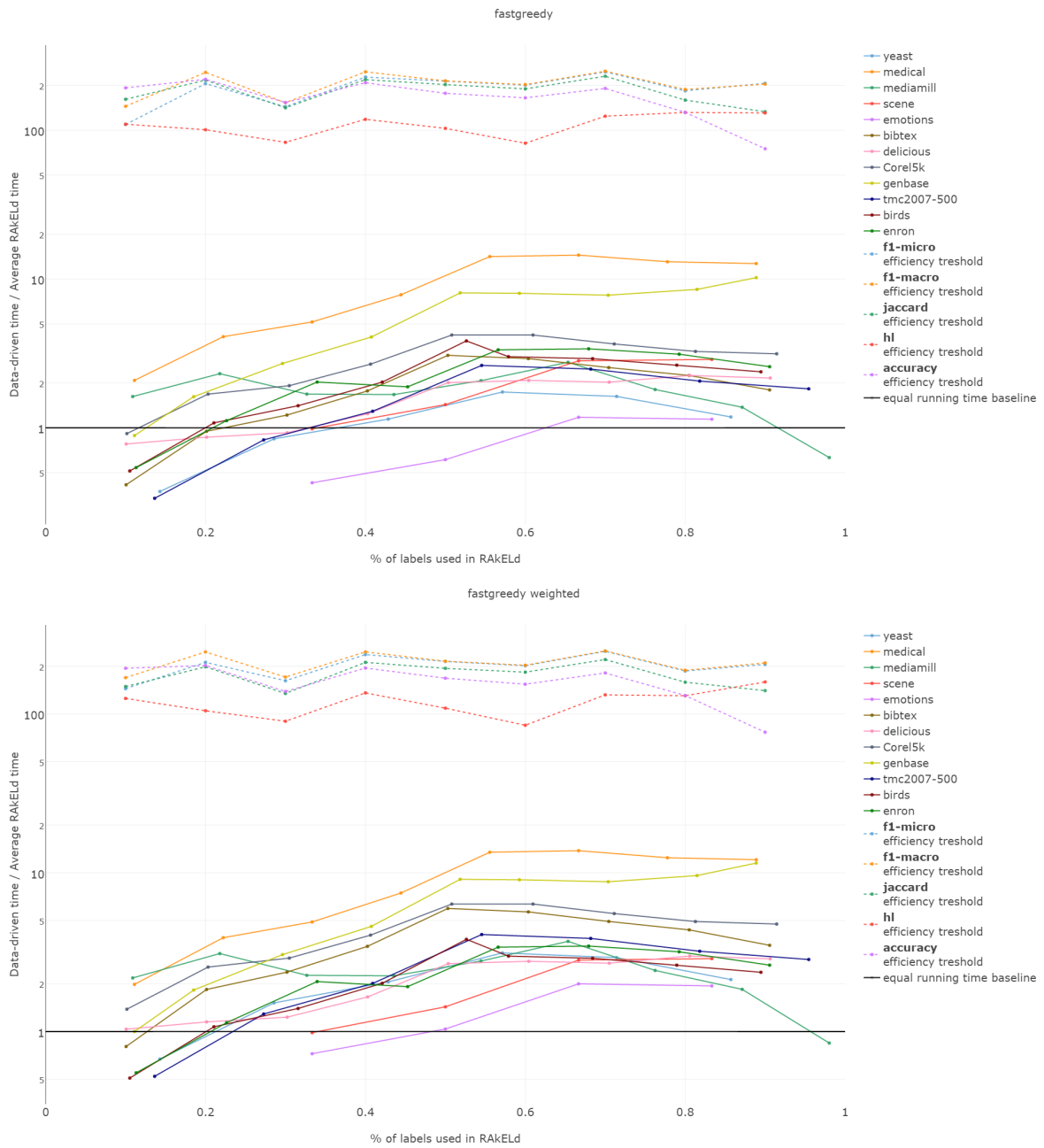


Figure B1. Efficiency of fast greedy modularity maximization data-driven approach against RAKEld.



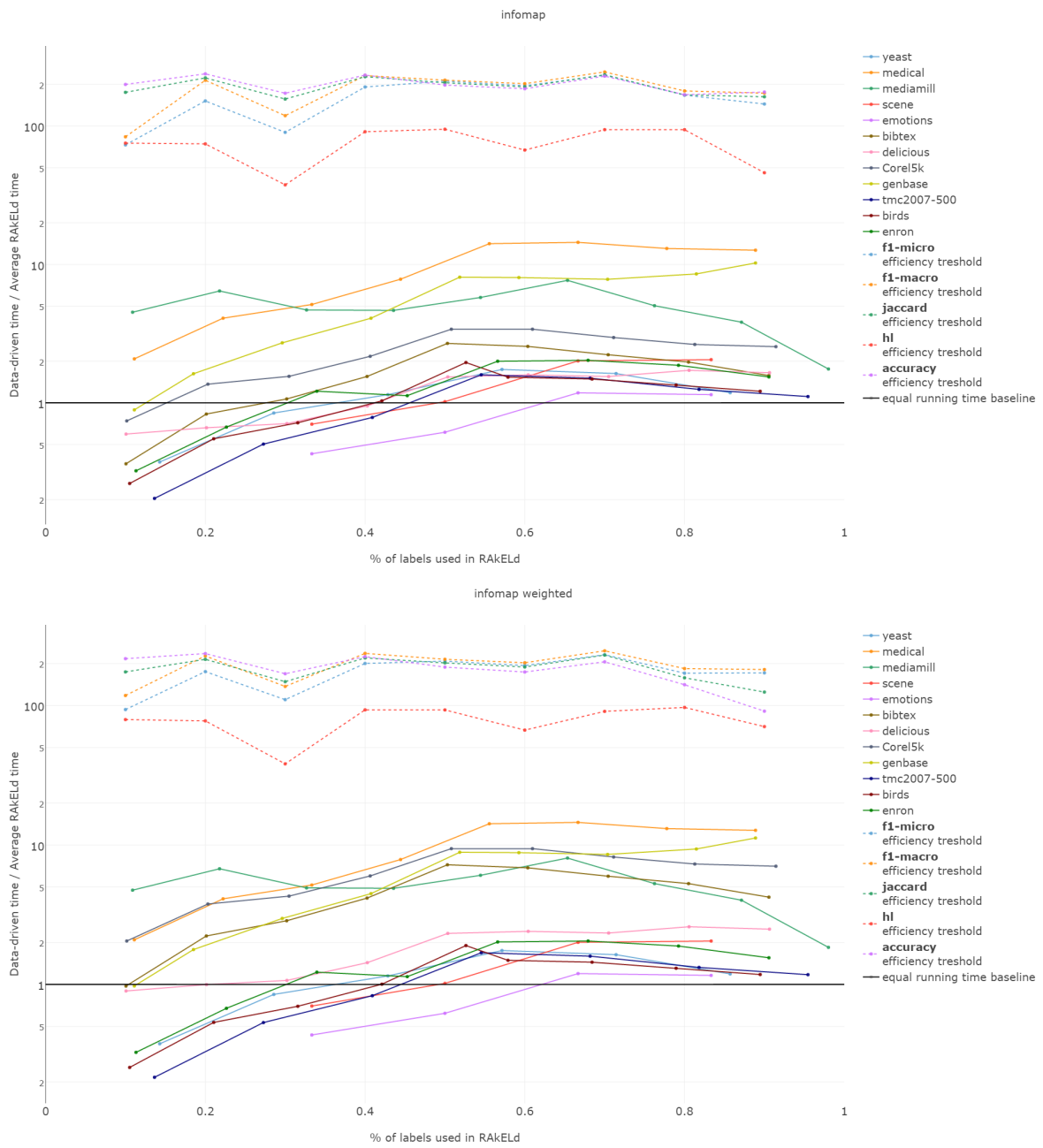


Figure B2. Efficiency of the infomap greedy data-driven approach against RAKEld.

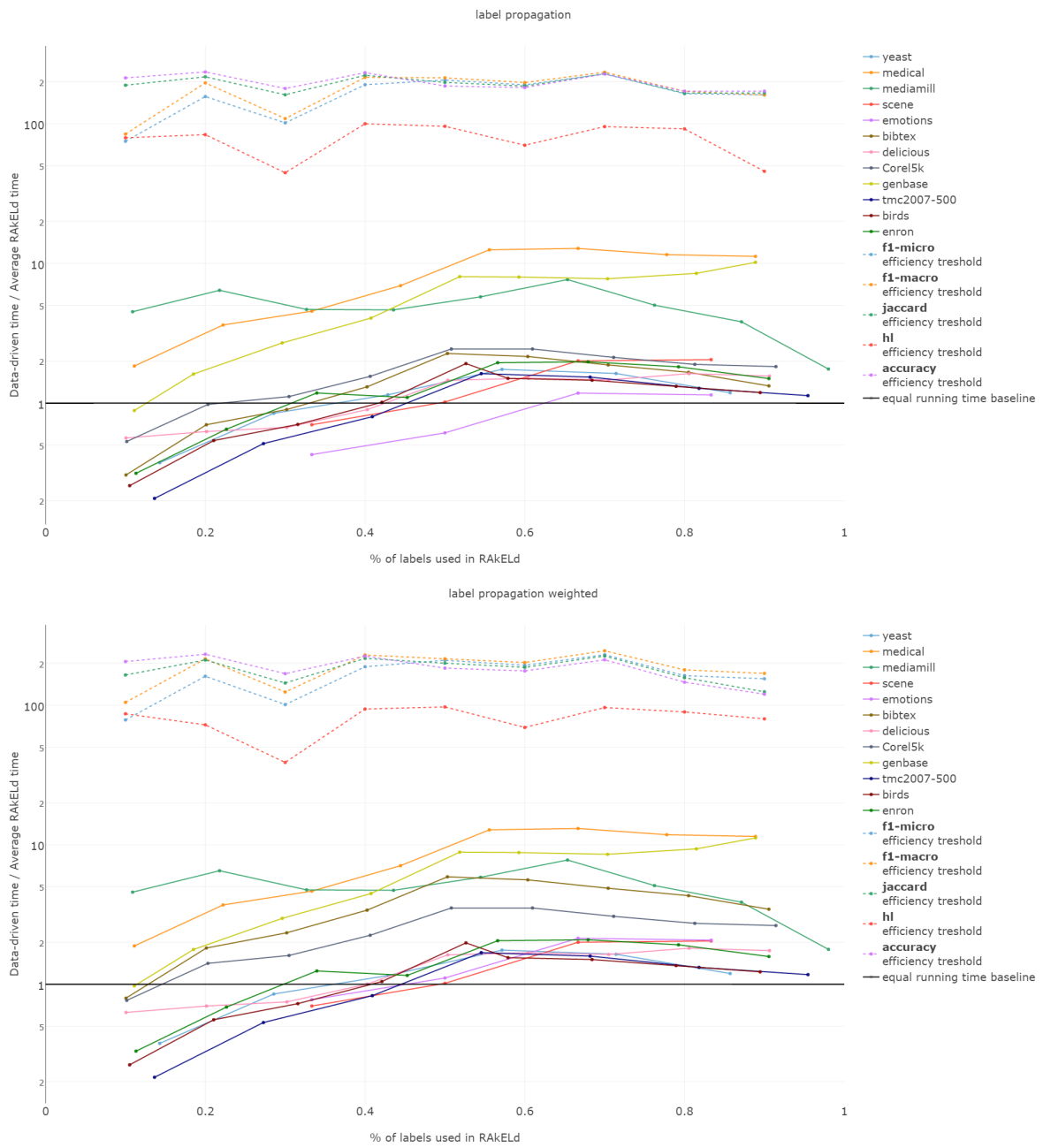
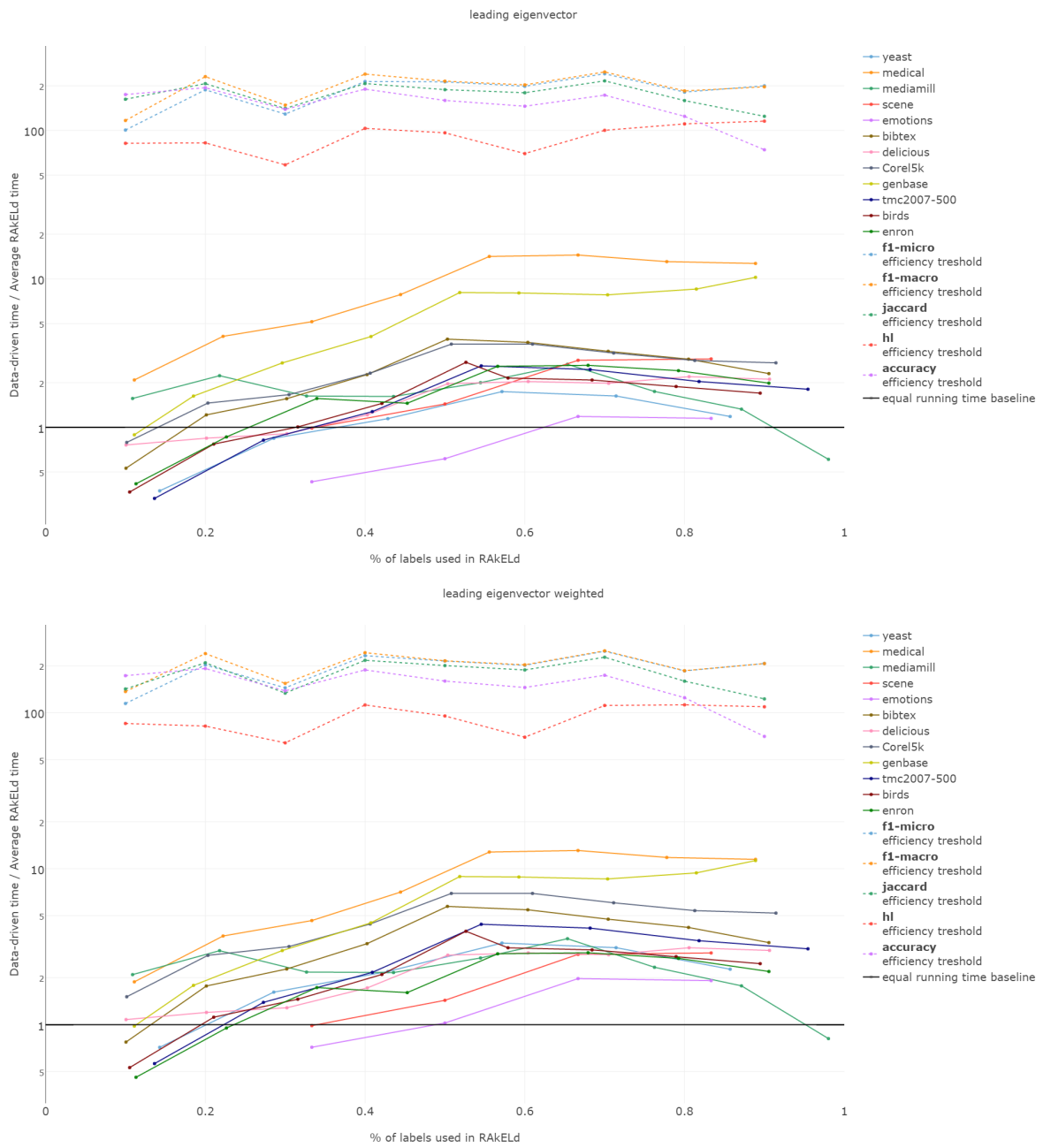


Figure B3. Efficiency of the label propagation data-driven approach against RAKELd.



**Figure B4.** Efficiency of the leading eigenvector modularity maximization data-driven approach against RAKELd.

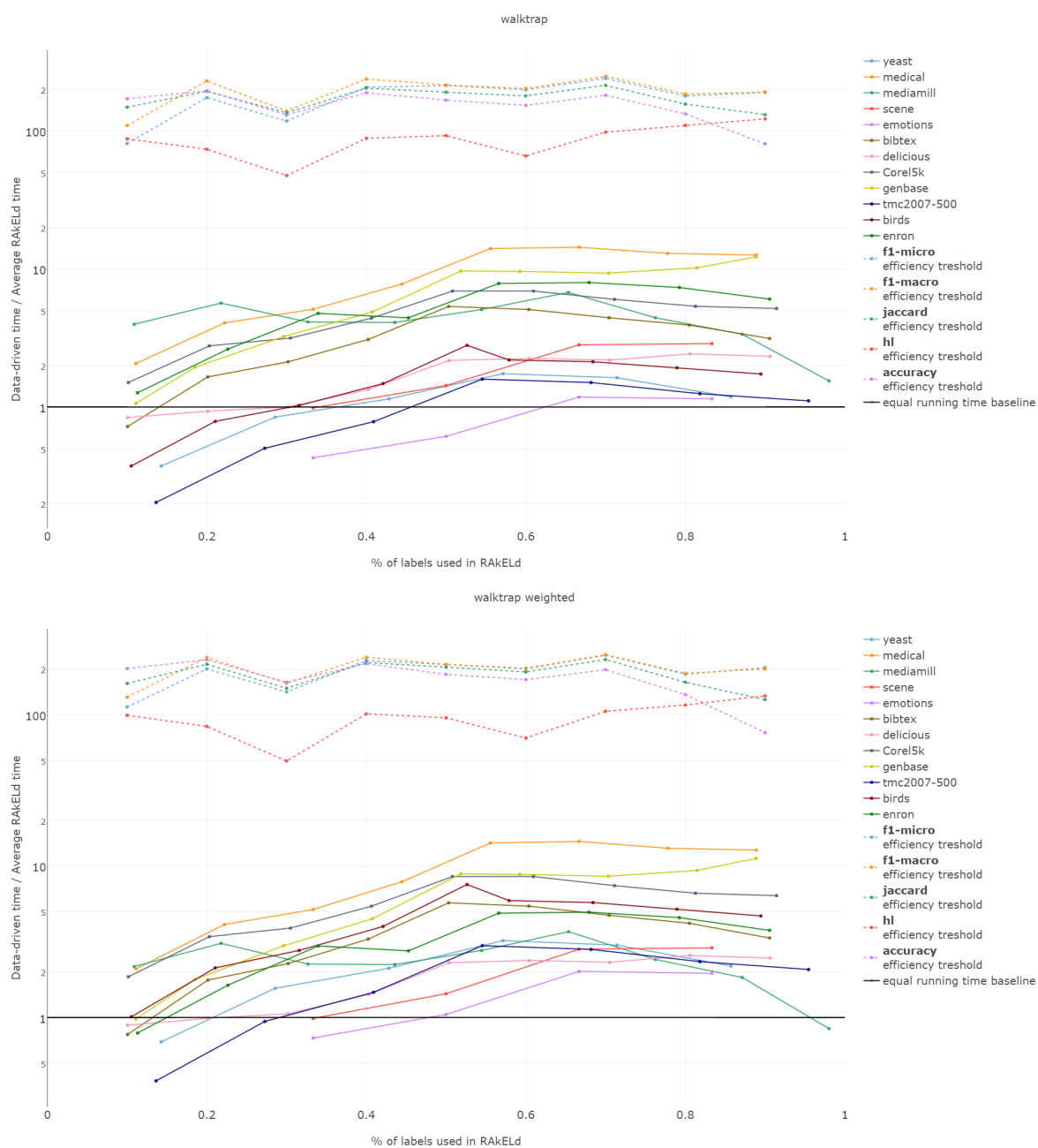


Figure B5. Efficiency of the walktrap data-driven approach against RAKELd.

References

1. Tsoumakas, G.; Katakis, I. Multi-label classification: An overview. *Int. J. Data Warehous. Min.* **2007**, *3*, 1–13.
2. Dembczyński, K.; Waegeman, W.; Cheng, W.; Hüllermeier, E. On label dependence and loss minimization in multi-label classification. *Mach. Learn.* **2012**, *88*, 5–45.
3. Tsoumakas, G.; Katakis, I.; Vlahavas, I. Random k-Labelsets for Multilabel Classification. *IEEE Trans. Knowl. Data Eng.* **2011**, *23*, 1079–1089.
4. Read, J.; Pfahringer, B.; Holmes, G.; Frank, E. Classifier chains for multi-label classification. *Mach. Learn.* **2011**, *85*, 333–359.
5. Dembczynski, K.; Cheng, W.; Hüllermeier, E. Bayes Optimal Multilabel Classification via Probabilistic Classifier Chains. In Proceedings of the 27th International Conference on Machine Learning (ICML-10), Haifa, Israel, 21–24 June 2010; pp. 279–286.

6. Tsoumakas, G.; Katakis, I.; Vlahavas, I. Effective and efficient multilabel classification in domains with large number of labels. In Proceedings of the ECML/PKDD 2008 Workshop on Mining Multidimensional Data (MMD '08), Antwerp, Belgium, 19 September 2008; pp. 30–44.
7. Madjarov, G.; Kocev, D.; Gjorgjevikj, D.; Džeroski, S. An extensive experimental comparison of methods for multi-label learning. *Pattern Recognit.* **2012**, *45*, 3084–3104.
8. Zhang, M.L.; Zhou, Z.H. A review on multi-label learning algorithms. *IEEE Trans. Knowl. Data Eng.* **2014**, *26*, 1819–1837.
9. Clauset, A.; Newman, M.E.J.; Moore, C. Finding community structure in very large networks. *Phys. Rev. E* **2004**, *70*, 066111.
10. Newman, M.E.J.; Girvan, M. Finding and evaluating community structure in networks. *Phys. Rev. E* **2004**, *69*, 026113.
11. Newman, M.E. Analysis of weighted networks. *Phys. Rev. E* **2004**, *70*, 056131.
12. Brandes, U.; Delling, D.; Gaertler, M.; Görke, R.; Hofer, M.; Nikoloski, Z.; Wagner, D. On Modularity Clustering. *IEEE Trans. Knowl. Data E* **2008**, *20*, 172–188.
13. Newman, M.E.J. Finding community structure in networks using the eigenvectors of matrices. *Phys. Rev. E* **2006**, *74*, 036104.
14. Rosvall, M.; Axelsson, D.; Bergstrom, C.T. The map equation. *Eur. Phys. J. Spec. Top.* **2009**, *178*, 13–23.
15. Raghavan, U.N.; Albert, R.; Kumara, S. Near linear time algorithm to detect community structures in large-scale networks. *Phys. Rev. E* **2007**, *76*, 036106.
16. Pons, P.; Latapy, M. Computing communities in large networks using random walks (long version). **2005**, arXiv:physics/0512106.
17. MULAN. Available online: <http://mulan.sourceforge.net/datasets-mlc.html> (accessed on 21 July 2016).
18. Katakis, I.; Tsoumakas, G.; Vlahavas, I. Multilabel text classification for automated tag suggestion. Available online: [http://www.kde.cs.uni-kassel.de/ws/rsdc08/pdf/all\\_rsdv2.pdf#page=83](http://www.kde.cs.uni-kassel.de/ws/rsdc08/pdf/all_rsdv2.pdf#page=83) (accessed on 21 July 2016).
19. Klimt, B.; Yang, Y. The enron corpus: A new dataset for email classification research. In *Machine Learning: ECML 2004*; Springer: Berlin/Heidelberg, Germany, 2004; pp. 217–226.
20. UC Berkeley Enron Email Analysis. Available online: [http://bailando.sims.berkeley.edu/enron\\_email.html](http://bailando.sims.berkeley.edu/enron_email.html) (accessed on 21 July 2016).
21. Computationalmedicine.org. Available online: <http://www.computationalmedicine.org/challenge/> (accessed on 21 July 2016).
22. Boutell, M.R.; Luo, J.; Shen, X.; Brown, C.M. Learning multi-label scene classification. *Pattern Recognit.* **2004**, *37*, 1757–1771.
23. Briggs, F.; Lakshminarayanan, B.; Neal, L.; Fern, X.Z.; Raich, R.; Hadley, S.J.K.; Hadley, A.S.; Betts, M.G. Acoustic classification of multiple simultaneous bird species: A multi-instance multi-label approach. *J. Acoust. Soc. Am.* **2012**, *131*, 4640–4650.
24. Briggs, F.; Huang, Y.; Raich, R.; Eftaxias, K.; Lei, Z.; Cukierski, W.; Hadley, S.; Hadley, A.; Betts, M.; Fern, X.; et al. The 9th annual MLSP competition: New methods for acoustic classification of multiple simultaneous bird species in a noisy environment. In Proceedings of the 2013 IEEE International Workshop on Machine Learning for Signal Processing (MLSP '13), Southampton, UK, 22–25 September 2013; pp. 1–8.
25. Duygulu, P.; Barnard, K.; Freitas, J.F.G.D.; Forsyth, D.A. Object Recognition as Machine Translation: Learning a Lexicon for a Fixed Image Vocabulary. In Proceedings of the 7th European Conference on Computer Vision-Part IV (ECCV '02), Copenhagen, Denmark, 28–31 May 2012; pp. 97–112.
26. Snoek, C.G.M.; Worring, M.; Gemert, J.C.V.; Geusebroek, J.M.; Smeulders, A.W.M. The challenge problem for automated detection of 101 semantic concepts in multimedia. In Proceedings of the ACM International Conference on Multimedia, Santa Barbara, CA, USA, 23–27 October 2006; pp. 421–430.
27. MediaMill, Research on Visual Search. Available online: <http://www.science.uva.nl/research/mediamill/challenge/> (accessed on 21 July 2016).
28. Trohidis, K.; Tsoumakas, G.; Kalliris, G.; Vlahavas, I.P. Multi-Label Classification of Music into Emotions. In Proceedings of the Ninth International Conference on Music Information Retrieval (ISMIR 2008), Philadelphia, PA, USA, 14–18 September 2008; Volume 8, pp. 325–330.
29. Elisseeff, A.; Weston, J. A Kernel Method for Multi-Labelled Classification. In *Advances in Neural Information Processing Systems 14*; MIT Press: Cambridge, MA, USA, 2001; pp. 681–687.

30. Diplaris, S.; Tsoumakas, G.; Mitkas, P.A.; Vlahavas, I. Protein Classification with Multiple Algorithms. In *Advances in Informatics*; Springer: Berlin/Heidelberg, Germany, 2005; pp. 448–456.
31. Derrac, J.; García, S.; Molina, D.; Herrera, F. A practical tutorial on the use of nonparametric statistical tests as a methodology for comparing evolutionary and swarm intelligence algorithms. *Swarm Evol. Comput.* **2011**, *1*, 3–18.
32. Demšar, J. Statistical Comparisons of Classifiers over Multiple Data Sets. *J. Mach. Learn. Res.* **2006**, *7*, 1–30.
33. Scikit-Multilearn. Available online: <http://scikit-multilearn.github.io/> (accessed on 21 July 2016).
34. Pedregosa, F.; Varoquaux, G.; Gramfort, A.; Michel, V.; Thirion, B.; Grisel, O.; Blondel, M.; Prettenhofer, P.; Weiss, R.; Dubourg, V.; et al. Scikit-learn: Machine Learning in Python. *J. Mach. Learn. Res.* **2011**, *12*, 2825–2830.
35. Tsoumakas, G.; Spyromitros-Xioufis, E.; Vilcek, J.; Vlahavas, I. Mulan: A Java Library for Multi-Label Learning. *J. Mach. Learn. Res.* **2011**, *12*, 2411–2414.
36. Csardi, G.; Nepusz, T. The igraph software package for complex network research. *Inter. J. Complex Syst.* **2006**, 1695, 1–9.
37. Yang, Y. An Evaluation of Statistical Approaches to Text Categorization. *Inf. Retr.* **1999**, *1*, 69–90.



© 2016 by the authors; licensee MDPI, Basel, Switzerland. This article is an open access article distributed under the terms and conditions of the Creative Commons Attribution (CC-BY) license (<http://creativecommons.org/licenses/by/4.0/>).