

# Kent Academic Repository

## Full text document (pdf)

### Citation for published version

Xavier-Junior, Joao C. and Freitas, Alex A. and Ludermir, T.B. and Feitosa-Neto, A. and Barreto, C.A.S. (2019) An evolutionary algorithm for automated machine learning focusing on classifier ensembles: an improved algorithm and extended results. *Theoretical Computer Science*. ISSN 0304-3975. (In press)

### DOI

<https://doi.org/10.1016/j.tcs.2019.12.002>

### Link to record in KAR

<https://kar.kent.ac.uk/79330/>

### Document Version

Author's Accepted Manuscript

#### Copyright & reuse

Content in the Kent Academic Repository is made available for research purposes. Unless otherwise stated all content is protected by copyright and in the absence of an open licence (eg Creative Commons), permissions for further reuse of content should be sought from the publisher, author or other copyright holder.

#### Versions of research

The version in the Kent Academic Repository may differ from the final published version.

Users are advised to check <http://kar.kent.ac.uk> for the status of the paper. **Users should always cite the published version of record.**

#### Enquiries

For any further enquiries regarding the licence status of this document, please contact:

[researchsupport@kent.ac.uk](mailto:researchsupport@kent.ac.uk)

If you believe this document infringes copyright then please contact the KAR admin team with the take-down information provided at <http://kar.kent.ac.uk/contact.html>

# An Evolutionary Algorithm for Automated Machine Learning Focusing on Classifier Ensembles: an improved algorithm and extended results

João C. Xavier-Júnior<sup>1</sup>, Alex A. Freitas<sup>4</sup>, Teresa B. Ludermir<sup>3</sup>, Antonino Feitosa-Neto<sup>2</sup>, Cephax A. S. Barreto<sup>2</sup>

---

## Abstract

A large number of classification algorithms have been proposed in the machine learning literature. These algorithms have different pros and cons, and no algorithm is the best for all datasets. Hence, a challenging problem consists of choosing the best classification algorithm with its best hyper-parameter settings for a given input dataset. In the last few years, Automated Machine Learning (Auto-ML) has emerged as a promising approach for tackling this problem, by doing a heuristic search in a large space of candidate classification algorithms and their hyper-parameter settings. In this work we propose an improved version of our previous Evolutionary Algorithm (EA) – more precisely, an Estimation of Distribution Algorithm – for the Auto-ML task of automatically selecting the best classifier ensemble and its best hyper-parameter settings for an input dataset. The new version of this EA was compared against its previous version, as well as against a random forest algorithm (a strong ensemble algorithm) and a version of the well-known Auto-ML method Auto-WEKA adapted to search in the same space of classifier ensembles as the proposed EA. In general, in ex-

---

*Email addresses:* [jcxavier@imd.ufrn.br](mailto:jcxavier@imd.ufrn.br) (João C. Xavier-Júnior), [a.a.freitas@kent.ac.uk](mailto:a.a.freitas@kent.ac.uk) (Alex A. Freitas), [tbl@cin.ufpe.br](mailto:tbl@cin.ufpe.br) (Teresa B. Ludermir), [antonino\\_feitosa@yahoo.com](mailto:antonino_feitosa@yahoo.com) (Antonino Feitosa-Neto), [cephasax@gmail.com](mailto:cephasax@gmail.com) (Cephax A. S. Barreto)

<sup>1</sup>Digital Metropolis Institute - Federal University of Rio Grande do Norte, Natal, Brazil

<sup>2</sup>Department of Informatics and Applied Mathematics - Federal University of Rio Grande do Norte, Natal, Brazil

<sup>3</sup>Center for Information Technology - Federal University of Pernambuco, Recife, Brazil

<sup>4</sup>School of Computing - University of Kent, Canterbury, United Kingdom

periments with 21 datasets, the new EA version obtained the best results among all methods in terms of four popular predictive accuracy measures: error rate, precision, recall and F-measure.

*Keywords:* Automated Machine Learning (Auto-ML), Classification, Evolutionary Algorithms, Estimation of Distribution Algorithms

---

## 1. Introduction

Classification is a very popular Machine Learning task where each instance (object being classified) consists of a set of predictive features and a nominal (or discrete) class variable. In essence, the goal of a classification algorithm is to learn a classification model that can be used to predict the class value (label) of a new instance, based on the values of the features of that instance. Several decades of classification research have produced a large number of classification algorithms [20]. In practice, none of these algorithms is the best for all possible datasets, since the predictive performance of an algorithm is strongly dependent on characteristics of the input dataset [10], as well as on the hyper-parameter settings of the algorithm. This creates the very difficult problem of how to choose the best classification algorithm for the dataset at hand.

A promising and relatively recent approach for tackling this problem is the Automated Machine Learning (Auto-ML) approach [4], [6], [24]. This approach is promising because it uses a search and optimization method to automatically perform a search in a very large space of candidate algorithms and hyper-parameter settings, in order to find the best combination of a classification algorithm and its hyper-parameter settings for the problem at hand. Hence, this approach relieves the user from the task of performing ad-hoc, tedious and very time-consuming experiments with different algorithms and their hyper-parameter settings.

This work follows the Auto-ML approach, focusing on a broad type of classification algorithms called classifier ensembles. An ensemble combines the outputs of many base classifiers (e.g. by majority voting), which tends to improve

25 predictive accuracy by comparison with the use of a single base classifier [1].  
Ensembles are usually considered one of the state-of-the-art types of classifica-  
tion algorithms in terms of predictive accuracy. For instance, a relatively recent  
study [9] compared the predictive accuracy of 179 classification algorithms across  
121 datasets, and concluded that overall the best algorithms were versions of  
30 random forests, which are ensembles of decision-tree classifiers. Even focusing  
on ensembles, however, there are still many different types of ensembles, and  
their predictive accuracies also depend on both the characteristics of the input  
dataset and their hyper-parameter settings.

In this context, the main contribution of this work is to propose an improved  
35 version of our previous Evolutionary Algorithm (EA) [28] for the difficult Auto-  
ML problem of automatically selecting the best ensemble method and its best  
hyper-parameter settings for an input dataset. EAs have the advantages of  
performing a global search in the space of candidate solutions (less likely to get  
trapped into a local optimum than a greedy search) and being robust to noise  
40 [12, 13].

More precisely, the type of EA proposed in this work (as well as its previ-  
ous version) is an Estimation of Distribution Algorithm (EDA) [7]. Unlike  
most EAs, where new candidate solutions are produced by genetic operators  
like crossover and mutation, EDAs evolve a probabilistic model of the best so-  
45 lutions and their components, and at each generation (iteration) they use the  
current probabilistic model to generate new candidate solutions, as discussed  
in Section 2.3. Hence, EDAs combine methods and concepts from both EAs  
and probability theory, which arguably gives them a sounder mathematical ba-  
sis than conventional EAs. In addition, EDAs avoid the need to decide which  
50 genetic operators to use, and so avoid the need for time-consuming experiments  
for optimizing the parameters of genetic operators (like the crossover probability  
and mutation probability).

The proposed improved version of our EDA was compared against its previ-  
ous version [28] and against two other methods, namely a random forest (a  
55 well-known type of ensemble) method and the well-known Auto-ML method

Auto-WEKA [4], which was adapted to search in the same space of classifier ensembles and their configurations (hyper-parameter settings) as the two EDA versions.

This current paper extends the experimental results reported in [28] in three  
60 ways. First, the number of datasets used in the experiments was increased from 15 to 21. Second, whilst in our previous work we reported results for a single predictive accuracy measure, the classification error rate, in the current work we report both the error rate and the values of the precision, recall and F-measure. Third, this current work extends the experiments to include the aforementioned  
65 random forest algorithm.

The results of these new experiments across 21 datasets showed that, overall, the proposed new version of the EDA obtained better predictive performance than its previous version and the other two methods (Auto-WEKA and random forest) for all the four predictive accuracy measures used in our experiments,  
70 namely error rate, precision, recall and F-measure. In addition, the new proposed EDA version obtained statistically significantly better results than Auto-WEKA and random forest in most cases. More precisely, the new EDA was significantly better than Auto-WEKA for three out of the four measures (viz., error rate, recall and F-measure); and the new EDA was also significantly better  
75 than random forest for three measures: precision, recall and F-measure.

The remainder of this paper is organised as follows. Section 2 discusses background on classifier ensembles, Automated Machine Learning (Auto-ML) and EDAs. Section 3 discusses related work on the automated selection of ensemble methods and on the use of EAs for Auto-ML. Section 4 presents the  
80 proposed new version of our EDA, and describes how it differs from its previous version. Section 5 describes the experimental methodology, and Section 6 presents the computational results. Finally, Section 7 presents our conclusions and a direction for future work.

## 2. Background

### 85 2.1. Classification and Classifier Ensembles

In the classification task of machine learning, each instance (object) in the input dataset is represented by a set of features (characteristics) and a class attribute. A classification algorithm has access to the class values of instances in the training set, but not in the test set. Hence, the goal is to learn a model  
90 from the training set that is able to predict the class value of each instance in the test set (with instances unseen during training), based on the feature values for that instance.

Classifier ensembles learn a classification model consisting of a set of base classifiers. Such ensembles have a two-layer structure. In the first layer, each of  
95 the base classifiers receives input data and predicts a class for a new instance. These predictions are sent to a combination module in the second layer, which combines all received predictions into a single predicted class for each instance (e.g. via majority voting). Combining the results of different base classifiers often outperforms a single base classifier [1], [2], since an ensemble's predictions  
100 are usually more robust than the predictions performed by a single classifier.

The two main aspects of the design of classifier ensembles are the selection of the type(s) of classifiers to be used as base classifiers and the combination method. Regarding the choice of classifier types, classifier ensembles can be categorized into homogeneous ensembles, which use multiple base classifiers of  
105 the same type, or heterogeneous ensembles, which use different types of base classifiers. Regarding the type of method used to combine the predictions of the base classifiers, several methods have been proposed [3], such as: simple majority voting, weighted voting (where the vote for a class is weighted by its estimated probability), using a full classification algorithm – treating the  
110 classes predicted by the base classifiers as features for predicting the class at the meta-level, etc.

## 2.2. Automated Machine Learning (Auto-ML)

Many types of classification algorithms have been proposed and are often used, such as Decision Trees, Neural Networks, Support Vector Machines, among  
115 many others [9]. However, in general different types of classification algorithm have different pros and cons, and different biases; therefore no single type of classification algorithm can be considered the best for all datasets or application domains. In practice the predictive accuracy of a classification algorithm strongly depends on two major factors: (a) the characteristics of the input  
120 dataset [10]; and (b) the algorithm’s hyper-parameter settings. This leads to the challenging optimization problem of how to select the best classification algorithm and its corresponding best hyper-parameter settings for each input dataset provided by a user. An emergent approach to solve this problem involves Automated Machine Learning (Auto-ML) methods, which automatically search  
125 for the combination of classification algorithm and hyper-parameter settings that maximizes predictive accuracy in an input dataset.

Recent research on Auto-ML has provided some off-the-shelf Auto-ML tools for machine learning researchers and practitioners; such as Auto-sklearn [6] and Auto-WEKA [4]. Here we briefly review only Auto-WEKA, which is used as a  
130 strong baseline method in our experiments reported later.

Auto-WEKA, which can be easily installed within the well-known WEKA tool [5], is a method for automatically selecting the best combination of a machine learning algorithm and its hyper-parameter settings, as proposed in [4]. Auto-WEKA uses a Bayesian optimization search method called SMAC (Se-  
135 quential Model-based Algorithm Configuration) to automatically search through the joint space of WEKA’s learning algorithms and their respective hyper-parameter settings, with the goal of maximizing predictive accuracy. Auto-WEKA has been shown to perform well for a wide variety of data sets [11].

## 2.3. Estimation of Distribution Algorithms

140 In this subsection we assume the reader is broadly familiar with Evolutionary Algorithms (EAs), and focus on discussing the specific type of EA used as the

basis of the Auto-ML method for ensembles proposed in this work, namely Estimation of Distribution Algorithms (EDAs), as well as discussing the main differences between EDAs and more conventional EAs. For a general discussion  
145 of EAs, the reader is referred to [12, 13].

EDAs are a type of EA which explore the space of potential solutions by building and sampling explicit probabilistic models of promising candidate solutions [7]. EDAs have been applied to several types of machine learning tasks, including classification [17] and feature selection [14], [15], [18].

150 EDAs iteratively generate and evaluate a population of candidate solutions (individuals) to a problem. The initial population is generated at random, using a uniform distribution over all possible candidate solutions. Then, each generated individual has its quality evaluated by a fitness function. Next, individuals are ranked based on their fitness values, and a subset of the best individuals  
155 (usually the 50% best) are selected. Then, a probabilistic model is constructed aiming to estimate the probability distribution of the selected individuals (candidate solutions). Once the model is constructed, new individuals are generated by sampling the distribution encoded by this model. The fitness of each new individual is evaluated, and so on. This process is repeated until some termination  
160 criterion is met, as usual in EAs in general.

The crucial difference between EDAs and other EAs is how they generate new individuals at each generation (iteration), as follows. In EDAs the selected individuals are used to update a probabilistic model, from which new individuals will be probabilistically sampled in the next generation. By contrast, in  
165 other EAs the next generation’s individuals are generated by applying solution-alteration operators like crossover and mutation to the selected individuals of the current generation. EDAs explicitly maintain and evolve a probabilistic model of the best solutions evaluated so far, unlike other EAs. Hence, an advantage of EDAs, by comparison with more conventional EAs, is that EDAs  
170 directly use sound concepts of probability theory to guide the evolutionary process. Another advantage of EDAs is that they require fewer parameters than most EAs. In particular, most EAs require the user to choose which type of



crossover and mutation operators should be used to create new solutions, as well as choosing the corresponding crossover and mutation probability rates. 175 EDAs relieve the user from such concerns, since they do not use any operator to generate new solutions, and simply sample new individuals from the currently available probabilistic model, which is gradually evolved along the search.

The Population-Based Incremental Learning (PBIL) algorithm, proposed in [8] and recently reviewed in [16], is an EDA that evolves a probability vector, where each vector component represents the probability of that component 180 being selected for inclusion in a candidate solution. The vector components' values are usually initialized with a probability of 0.5. Then, at each generation, a population of individuals (candidate solutions) are sampled from the probability vector based on its probability values, and each individual is evaluated using 185 a fitness function, which measures the predictive accuracy of each individual. A predefined number of the best individuals (based on fitness) in the current generation are selected, and the relative frequencies of solution components in those selected individuals are used to update the probability vector, by increasing the probability values for the solution components that occurred most often 190 in the selected individuals. The amount of increase is controlled by a learning rate parameter. More precisely, the probability of each  $i$ -th component of the probability vector at the current  $g$ -th generation (iteration), denoted by  $p[i]$ , is updated with the equation  $p[i] = (1 - LR) \times p[i]_{g-1} + LR \times RF[i]_g$ , where  $LR$  is the learning rate and  $RF$  is the relative frequency of the  $i$ -th component among 195 all individuals selected in the current  $g$ -th generation. Hence, the probability vector gradually evolves towards components with probability values closer to 1 or 0, depending on whether or not, respectively, the component has been used in the best candidate solutions evaluated along the generations.

The PBIL algorithm has only 3 parameters, namely: (a) the population 200 size, i.e., the number of individuals sampled from the probability vector at each generation; (b) the Learning Rate, which specifies how large the steps towards good solutions are; (c) the number of best individuals selected for updating the probability vector at each generation.

### 3. Related Work

205 This section reviews related work on the automated selection of classifier ensembles (Subsection 3.1) and Evolutionary Algorithms (EAs) used for Auto-ML purposes (Subsection 3.2). Note that both these subsections focus on Auto-ML methods. That is, we do not review here conventional methods for learning classifier ensembles without using Auto-ML concepts and methods, since such  
210 conventional ensemble learning methods are already extensively discussed in the literature – see e.g. some relevant reviews in [29], [30], [31].

#### 3.1. Automated Selection of Ensembles Methods

Current Auto-ML methods typically use a search space with many types of classification algorithms, without focusing on ensembles as in this work.  
215 However, some studies have proposed different approaches for automating the creation of classifier ensembles (base classifiers and their hyper-parameter settings) [21], [22], [23]. In particular, Wistuba et al. [21] proposed an automatic approach to generate ensembles with several layers, called Automatic Frankenstein, where a Bayesian Optimization method is used to select base classifiers  
220 and their settings using a bagging strategy.

In fact, most Auto-ML methods are based on Bayesian optimization. For instance, Lacoste et al. [23] proposes an extension of SMBO (Sequential Model-Based Optimization) to optimize the selection of ensemble members based on their performance on randomly selected subsets of the validation data produced  
225 by a bootstrap method. In addition, Levesque et al [22] propose an approach to build fixed-size ensembles, optimizing the configuration of one base classifier of the ensemble at each iteration of the hyper-parameter optimization algorithm, considering the interaction with other models when evaluating performance. In this way, the Bayesian optimization method estimates which prediction model  
230 is the best candidate to be added to the ensemble.

It is important to emphasize that all three aforementioned methods for automating the selection and configuration of classifier ensembles use the Bayesian

optimization method, whereas our proposed method is based on an Estimation of Distribution Algorithm (EDA) – a type of Evolutionary Algorithm (EA). Unlike the sequential nature of the search performed by the Bayesian optimization method, EAs evolve a population of candidate classifier ensembles, performing a more global, broader search (conceptually a parallel search) in the space of candidate solutions.

### 3.2. The Use of Evolutionary Algorithms for Auto-ML

Several Evolutionary Algorithms (EAs) have been proposed for Auto-ML, such as [24], [26], [25], [19]. For example, in [25] the authors proposed the use of a genetic algorithm for searching a very large search space of many different multi-label classification algorithms and their hyper-parameters; whilst in [24] and [26] the authors proposed a genetic programming method for automating the selection and configuration of both classification algorithms and data pre-processing methods (classification pipelines).

On the other hand, in a very recent work [19], the authors proposed the use of a genetic programming (GP) method to search the space of possible architectures of hierarchical ensembles and to optimize their hyper-parameters. Broadly speaking, the GP method proposed in [19] addresses the same type of problem addressed in our work (the automatic creation of ensembles), but using GP, a type of EA that is very different from the EDA proposed in this work – for a brief review of the differences between EDAs and other types of EAs, see Section 2.3. In addition, the work in [19] focuses more on transfer learning and meta-learning, which is not the focus of this current work.

Hence, to the best of our knowledge, our recent work in [28] was the first work to propose an EDA for the Auto-ML task of optimizing the selection and configuration of classifier ensembles. As mentioned in the Introduction, this current work extends our previous work by proposing a new version of that EDA (described in detail in the next section), as well as performing extended computational experiments with more datasets and more classification methods.

#### 4. The Proposed PBIL-Auto-Ens Method for Auto-ML Focusing on Ensembles

As mentioned earlier, the main contribution of this work is to propose an improved version of our previous Estimation of Distribution Algorithm (EDA) for the Auto-ML problem of automatically selecting the best ensemble method and its best configuration (hyper-parameter settings) for an input dataset. Both the proposed new version and the previous version of our EDA are based on the general PBIL algorithm described in Section 2.3.

The first version of our PBIL for the aforementioned Auto-ML problem was proposed in [28], where it was called PBIL-Auto-Ens (PBIL for Auto-ML focusing on Ensembles). That version is hereafter referred to as PBIL-Auto-Ens-v1 (Version 1), since in this current paper we propose the second, improved version of this method, hereafter referred to as PBIL-Auto-Ens-v2 (Version 2). Next, we first focus on describing in detail PBIL-Auto-Ens-v2, and briefly discuss later the main differences between PBIL-Auto-Ens-v2 and PBIL-Auto-Ens-v1.

The proposed PBIL-Auto-Ens-v2 extends the original PBIL algorithm [8] and its more recent variants [16] in two major ways. First, while a standard PBIL typically has a single probability vector, PBIL-Auto-Ens-v2 has many probability vectors (PVs for short), which are organized into a hierarchical structure. Second, the creation of individuals by sampling solution components from the probability vectors is adapted to follow the hierarchical structure of the set of probability vectors, as described later.

As an overview of the proposed PBIL-Auto-Ens-v2, it consists of the following main steps. First, we initialize a number of probability vectors, whose components contain the probabilities of different components of a candidate solution. In this work, these are essentially the probabilities of selecting each type of ensemble method, the probabilities of selecting each parameter setting for each ensemble method, the probabilities of selecting each base classification algorithm within each ensemble method, and the probabilities of selecting each parameter setting for each base classification algorithm. Second, PBIL-Auto-

Ens-v2 creates a population of individuals (candidate solutions) by sampling solution components from those probability vectors, so that each individual consists of a complete specification of an ensemble method, i.e., with all its parameter settings, its base classification algorithm, and the latter's parameter settings. Third, each individual is evaluated according to a fitness function, which measures the predictive accuracy of the ensemble represented by that individual. Fourth, the best (highest fitness) individuals of the current generation are selected, and their solution components are used to update the probability vector. The basic idea is that, if a solution component has been used very often in the selected individuals, the probability of that solution component will be increased in the corresponding probability vector; therefore, that solution component will be more likely to be sampled for creating new individuals in future generations, leading to an improvement of the candidate solutions over time. This iterative process is repeated until a stopping criterion (like a runtime limit) is satisfied. The PBIL-Auto-Ens-v2 method is described in detail in the following sections.

#### *4.1. The Hierarchical Structure of the Set of Probability Vectors*

Let us first describe in detail the hierarchical structure of the set of probability vectors used by PBIL-Auto-Ens-v2, shown in Figure 1. At the first level, there is a PV for selecting the type of ensemble used. This PV has 7 components, representing the probabilities of selecting each of the following ensemble types: Random Committee (RC), AdaBoost (AD), Bagging (BA), Random Forest (RF), Stacking (ST), Vote (VT), and No Ensemble (NE). The latter gives PBIL-Auto-Ens-v2 the option of producing a candidate solution (individual) using only a single base classifier, without any ensemble method.

The second level is the ensembles' hyper-parameter optimization level. At this level, there is in general one PV for each hyperparameter of each of the ensemble methods at level 1. The exceptions are the nodes at level 2 indicating the selection of a base classifier among the corresponding child nodes, where that selection is implemented by PVs at level 3, as discussed below. For instance,

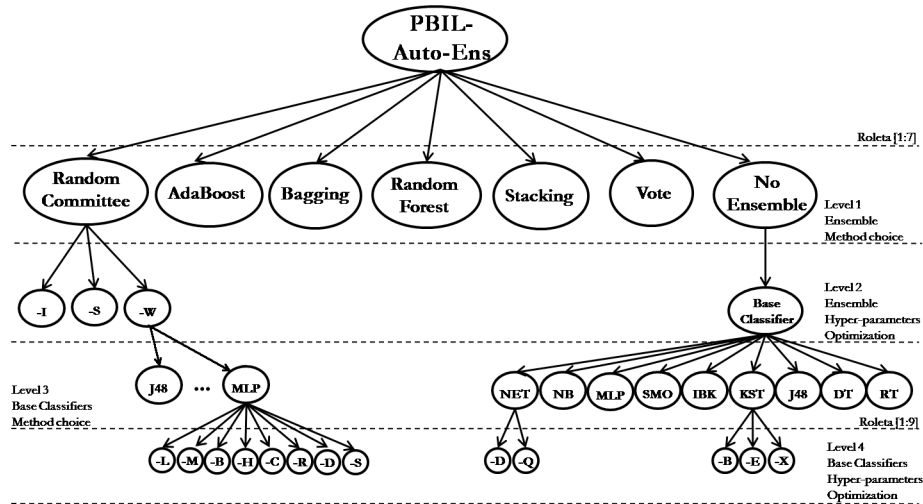


Figure 1: The general structure of PBIL-Auto-Ens-v2’s search space

for the No Ensemble option at level 1, its child node at level 2 is a node called Base Classifier, which is not associated with any PV by itself at level 2, and is placed in the hierarchy only as a bridge between the No Ensemble node at level 1 and the list of classifiers candidate for selection at level 3. Analogously, the node W at level 2 (a child of the RC node) indicates a selection between base classifiers in level 3, and this selection is implemented by a PV at level 3.

Note that in Figure 1 the “...” between the J48 and MLP nodes is a shorthand notation to simplify the figure, referring to all other base classifiers. I.e., the set of base classifiers available for RC (as well as for all other ensemble types) is the same set of 9 base classifiers shown at the right-hand side of level 3 in that figure. Note also that, in order to further simplify Figure 1, this figure shows only the hyper-parameters for the RC ensemble method, but this level 2 also includes hyper-parameters for all other ensemble methods at level 1.

Table 1 shows the full set of PVs at level 2. The number of PVs for each ensemble type (shown in the second column) is also the number of hyper-parameters for that ensemble type – not counting the hyper-parameter W specifying the base classifier, which is associated with a PV at the level 3, as explained

earlier. In the next three columns, the row for each ensemble type is divided  
340 into several sub-rows – one row for each PV, i.e., one row for each of the hyper-  
parameters being optimized. The third column shows the PV name, a string  
of the form L2-XX-Y, where L2 indicates that the current PV is at level 2, XX  
is a two-character variable whose value denotes the type of ensemble method  
(e.g. RC for Random Committee) and Y is a single character variable whose  
345 value identifies the hyper-parameter being optimized (e.g. I for number of it-  
erations). The fourth column shows the number of components of each PV,  
which is the number of candidate discrete values for the corresponding hyper-  
parameter. The last column shows the corresponding candidate discrete values  
for that hyper-parameter.

350 At the third level, there is a PV for selecting the type of base classifier  
used. This leads to 6 PVs at this level, as shown in Table 2, where each row  
describes the characteristics of a PV. In the second column, the PV name is  
a string of the form L3-XX-BC-sel, where L3 indicates that the current PV is  
at level 3 in Figure 1, XX is a two-character variable whose value denotes the  
355 type of ensemble method (e.g. RC for Random Committee, and NE for No  
Ensemble) and BC-sel denotes base classifier selection. Note that all PVs in  
this table have the substrings L3 and BC-sel in their name, since all refer to the  
selection of a base classifier at level 3 of Figure 1. The number of components  
in the PV (in the third column) is the number of candidate base classifiers for  
360 the corresponding type of ensemble. The list of such classifiers is show in the  
last column. Note that all 6 ensemble types in this table use the same set of  
9 candidate base classifiers. Note also that there is no PV for selecting base  
classifiers when the RF ensemble method is used, since RFs always use DTs as  
the base classifier.

365 The fourth level is the base classifiers’ hyper-parameter optimization level.  
At this level, there is again one PV for each hyperparameter of each of the base  
classifiers at the third level. Table 3 shows the full set of PVs at level 4. This  
table has a structure analogous to the one of Table 1. Note that at level 4 of  
Figure 1 each PV is used to optimize the hyper-parameters of a base classifier

Table 1: Probability Vectors (PVs) for ensemble methods hyper-parameter optimization at level 2 of Figure 1. The columns of the table refer to the type of the ensemble at level 1, the number of PVs for the ensemble, the name of the PV, the number of components in the PV, and the values of the hyper-parameters, respectively.

<b>Ens. type L1</b>	<b># PVs Ens.</b>	<b>PV name</b>	<b># Comp. in PV</b>	<b>Hyper-parameter values</b>
RC	2	L2-RC-I	63	from 2 to 64
		L2-RC-S	255	from 1 to 255
AD	4	L2-AD-Q	2	true / false
		L2-AD-P	51	from 50 to 100
		L2-AD-I	127	from 2 to 128
		L2-AD-S	255	from 1 to 255
BA	4	L2-BA-P	91	from 10 to 100
		L2-BA-I	127	from 2 to 128
		L2-BA-S	255	from 1 to 255
		L2-BA-O	2	true / false
RF	3	L2-RF-I	255	from 2 to 256
		L2-RF-K	32	from 1 to 32
		L2-RF-W	20	from 1 to 20
ST	4	L2-ST-X	10	from 1 to 10
		L2-ST-S	255	from 1 to 255
		L2-ST-B	9	from 1 to 9
		L2-ST-NBC	10	from 1 to 10
VT	4	L2-VT-S	255	from 1 to 255
		L2-VT-R	6	AVG, PROD, MAJ, MIN, MAX, MED
		L2-VT-B	9	from 1 to 9
		L2-VT-NBC	10	from 1 to 10



Table 2: Probability Vectors (PVs) for base classifier selection at level 3 of Figure 1. The columns of the table refer to the type of the ensemble at level 1, the number of PVs for the ensemble, the name of the PV, the number of components in the PV, and the candidate base classifiers (one PV component for each of them), respectively.

<b>Ensem. type at Level 1</b>	<b>PV name</b>	<b># Comp. in the PV</b>	<b>Candidate base classifiers</b>
Rand. Comm.	L3-RC-BC-sel	9	NET, NB, MLP, SMO, IBK, KST, J48, DT, RT
AdaBoost	L3-AD-BC-sel	9	NET, NB, MLP, SMO, IBK, KST, J48, DT, RT
Bagging	L3-BA-BC-sel	9	NET, NB, MLP, SMO, IBK, KST, J48, DT, RT
Stacking	L3-ST-BC-sel	9	NET, NB, MLP, SMO, IBK, KST, J48, DT, RT
Vote	L3-VT-BC-sel	9	NET, NB, MLP, SMO, IBK, KST, J48, DT, RT
No ensemble	L3-NE-BC-sel	9	NET, NB, MLP, SMO, IBK, KST, J48, DT, RT

370 regardless of which type of ensemble (at level 2) is using that base classifier.

This approach has the disadvantage of being a relatively coarse-grained approach for hyper-parameter optimization, limiting PBIL’s ability to find fine-grained hyper-parameters settings of a base classifier that would be particularly tailored for a specific type of ensemble method (with its specific hyper-parameter settings). For example, intuitively, the optimal hyperparameter settings for J48  
375 would depend on whether it is being used as a base classifier for AdaBoost or for Stacking (as well as on their hyper-parameter settings).

This coarse-grained hyper-parameter optimization tends to occur more strongly in the early generations of PBIL, when a given base classifier would tend to be  
380 used as part of several different ensemble types in different individuals (since initially all ensemble types have the same probability of being selected to be used in an individual). As generations pass by, the problem should be to some extent mitigated (although not completely eliminated) as the PV for selecting the ensemble type at level 1 is expected to gradually converge to the best en-

Table 3: Probability Vectors (PVs) for the base classifiers' hyper-parameter optimization at level 4 of Figure 1. The columns of the table refer to the type of the base classifier at level 3, the number of PVs for that base classifier, the name of the PV, the number of components in the PV, and the values of the hyper-parameters, respectively.

BC L2	# PVs BCs.	PV name	# Comp. in PV	Hyper-parameter values
DT	4	LV4-DT-E	4	acc, rmse, mae, auc
		LV4-DT-I	2	true / false
		LV4-DT-S	2	BestFirst, GreedyStepWise
		LV4-DT-X	1	from 1 to 4
IBK	5	LV4-IBK-E	2	true / false
		LV4-IBK-K	64	from 1 to 64
		LV4-IBK-X	2	true / false
		LV4-IBK-F	2	true / false
		LV4-IBK-I	2	true / false
J48	8	LV4-J48-O	2	true / false
		LV4-J48-U	2	true / false
		LV4-J48-B	2	true / false
		LV4-J48-J	2	true / false
		LV4-J48-A	2	true / false
		LV4-J48-S	2	true / false
		LV4-J48-M	64	from 1 to 64
		LV4-J48-C	95	from 0,05 to 5,0
KST	3	LV4-KST-B	100	from 1 to 100
		LV4-KST-E	2	true / false
		LV4-KST-X	4	a, d, m n
MLP	8	LV4-MLP-L	10	from 0,1 to 1,0
		LV4-MLP-M	10	from 0,1 to 1,0
		LV4-MLP-B	2	true / false
		LV4-MLP-H	4	a, i, o, t
		LV4-MLP-C	2	true / false
		LV4-MLP-R	2	true / false
		LV4-MLP-D	2	true / false
		LV4-MLP-S	255	from 1 to 255
NB	2	LV4-NB-D	2	true / false
		LV4-NB-K	2	true / false
NET	2	LV4-NET-Q	6	K2, HC, LHC, SA, TS, TAN
		LV4-NET-D	2	true / false
RT	5	LV4-RT-M	64	from 1 to 64
		LV4-RT-K	33	from 0 to 32
		LV4-RT-depth	21	from 0 to 20
		LV4-RT-N	6	from 0 to 5
		LV4-RT-U	2	true / false
SMO	1	LV4-SMO-SEL	4	Def., N.P.Kernel, P.Kernel, Puk, RBFKernel

Table 4: Probability Vectors (PVs) for the SVM base classifier’s hyper-parameter optimization.

<b>SMOs L4, L5</b>	<b>PVs SMOs.</b>	<b>PV name</b>	<b>Comp. in PV</b>	<b>Hyper-parameter values</b>
def.	4	LV5-default-C LV5-default-N LV5-default-M LV5-default-K	11 3 2 4	from 0.5 to 1.5 from 0 to 2 true / false N.P.Kernel, P.Kernel, Puk, RBFKernel
N.P.Kernel	6	LV5-NPKernel-C LV5-NPKernel-N LV5-NPKernel-M LV5-NPKernel-K LV5-NPKernel-E LV5-NPKernel-L	11 3 2 4 49 2	from 0.5 to 1.5 from 0 to 2 true / false N.P.Kernel, P.Kernel, Puk, RBFKernel from 0.2 to 5.0 true / false
PolyKernel	6	LV5-PKernel-C LV5-PKernel-N LV5-PKernel-M LV5-PKernel-K LV5-PKernel-E LV5-PKernel-L	11 3 2 4 49 2	from 0.5 to 1.5 from 0 to 2 true / false N.P.Kernel, P.Kernel, Puk, RBFKernel from 0.2 to 5.0 true / false
Puk	6	LV5-Puk-C LV5-Puk-N LV5-Puk-M LV5-Puk-K LV5-Puk-S LV5-Puk-O	11 3 2 4 100 10	from 0.5 to 1.5 from 0 to 2 true / false N.P.Kernel, P.Kernel, Puk, RBFKernel from 0.1 to 10.0 from 0.1 to 1.0
RBFKernel	5	LV5-RBFKernel-C LV5-RBFKernel-N LV5-RBFKernel-M LV5-RBFKernel-K LV5-RBFKernel-G	11 3 2 4 1000	from 0.5 to 1.5 from 0 to 2 true / false N.P.Kernel, P.Kernel, Puk, RBFKernel from 0.001 to 1.0

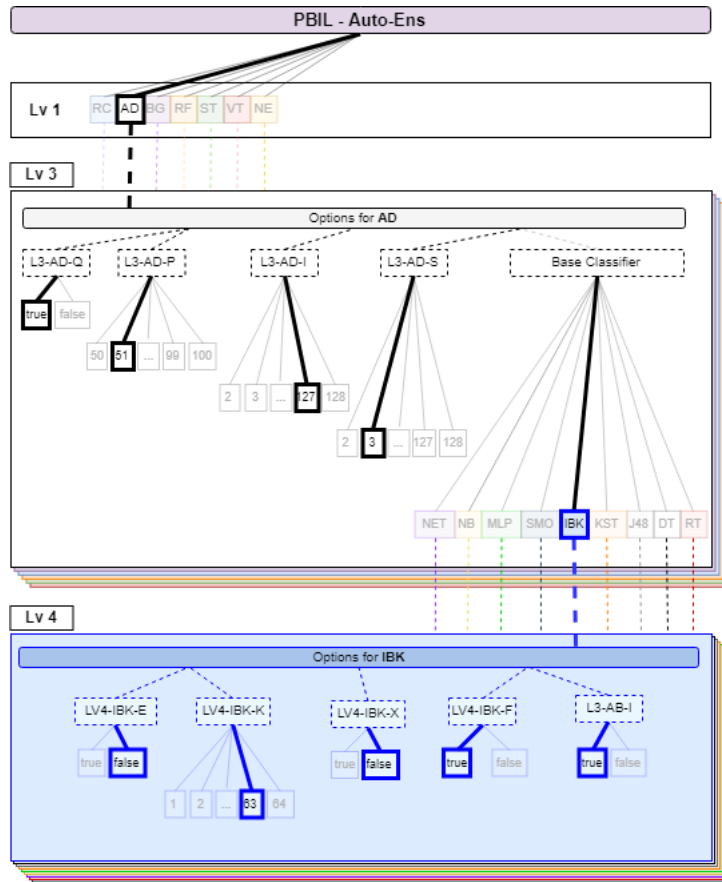


Figure 2: An example of how an individual is generated

semble type. Hence, in late generations a given base classifier should be used  
 385 in different individuals mainly as part of the same best ensemble type, giving  
 more opportunities for the algorithm to focus on a finer-grained optimization  
 of hyper-parameter settings of that base classifier, i.e. finding hyper-parameter  
 settings that are more tailored to that particular best ensemble type.

390 To compensate for the above disadvantage, however, this approach of hav-  
 ing a single PV for each base classifier’s hyper-parameter at level 4 has the  
 important advantage that it drastically reduces the number of PVs that need  
 to be optimized by the PBIL, which should drastically reduce the time for the  
 algorithm to converge to a near-optimal solution. In addition, the much smaller

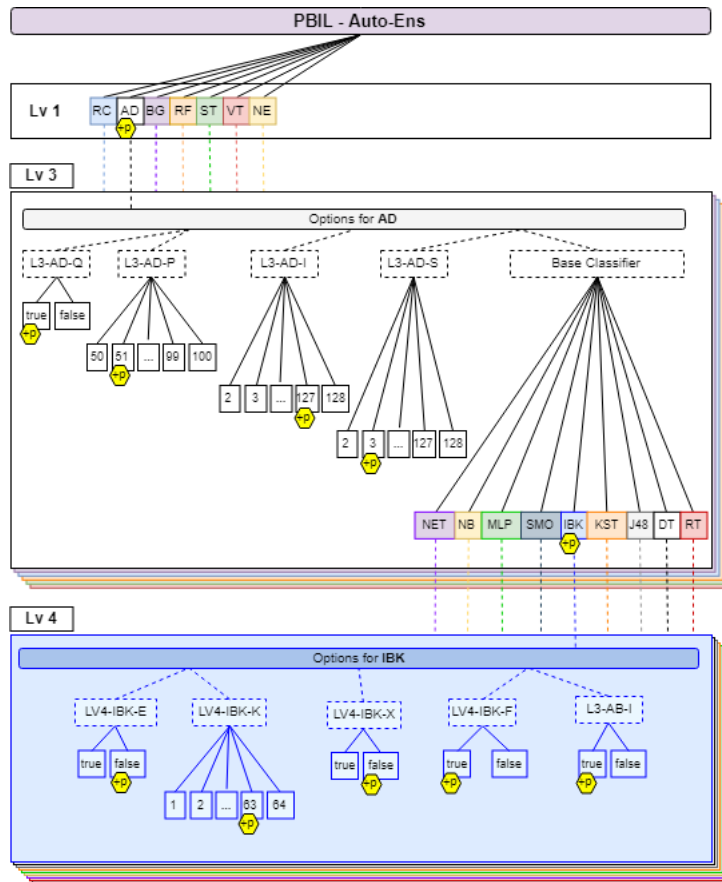


Figure 3: An example of probability vector updates after selecting the individual shown in Figure 2

395 number of PVs associated with this approach can also help to reduce the chances  
of overfitting, particularly on datasets that are not very large.

Finally, there is also a set of PVs referring to hyper-parameters that are  
specific to the base classifier SMO, a type of Support Vector Machine. These  
PVs are specified in Table 4, and they are PVs at the fifth level of the PV  
400 hierarchy. This fifth level was not shown in Figure 1 in order to keep the figure  
relatively simple.

At the start of the evolution (generation 0), all solution components in each  
of the PVs are initialized with a uniform probability distribution, so that each  
component is equally likely to be sampled from each PV. During the evolution-  
405 ary search, at each generation (iteration) individuals are generated by sampling  
solution components from the PVs; then the generated individuals are evalu-  
ated by the fitness function, and the best 50% of the individuals in the current  
generation are selected for updating the probabilities of individual solution com-  
ponents in the PVs. These processes are described in the next subsections.

#### 410 *4.2. The Procedure for Generating Individuals by Sampling from the Probability Vectors*

We now describe how PBIL-Auto-Ens-v2 generates individuals at each gen-  
eration, by sampling solution components from the hierarchy of Probability  
Vectors (PVs) shown in Figure 1.

415 An individual represents a full candidate solution, specifying the choice of  
an ensemble method and its hyper-parameter settings (or No Ensemble if this  
option was chosen), as well as the choice of a base classifier and its hyper-  
parameter settings.

To generate an individual, PBIL-Auto-Ens-v2 first samples a value from the  
420 PV determining the choice of an ensemble method, at level 1 of Figure 1. The  
chosen ensemble method is then used to determine which probability vector will  
be sampled at level 2, and so on, so that the generation of an individual can  
be conceptualized as following a path in the graph representing the hierarchical  
structure of the set of PVs, sampling from one PV at each level, until a full

425 candidate solution is specified.

An example of how an individual is generated is shown in Figure 2. In this example, the individual is generated by the following sequence of PV samplings and solution component choices:

- At level 1: sampling from the PV for ensemble method, choosing AdaBoost  
430 (AD) ensemble (represented by value 2);
- At level 2: sampling from the PVs for hyper-parameters Q, P, I and S of AD, choosing values True, 51, 127 and 3, respectively;
- At level 3: sampling from the PV for hyper-parameter BC-Sel (Base Classifier Selection) of AD, choosing base classifier IBK (represented by value  
435 5);
- At level 4: sampling from the PVs for hyper-parameters E, K, X, F, I of IBK, choosing values False, 63, False, True, True, respectively.

Hence, an individual contains a variable-length list of pairs of the form (i-th PV, i-th Index), where in each pair, the first element denotes the id (name) of a  
440 PV used to generate the individual, and the second element denotes the index of the component of that PV that was sampled; i.e., the index of the method or of the hyper-parameter setting, among the components of the PV.

For instance, the individual generated by the above sequence of PV samplings would be represented by the following list of pairs (we show next only  
445 the first two pairs and the last pair of this list, to simplify):

(L1-EnsType, 2), (L2-AD-Q, 1), . . . , (L4-IBK-I, 1).

where, in the first pair, L1-EnsType is the PV at level 1 for selecting the ensemble method and the value 2 is the index of the component AdaBoost of that PV; in the second pair, L2-AD-Q is the PV at level 2 for selecting the value  
450 of hyper-parameter Q of the AD method and 1 is the index of the value True, which was chosen (sampled) for this hyper-parameter, etc.

### 4.3. Fitness Computation

At each generation, the fitness (quality measure) of each just-created individual is evaluated by applying the classifier ensemble represented by that individual to the training set, using an internal 10-fold cross-validation procedure (a well-known evaluation procedure in machine learning) to estimate the error rate of that ensemble. That is, the fitness of an individual is the mean of the error rates over the 10 folds of the internal cross-validation procedure. We emphasize that fitness is computed using only the training set, without any access to the test set, which is reserved for the final evaluation of the predictive accuracy of the best ensemble returned by the algorithm.

### 4.4. Selection of the Best Individuals and Updating of the Probability Vectors

At the end of each generation, the 50% best individuals in terms of fitness of that generation (i.e. the 50% individuals with the smallest estimated error rates on the training set) are selected and their solution components are used to update the corresponding component probabilities in all the PVs that were used to create the selected individuals. This updating consists of increasing the probability for each solution component in proportion to the relative frequency of use of that component among the selected individuals (the 50% best ones), and also in proportion to the learning rate parameter. More precisely, for each  $i$ -th solution component, its probability at the end of generation  $g$ , denoted by  $p[i]_g$ , is updated using the formula:

$$p[i]_g = (1 - LR) * p[i]_{g-1} + LR * RelFreq[i]_g,$$

where  $LR$  is the learning rate and  $RelFreq[i]_g$  is the relative frequency of the  $i$ -th solution component among the individuals selected at generation  $g$ . Hence, by iteratively selecting the best candidate solutions and increasing the probabilities of each of their components in the probability vector, gradually the probability vector evolves to contain higher probability values for the best solution components, leading to the creation of better and better classifier ensembles. This evolutionary process terminates when a predefined runtime limit is reached.



Figure 3 shows, in diagrammatic form, an example of how the probability vectors are updated after selecting the individual shown in Figure 2. In Figure 3, the small yellow circles with the symbol “+p” identify the probability vectors’ components whose probability values are increased, due to the fact that those components were included in the candidate solution represented by the individual shown in Figure 2.

PBIL-Auto-Ens-v2 differs from its predecessor (PBIL-Auto-Ens-v1) [28] in two major ways. First, PBIL-Auto-Ens-v2 uses a more elaborated hierarchical structure for the set of PVs. In particular, in PBIL-Auto-Ens-v1, at the first level of the hierarchy there is a mixture of classifier ensemble methods and single base classifiers (which directly compete with each other for selection at the first level), whilst in PBIL-Auto-Ens-v2 classifier ensembles and base classifiers are kept in separate hierarchical levels. Hence, there are more modular competitions between the classifier ensembles and their base classifiers, with separate competitions for selection within each of these two groups of classifiers. Second, in PBIL-Auto-Ens-v1 each individual is represented by a variable-length binary vector, with a sub-optimal procedure for decoding those bits into choices of ensembles, base classifiers, and hyper-parameter settings for both types of methods. By contrast, BIL-Auto-Ens-v2 uses a simpler and more natural individual representation, using PV names and index values to directly encode choices of ensembles, base classifiers, and hyper-parameter settings.

#### 4.5. The Space Complexity of PBIL-Auto-Ens-v2

PBIL-Auto-Ens-v2’s space complexity can be calculated by considering the space taken by the set of probability vectors and the space taken by the current population of individuals, as follows. First, the algorithm stores in memory a set of probability vectors, whose total size, denoted PVSize is given by the total number of components in all probability vectors. Second, the space taken by the current population of individuals (candidate solutions) equals the population size (PopSize) times the (average) size of each individual (IndSize). Hence, the space complexity is  $O(\text{PVSize} + \text{PopSize} \times \text{IndSize})$ . To make this formula

more concrete, PVSizE is given by the summation of all values in the columns “Comp. in PV” in Tables 1, 2, 3, 4 plus 9 components for selecting the type of ensemble method at level 1, which is in total 4,261. PopSize is fixed at 50, but IndSize varies across individuals. In practice, IndSize is substantially less  
515 than 50 for any given individual (since each individual needs to store only the chosen parameter settings of chosen algorithms, a small subset of the entire set of PVs), so this space complexity is dominated by PVSizE.

## 5. Experimental Methodology

### 5.1. Datasets Used in the Experiments

520 The proposed PBIL-Auto-Ens method was evaluated on 21 classification datasets, available for download from the well known UCI machine learning repository. Most of these datasets have also been used in recent Auto-ML studies [4], [6], [21]. Table 5 shows the number of instances, attributes (separately for discrete and continuous attributes) and classes in each of the datasets.

### 525 5.2. Methods Compared in the Experiments

The two versions of PBIL-Auto-Ens were compared against two strong baseline methods: random forests (a well-known ensemble method) and an Auto-ML method, namely Auto-WEKA [4]. All experiments were done using a well-known 5-fold cross-validation (CV) procedure. For the Auto-ML methods, i.e. the two  
530 PBIL-Auto-Ens versions and Auto-WEKA, we used a nested version of the CV procedure. More precisely, these Auto-ML methods have been run with an external 5-fold CV, and an internal 10-fold CV. Hence, at the external CV level, the input dataset is randomly divided into 5 folds (each with about 20% of the instances), and then the two versions of PBIL-Auto-Ens and Auto-WEKA are  
535 run 5 times, each using a different fold as the test set and the other 4 folds as the training set. In each of those 5 runs, the training set (80% of the full dataset) is randomly divided into 10 folds, each with about 10% of the training set (i.e., about 8% of the full dataset). Then, whenever a candidate solution (ensemble method or base classifier) is generated by PBIL-Auto-Ens or Auto-WEKA,

Table 5: Main characteristics of the datasets used in the experiments.

<b>Id</b>	<b>Dataset</b>	<b># Instances</b>	<b># Disc. Attr.</b>	<b># Cont. Attr.</b>	<b># Classes</b>
d1	Abalone	4,177	1	7	28
d2	Adult	32,561	8	6	2
d3	Arrhythmia	452	0	260	13
d4	Automobile	205	11	15	7
d5	Car	1,728	0	6	4
d6	Dermatology	366	1	33	6
d7	Ecoli	336	0	7	8
d8	Flags	194	20	10	8
d9	GermanCredit	1,000	13	7	2
d10	Glass Identificaton	214	0	10	7
d11	Image Segmentation	2,310	0	19	7
d12	KR-vs-KP	3,196	36	2	0
d13	Madelon	2,600	0	500	2
d14	Nursery	12,960	8	0	5
d15	Secom	1,567	0	590	2
d16	Semeion	1,593	0	256	10
d17	SolarFlare1	323	13	0	8
d18	Sonar	208	0	60	2
d19	Waveform	5,000	0	40	3
d20	Wine	4,898	0	11	11
d21	Yeast	1,484	0	8	10

540 that solution is evaluated by running its configuration using the internal 10-fold  
CV, so that each of the 10 runs of that candidate solution uses 9 internal folds  
(72% of the full dataset) as a learning set (to learn the classification model)  
and one internal fold (8% of the full dataset) as a validation set to evaluate the  
predictive accuracy of the learned model. The quality measure of that candi-  
545 date solution is given by the mean error rate of the learned model over the 10  
internal validation sets. Hence, the evaluation of each candidate solution uses  
only the training set, not the external test set, which is reserved for measur-  
ing the predictive accuracy of the best solution returned by the two versions of  
PBIL-Auto-Ens and Auto-WEKA.

550 Note that both PBIL-Auto-Ens and Auto-WEKA are non-deterministic search  
methods, i.e. their results depend on a seed number used to randomly initialize  
the candidate solutions. For each method, we report its mean result over exper-  
iments with 5 random seeds (the same seeds are used by all methods), running  
an external 5-fold CV for each seed as explained above – i.e., each reported  
555 result is the mean over 25 results (with 25 different test sets). Random forest is  
also non-deterministic, and it was also run with the same 5 random seeds and  
5-fold cross-validation – i.e. its results are also the mean over 25 test sets.

We used all default parameter settings of Auto-WEKA, including the clas-  
sification error rate (the proportion of incorrectly classified instances) as the  
560 evaluation function to be optimized during training. To make the comparison  
between both PBIL-Auto-Ens versions and Auto-WEKA fair, we also used error  
rate as the fitness function of them. The three Auto-ML methods used the same  
runtime limit, 60 minutes for each run, where one run means one execution of  
one iteration of the external 5-fold CV (applying the method to the training set  
565 of that iteration) for a single value of the random seed, for each dataset. Hence,  
in total each Auto-ML method was run for 525 hours, considering 5 external  
CV iterations times 5 seeds times 21 datasets.

Both PBIL-Auto-Ens versions were run with the following parameter settings  
(as summarized in Table 6): population size of 50, learning rate of 0.5, and  
570 50% of the best individuals of the current generation selected for updating the

Table 6: Configurations for both PBIL-Auto-Ens versions

Population size	Learning rate	% of best individuals selected
50	0.5	50%

probability model. The latter two parameter settings are relatively standard in the PBIL literature, whilst the population size was set based on preliminary experiments. Note that, unlike most PBIL (and EA) implementations, the two PBIL-Auto-Ens versions do not have a parameter for the number of generations, because their stopping criterion is a runtime limit (like in Auto-WEKA). The random forest algorithm was run with its default parameter settings in WEKA.

Both PBIL-Auto-Ens versions have been implemented in Java using the WEKA API. We run the experiments on a desktop PC with Ubuntu 16.04 64 bit operating system driven by an Intel(R) Xeon(R) CPU E5-4610 v4 @ 1.80GHz, 6 core, and RAM with 6 Gb. The program code of PBIL-Auto-Ens-v2 is freely available at: <https://github.com/ml-imd/PBIL-AutoEns-v2>.

### 5.3. Predictive accuracy measures

All methods were evaluated based on 4 predictive accuracy measures, which can be divided into two groups: (a) the classification error rate; and (b) the precision, recall and F-measure. The error rate is simply the ratio of the number of misclassified instances over the total number of instances. The error rate does not distinguish between different types of misclassifications, so it tends to give more importance to the correct classification of instances belonging to the most frequent class in the dataset than to the correct classification of instances belonging to the other class(es).

By contrast, the other 3 measures, used as a whole, tend to evaluate the predictive accuracy in a more balanced way across all classes. Precision is the ratio of the number of correctly classified positive instances over the total number of instances classified as positive (regardless of they belonging to the positive or negative class). Recall is the ratio of the number of correctly classified positive instances over the total number of positive instances (regardless of they

being correctly or wrongly classified). The F-measure is the harmonic average between precision (prec) and recall (rec), defined as:

$$F - measure = \frac{(2 * precision * recall)}{(precision + recall)} \quad (1)$$

Note that the definitions of precision, recall and F-measure consider one  
600 class as the positive class and the other class(es) as the negative class. Hence,  
to compute average values of these measures across classes, we need to consider  
one class at a time as the positive class. More precisely, the average values for  
these measures are computed as follows. For each run of a classification method,  
the average precision and recall were computed by first measuring the precision  
605 and recall per class by considering each class in turn as the positive class, then  
computing the arithmetic mean of those precision and recall values over all  
classes, and finally averaging over all runs of the method for that dataset. The  
average F-measure was computed by applying the above F-measure formula to  
the average values of precision and recall.

## 610 **6. Experimental Results**

This section presents experimental results comparing the predictive performance of the proposed PBIL-Auto-Ens-v2 against three other methods, namely:  
(a) its previous version (PBIL-Auto-Ens-v1); (b) the Auto-WEKA version adapted  
to focus on ensembles, using the same search space as both versions of PBIL-  
615 Auto-Ens; and (c) a random forest algorithm with default hyper-parameter  
settings, as a strong baseline ensemble method. The results are discussed in  
terms of two types of predictive performance measures: the error rate; and the  
precision, recall and F-measure.

### *6.1. Results for the Mean Error Rate*

620 Table 7 presents the mean error rates for Auto-WEKA, PBIL-Auto-Ens-v1  
and PBIL-Auto-Ens-v2. In this table, the best result for each dataset is shown  
in boldface. In addition, for each method, its number of wins (i.e., the number

of datasets where it obtained the best result) and its average rank are shown at the bottom of the table. The lower the average rank of a method, the better its  
625 predictive performance.

As shown at the bottom of Table 7, PBIL-Auto-Ens-v2 obtained the best (smallest) average rank (1.95), with PBIL-Auto-Ens-v1 in the second place (rank 2.19). In addition, PBIL-Auto-Ens-v2 and PBIL-Auto-Ens-v1 achieved the smallest error among all methods in 8 and 7 of the 21 datasets, respectively,  
630 whilst Auto-WEKA and random forest were the winner in only 3 datasets each.

In order to conduct a statistical analysis of the results, the Friedman test and Nemenyi post-hoc test were used (as recommended in [27]) to determine whether or not there is a statistically significant difference between the predictive accuracies of the methods across the 21 datasets. Both tests are applied at the  
635 conventional significance level of 5%. The Friedman test was chosen because it is non-parametric (avoiding the assumption of normality), being based on the average rank of the four methods across all datasets. Its null hypothesis is that there is no difference in the average ranks of the four methods. If this null hypothesis is rejected, we apply the post-hoc Nemenyi test to evaluate if there is  
640 a significant difference between each pair of methods. This is necessary because the Friedman test compares the four methods as a whole, without indicating which pairs of methods have significantly difference performance.

The Friedman test produced the  $p$ -value = 0.0208, therefore the difference between the error rates of the four methods is statistically significant. The  
645 pairwise comparisons using the Nemenyi post-hoc test produced only one statistically significant result: PBIL-Auto-Ens-v2 obtained a significantly better average rank of error rates than Auto-WEKA ( $p$ -value = 0.0420). That is, there is no significant difference between the error rates of other pairs of methods.

### 6.2. Results for the Average Precision, Recall and F-measure

650 Tables 8, 9 and 10 present the average values of precision, recall, and F-measure, respectively, for all methods being compared. As shown at the bottom of Table 8, PBIL-Auto-Ens-v2 obtained the best (smallest) average rank (1.95),

Table 7: Error Rate on the test set (mean over 25 runs for each dataset).

<b>Id</b>	<b>Dataset</b>	<b>Auto-WEKA</b>	<b>Random Forest</b>	<b>PBIL-Auto-Ens-v1</b>	<b>PBIL-Auto-Ens-v2</b>
d1	Abalone	0.7316	0.7616	<b>0.7312</b>	0.7459
d2	Adult	0.1452	0.1504	0.1427	<b>0.1426</b>
d3	Arrhythmia	0.2769	0.3221	<b>0.2743</b>	0.2751
d4	Automobile	0.1912	<b>0.1824</b>	0.1834	0.1842
d5	Car	<b>0.0028</b>	0.0576	0.0102	0.0082
d6	Dermatology	0.0294	0.0339	<b>0.0268</b>	0.0284
d7	Ecoli	0.1387	0.1542	0.1420	<b>0.1363</b>
d8	Flags	0.3352	<b>0.3155</b>	0.3455	0.3537
d9	German-Credit	0.2716	0.2612	0.2700	<b>0.2602</b>
d10	Glass Identification	0.2383	<b>0.2142</b>	0.2392	0.2195
d11	Image Segmentation	<b>0.0199</b>	0.0203	0.0239	0.0231
d12	KR-vs-KP	0.0527	0.0091	0.0154	<b>0.0057</b>
d13	Madelon	0.3009	0.3541	<b>0.2869</b>	0.2991
d14	Nursery	0.0228	0.0121	0.0098	<b>0.0085</b>
d15	Secom	0.0779	0.0664	0.0691	<b>0.0658</b>
d16	Semeion	0.0989	0.0637	<b>0.0588</b>	0.0694
d17	SolarFlare1	0.1182	0.1251	0.1140	<b>0.1126</b>
d18	Sonar	0.2153	0.1788	<b>0.1539</b>	0.1701
d19	Waveform	<b>0.1333</b>	0.1485	0.1476	0.1349
d20	Wine-quality	0.3393	0.3220	<b>0.3210</b>	0.3384
d21	Yeast	0.3980	0.3954	0.3952	<b>0.3950</b>
<b>Number of wins</b>		3/21	3/21	7/21	<b>8/21</b>
<b>Average Rank</b>		2.95	2.86	2.19	<b>1.95</b>

with PBIL-Auto-Ens-v1 in the second place (rank 2.14). In addition, PBIL-Auto-Ens-v2 and PBIL-Auto-Ens-v1 achieved the highest precision among all methods in 8 and 7 of the 21 datasets, respectively; whilst Auto-WEKA and



random forest were the winner in only 4 and 2 datasets, respectively.

In Table 9 (for recall), PBIL-Auto-Ens-v2 obtained the best average rank (1.76), with PBIL-Auto-Ens-v1 in the second place (rank 2.24). In addition, PBIL-Auto-Ens-v2 and PBIL-Auto-Ens-v1 achieved the highest recall among  
660 the three methods in 9 and 6 of the 21 datasets, respectively; whilst Auto-WEKA and random forest were the winner in only 4 and 2 datasets, respectively.

Regarding the F-measure (involving a trade-off between precision and recall), as shown in Table 10, PBIL-Auto-Ens-v2 obtained again the best average rank (1.81), with PBIL-Auto-Ens-v1 again in the second place (rank 2.19). In  
665 addition, PBIL-Auto-Ens-v2 and PBIL-Auto-Ens-v1 achieved the highest F-measure in 9 and 6 of the 21 datasets, respectively; whilst Auto-WEKA and random forest were the winner in only 4 and 2 datasets, respectively.

We also applied the aforementioned Friedman and post-hoc Nemenyi tests to the methods' results for precision, recall and F-measure, based on the methods' average ranks shown in Tables 8, 9 and 10, respectively. For both statistical  
670 tests, we used again the conventional significance level of 5%. For precision, recall and F-Measure, the Friedman test produced the  $p$ -values of 0.0114, 0.0023 and 0.0032, respectively. All these results are statistically significant, so we used the Nemenyi post-hoc test to compare the precision, recall and F-measure for  
675 each of the six pairs of methods.

Regarding the precision measure, the pairwise comparisons using the Nemenyi post-hoc test produced only one statistically significant result: PBIL-Auto-Ens-v2 obtained a significantly better result than the random forest algorithm ( $p$ -value = 0.0210).

680 Regarding the recall measure, the Nemenyi post-hoc test produced two statistically significant results: PBIL-Auto-Ens-v2 was significantly better than both Auto-WEKA ( $p$ -value = 0.0256) and random forest ( $p$ -value = 0.0037).

Regarding the F-measure, the Nemenyi post-hoc test produced again two statistically significant results: PBIL-Auto-Ens-v2 was significantly better than  
685 both Auto-WEKA ( $p$ -value = 0.0304) and random forest ( $p$ -value = 0.0069).

These results can be summarized as follows. PBIL-Auto-Ens-v2 consistently

Table 8: Average Precision on the test set (mean over 25 runs for each dataset).

<b>Id</b>	<b>Dataset</b>	<b>Auto-WEKA</b>	<b>Random Forest</b>	<b>PBIL-Auto-Ens-v1</b>	<b>PBIL-Auto-Ens-v2</b>
d1	Abalone	0.0999	0.1030	0.1075	<b>0.1078</b>
d2	Adult	0.8181	0.8038	<b>0.8208</b>	0.8197
d3	Arrhythmia	0.4407	0.3061	<b>0.4505</b>	0.4467
d4	Automobile	0.6662	0.6114	0.6578	<b>0.6861</b>
d5	Car	<b>0.9944</b>	0.8589	0.9913	0.9925
d6	Dermatology	0.9614	0.9684	<b>0.9730</b>	0.9700
d7	Ecoli	0.6259	0.5502	0.6227	<b>0.6296</b>
d8	Flags	<b>0.4728</b>	0.4412	0.4443	0.4101
d9	GermanCredit	0.6711	0.6911	0.6693	<b>0.6916</b>
d10	Glass	0.7334	0.6950	0.7587	<b>0.7636</b>
d11	Image Seg.	0.9768	<b>0.9801</b>	0.9780	0.9785
d12	KR-vs-KP	0.9494	0.9910	0.9937	<b>0.9944</b>
d13	Madelon	0.7022	0.6482	<b>0.7167</b>	0.7041
d14	Nursery	0.7547	0.7712	0.7772	<b>0.8021</b>
d15	Secom	<b>0.4962</b>	0.4834	0.4668	0.4773
d16	Semeion	0.9056	0.9367	<b>0.9435</b>	0.9339
d17	SolarFlare1	0.1351	0.1373	<b>0.1752</b>	0.1488
d18	Sonar	0.7895	0.8292	<b>0.8555</b>	0.8373
d19	Waveform	0.8581	0.8524	0.8534	<b>0.8658</b>
d20	Wine	0.3141	<b>0.3170</b>	0.3137	0.3092
d21	Yeast	<b>0.5512</b>	0.4751	0.5223	0.4985
<b>Number of wins</b>		4/21	2/21	7/21	<b>8/21</b>
<b>Average Rank</b>		2.81	3.10	2.14	<b>1.95</b>

obtained the best results (in terms of both average rank and number of wins) for all the four measures: error rate, Precision, Recall and F-measure. There was no statistically significant difference between the results of PBIL-Auto-Ens-v2

Table 9: Average Recall on the test set (mean over 25 runs for each dataset).

<b>Id</b>	<b>Dataset</b>	<b>Auto-WEKA</b>	<b>Random Forest</b>	<b>PBIL-Auto-Ens-v1</b>	<b>PBIL-Auto-Ens-v2</b>
d1	Abalone	0.1029	0.1028	0.1089	<b>0.1167</b>
d2	Adult	0.7562	0.7641	<b>0.7770</b>	0.7690
d3	Arrhythmia	0.4202	0.3081	0.4198	<b>0.4206</b>
d4	Automobile	0.6491	0.6090	0.6449	<b>0.6731</b>
d5	Car	<b>0.9891</b>	0.8591	0.9886	0.9888
d6	Dermatology	0.9568	0.9565	<b>0.9717</b>	0.9690
d7	Ecoli	0.6190	0.5483	0.6206	<b>0.6207</b>
d8	Flags	<b>0.4593</b>	0.4472	0.4441	0.4237
d9	GermanCredit	0.6348	0.6452	0.6292	<b>0.6599</b>
d10	Glass	0.7136	0.6931	0.7299	<b>0.7385</b>
d11	Image Seg.	0.9759	<b>0.9799</b>	0.9775	0.9791
d12	KR-vs-KP	0.9461	0.9907	0.9935	<b>0.9942</b>
d13	Madelon	0.6991	0.6470	<b>0.7131</b>	0.7013
d14	Nursery	0.7476	0.7713	0.7796	<b>0.8025</b>
d15	Secom	<b>0.5019</b>	0.4835	0.4997	0.5009
d16	Semeion	0.9006	0.9369	<b>0.9410</b>	0.9305
d17	SolarFlare1	0.1308	0.1370	<b>0.1372</b>	0.1319
d18	Sonar	0.7824	0.8185	<b>0.8423</b>	0.8262
d19	Waveform	0.8558	0.8520	0.8528	<b>0.8654</b>
d20	Wine	0.2284	<b>0.2659</b>	0.2400	0.2425
d21	Yeast	<b>0.5174</b>	0.4775	0.4881	0.4938
<b>Number of wins</b>		4/21	2/21	6/21	<b>9/21</b>
<b>Average Rank</b>		2.86	3.14	2.24	<b>1.76</b>

690 and PBIL-Auto-Ens-v1, for all measures. However, PBIL-Auto-Ens-v2's results were statistically significantly better than the results of both Auto-WEKA and random forest for three out of the four measures of predictive performance used

Table 10: Average F-Measure on the test set (mean over 25 runs for each dataset).

<b>Id</b>	<b>Dataset</b>	<b>Auto-WEKA</b>	<b>Random Forest</b>	<b>PBIL-Auto-Ens-v1</b>	<b>PBIL-Auto-Ens-v2</b>
d1	Abalone	0.1012	0.1039	0.1082	<b>0.1086</b>
d2	Adult	0.7877	0.7835	<b>0.7902</b>	0.7895
d3	Arrhythmia	0.4282	0.3075	0.4324	<b>0.4360</b>
d4	Automobile	0.6592	0.6130	0.6511	<b>0.6772</b>
d5	Car	<b>0.9917</b>	0.8586	0.9899	0.9900
d6	Dermatology	0.9601	0.9623	<b>0.9724</b>	0.9695
d7	Ecoli	0.6219	0.5506	0.6213	<b>0.6247</b>
d8	Flags	<b>0.4646</b>	0.4401	0.4433	0.4160
d9	GermanCredit	0.6523	0.6742	0.6483	<b>0.6748</b>
d10	Glass	0.7215	0.6947	0.7474	<b>0.7506</b>
d11	Image Seg.	0.9765	<b>0.9805</b>	0.9778	0.9789
d12	KR-vs-KP	0.9478	0.9909	0.9938	<b>0.9943</b>
d13	Madelon	0.7006	0.6476	<b>0.7149</b>	0.7051
d14	Nursery	0.7510	0.7722	0.7784	<b>0.8021</b>
d15	Secom	<b>0.5091</b>	0.4831	0.4827	0.4876
d16	Semeion	0.9031	0.9368	<b>0.9423</b>	0.9322
d17	SolarFlare1	0.1315	0.1377	<b>0.1513</b>	0.1370
d18	Sonar	0.7859	0.8237	<b>0.8488</b>	0.8317
d19	Waveform	0.8579	0.8522	0.8531	<b>0.8656</b>
d20	Wine	0.2641	<b>0.2795</b>	0.2767	0.2780
d21	Yeast	<b>0.5332</b>	0.4741	0.5036	0.4925
<b>Number of wins</b>		4/21	2/21	6/21	<b>9/21</b>
<b>Average Rank</b>		2.90	3.10	2.19	<b>1.81</b>

in our experiments. By contrast, although PBIL-Auto-Ens-v1 outperforms both Auto-WEKA and random forest for all those four measures, the difference of results is not statistically significant for any measure. Hence, PBIL-Auto-Ens-v2

represents a clear improvement over PBIL-Auto-Ens-v1, in terms of predictive performance.

Finally, it is worth noting that in general the results for all predictive accuracy measures, for all methods, tend to be worse in datasets with a large number of class labels. The most typical example of this scenario are the results for the Abalone dataset, which has by far the largest number of class labels (28). The results of all methods in this dataset are by far worse than their results in other datasets. To investigate in more detail the relationship between the number of class labels and the predictive accuracy results, we have measured the well-known Pearson’s linear correlation coefficient ( $r$ ) between the number of class labels and the values of each accuracy measure (error rate, precision, recall and F-measure) across the 21 datasets, for each of the four methods being compared in our experiments. These results are shown in Table 11.

As observed in this table, there is a strong positive correlation ( $r$  greater than 0.7) between the number of class labels and the error rate for all methods – i.e., in general larger numbers of class labels are associated with larger (worse) values of the error rate. Conversely, there is a strong negative correlation ( $r$  smaller than -0.6) between the number of class labels and the values of precision, recall and F-measure for all methods – i.e., in general larger numbers of class labels are associated with smaller (worse) values of precision, recall and F-measure. Note also that, for each predictive accuracy measure, there is little variation in the  $r$  values across the methods, i.e. all four methods are equally affected by the difficulty of predicting a large number of class labels.

Table 11: Pearson correlation between Number of Classes and Accuracy Measures.

<b>Pearson</b>	<b>Error</b>	<b>Precision</b>	<b>Recall</b>	<b>F-measure</b>
Auto-WEKA	0.7224	-0.6297	-0.6167	-0.6254
Random Forest	0.7308	-0.6305	-0.6197	-0.6253
PBIL-Auto-Ens-v1	0.7290	-0.6259	-0.6158	-0.6194
PBIL-Auto-Ens-v2	0.7355	-0.6253	-0.6146	-0.6209

### 6.3. An Analysis of the Best Solutions Returned by PBIL-Auto-Ens-v2

720 We now analyse the relative frequency with which different types of classifiers are returned by PBIL-Auto-Ens-v2 as the best solution found during its search. We perform this analysis only for PBIL-Auto-Ens-v2 because it obtained overall the highest predictive accuracy among the four methods, as discussed earlier. Although each returned solution consists of the name of a classifier and its  
725 configuration (i.e., its hyper-parameter settings), we report only the name of returned classifier, which is higher level information, much easier to interpret than the low-level information associated with hyper-parameter settings.

Figure 4 presents bar graphs displaying the 25 best solutions (classifiers) returned by PBIL-Auto-Ens-v2 (i.e., the best solution returned by each run  
730 of this method) for each of the 21 datasets. For each cell (dataset) in this figure, the horizontal axis shows the acronyms of the classifiers selected for that dataset, and the numbers at the top of each bar represent the number of times that the corresponding classifier was selected, out of the 25 runs. Recall that, although the search space of PBIL-Auto-Ens-v2 includes mainly classifier  
735 ensembles (and their configurations), it also includes the option of selecting and configuring only a single base classifier. The latter is a useful option, since in some datasets a single classifier can have a similar or perhaps even somewhat better predictive performance than an ensemble. In addition, a single classifier has the advantage of avoiding the need for the extra computational cost and  
740 complexity of an ensemble.

Overall, considering the results across all datasets in Figure 4, the classification algorithm most frequently selected by PBIL-Auto-Ens-v2 was by far AdaBoost (ADA), selected in 205 cases; followed by Bagging (BAG) and Random Forest (RF), selected in 65 and 60 cases, respectively. Hence, one of these  
745 three ensemble algorithms was selected as the best classifier in 330 cases (about 63% of all 525 cases).

Comparing the selection frequencies of ensembles and single classifiers as a whole, PBIL-Auto-Ens-v2 selected a classifier ensemble in 355 cases (about 68% of the cases), whilst it selected a single base classifier in the remaining 170 cases.

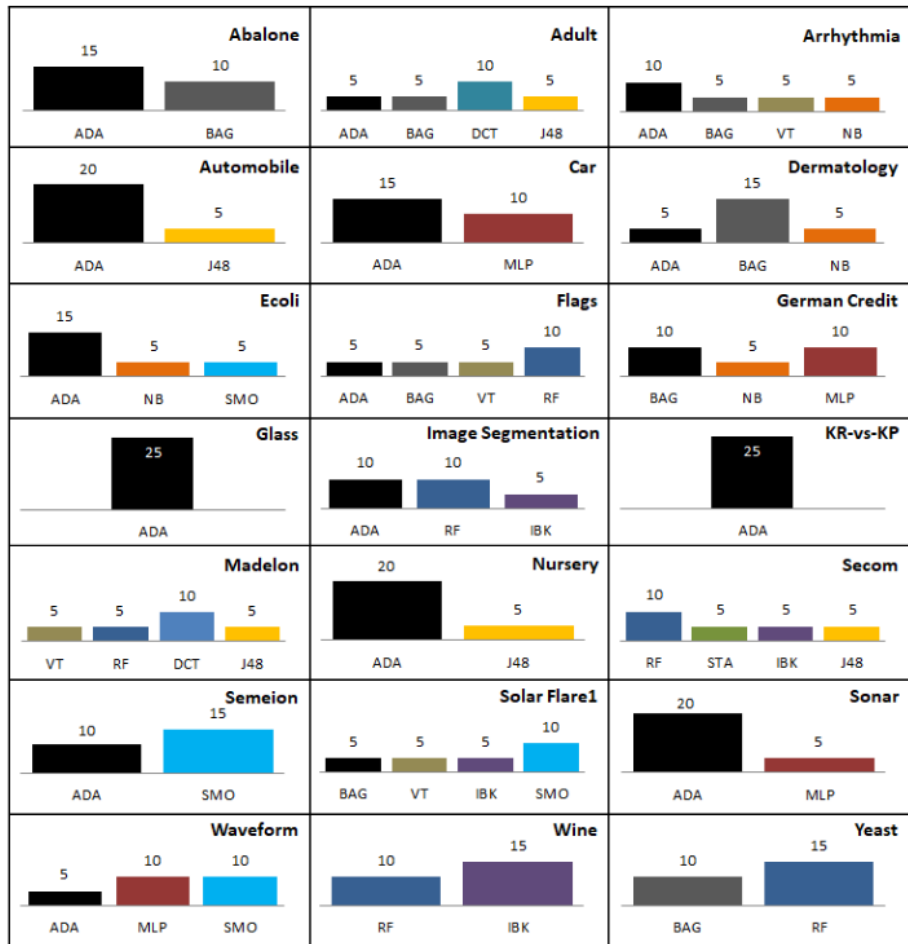


Figure 4: The best algorithms returned by PBIL-Auto-Ens-v2 for each dataset. The classifiers' acronyms: ADA, BAG, RF, STA, VT, NB, MLP, SMO, IBK, DCT and J48 refer to AdaBoost, Bagging, Random Forest, Stacking, Vote, Naive Bayes, Multilayer Perceptron, Support Vector Machine, K-Nearest Neighbors, Decision Table, and Decision Tree, respectively.

750 Moreover, for 5 datasets, PBIL-Auto-Ens-v2 selected an ensemble algorithm  
in all of its 25 runs. This was the case for the Abalone, Flags, Glass, KR-vs-KP,  
and Yeast datasets. By contrast, there was no dataset where PBIL-Auto-Ens-v2  
selected a single base classifier (rather than an ensemble) in all 25 runs.

755 However, for 7 of the 21 datasets, PBIL-Auto-Ens-v2 selected a single base  
classifier in the majority of its runs. This was the case for the Adult, German  
Credit, Madelon, Semeion, Solar Flare1, Waveform and Wine datasets.

Among the single base classifiers returned as best solutions across all datasets,  
the most frequently selected ones were a Support Vector Machine (SMO), Mul-  
tilayer Perceptron (MLP), and a K-Nearest Neighbor algorithm (IBK), which  
760 were selected in 40, 35 and 30 cases, respectively. Overall, there is less variation  
in the selection frequency of single classifiers than in the selection frequency of  
the ensemble algorithms, since the other three single classifiers, namely the J48  
decision tree, Naive Bayes (NB) and Decision Table (DCT), were selected in 25,  
20 and 20 cases, respectively.

765 In summary, despite the overall dominance of three classifier ensembles  
(ADA, BAG, and RF), PBIL-Auto-Ens-v2 is exhibiting great flexibility in se-  
lecting the best classification algorithm for each dataset, which is the core mo-  
tivation for Auto-ML.

## 7. Conclusion and Future Work

770 This work proposed a new version of our Estimation of Distribution Algo-  
rithm (a type of Evolutionary Algorithm), called PBIL-Auto-Ens-v2, for the  
Automated Machine Learning (Auto-ML) problem of automatically selecting  
the best classifier-ensemble method and its best hyper-parameter settings for  
an input dataset. PBIL-Auto-Ens-2 was compared against its previous version  
775 (PBIL-Auto-Ens-1) [28] and against two strong baseline methods: a random  
forest algorithm (a popular type of ensemble) and an adapted version of the  
well-known Auto-WEKA method [4] as an Auto-ML method. Both PBIL-Auto-  
Ens versions and the adapted Auto-WEKA version used the same search space



of candidate solutions (focusing on classifier ensembles) and the same evaluation  
780 function to guide their search. Hence, the differences in their predictive  
accuracies reflect mainly their different search methods, as discussed earlier.

In experiments using 21 classification datasets, overall, the proposed PBIL-  
Auto-Ens-2 obtained better predictive performance than its previous version and  
the other two methods (Auto-WEKA and random forest) for all the four predic-  
785 tive accuracy measures used in our experiments, namely error rate, precision,  
recall and F-measure. In addition, PBIL-Auto-Ens-2 significantly outperformed  
Auto-WEKA for three out of the four measures (viz., error rate, recall and F-  
measure); and PBIL-Auto-Ens-2 also significantly outperformed random forest  
for three measures (precision, recall and F-measure).

790 We also analysed the frequencies with which different types of classifica-  
tion algorithms were chosen as the best algorithm by the overall best Auto-ML  
method in our experiments (i.e., PBIL-Auto-Ens-v2). As the main result of this  
analysis, overall the three most frequently selected algorithms were Adaboost,  
Bagging and Random Forests, with Adaboost being the clear winner overall.

795 As a direction for future work, it would be interesting to extend the exper-  
iments to consider other Auto-ML methods, like the one recently proposed in  
[19]. Another future work direction would be to carry out experiments with  
much larger datasets, using more powerful computational resources.

## Acknowledgments

800 This work has been financially supported by the Brazilian Research Council  
(CNPq) under process number 150748/2017-5 and 407518/2018-5. We thank  
Mr. Douglas Felipe Alves Lima for helping code the PBIL-Auto-Ens-v2 algo-  
rithm.

## References

- 805 [1] Zhi-Hua Zhou, (2012). *Ensemble Methods: Foundations and Algorithms*.  
Chapman & Hall/CRC, 1st edition.

- [2] L.I. Kuncheva, (2004). *Classifier Ensembles for Changing Environments*. In: Roli F., Kittler J., Windeatt T. (eds) Multiple Classifier Systems. Lecture Notes in Computer Science 3077, pages 1–15. Springer.
- 810 [3] L.I. Kuncheva, (2004). *Combining Pattern Classifiers: Methods and Algorithms*. Wiley-Interscience.
- [4] C. Thornton, F. Hutter, H.H. Hoos and K. Leyton-Brown, (2013). *Auto-WEKA: Combined Selection and Hyperparameter Optimization of Classification Algorithms*. Proc. 19th ACM SIGKDD Int. Conf. on Knowledge  
815 Discovery and Data Mining, pages 847–855. ACM Press.
- [5] M. Hall, E. Frank, G. Holmes, B. Pfahringer, P. Reutemann, and I. Witten, (2009). *The WEKA data mining software: an update*. *ACM SIGKDD Explorations Newsletter*, volume 11(1), pages 10–18.
- [6] M. Feurer, A. Klein, K. Eggenberger, J. Springenberg, M. Blum, and  
820 F. Hutter, (2015). *Efficient and Robust Automated Machine Learning*. Advances in Neural Info. Processing Systems 28, pages 2962–2970.
- [7] P. Larrañaga, J.A. Lozano, (2002). *Estimation of Distribution Algorithms: A New Tool for Evolutionary Computation*. Kluwer, Boston, MA.
- [8] S. Baluja and R. Caruana, (1995). *Removing the Genetics from the Standard  
825 Genetic Algorithm*, Proc. 12th Int. Conf. on Machine Learning, pages 38–46, California, July, 1995.
- [9] M. Fernández-Delgado, E. Cernadas, S. Barro and D. Amorim, (2014). *Do We Need Hundreds of Classifiers to Solve Real World Classification Problems?*, *J. of Mach. Learn. Res.*, 15(1), pages 3133–3181.
- 830 [10] P. Brazdil, C. Giraud-Carrier, C. Soares and R. Vilalta, (2009). *Metalearning: Applications to Data Mining*, Springer.
- [11] L. Kotthoff, C. Thornton, H.H. Hoos, F. Hutter and K. Leyton-Brown, (2017). *Auto-WEKA 2.0: Automatic Model Selection and Hyperparameter Optimization in WEKA*, *J. of Mach. Learn. Res.*, 18(1), pages 826–830.

- 835 [12] A.A. Freitas, (2002). *Data Mining and Knowledge Discovery with Evolutionary Algorithms*, Springer.
- [13] A.E. Eiben and J.E. Smith, (2015). *Introduction to Evolutionary Computing*, 2nd edition, Springer.
- [14] I. Inza, P. Larrañaga and B. Sierra, (2002). *Feature subset selection by estimation of distribution algorithms*, Estimation of Distribution Algorithms, pages 269–293, Springer.
- 840 [15] K. Shelke, S. Jayaraman, S. Ghosh, J. Valadi, (2013). *Hybrid feature selection and peptide binding affinity prediction using an EDA based algorithm*, Proc. IEEE Congress on Evolutionary Computation (CEC), pages 2384–  
845 2389.
- [16] M. Zangari, R. Santana, A. Mendibury, A.T.R. Pozo, (2017). *Not all PBILs are the same: unveiling the different learning mechanisms of PBIL variants*, Applied Soft Computing, 53, pages 88–96, April 2017.
- [17] X. Yang, H. Dong and H. Zhang, (2009). *Naive Bayes Based on Estimation of Distribution Algorithms for Classification*, International Conference on  
850 Information Science and Engineering, pages 908-911.
- [18] Y. Saeys, S. Degroeve, D. Aeyels, P. Rouz and Y. Van de Peer, (2004). *Feature selection for splice site prediction: A new method using EDA-based feature ranking*, BMC Bioinformatics, 5(64), 11 pages, 2004.
- 855 [19] P. Kordík, Jan Černý and T. Frýda, (2018). *Discovering Predictive Ensembles for Transfer Learning and Meta-learning*, Machine Learning, 107(1), pages 177–207, January, 2018.
- [20] S. B. Kotsiantis, (2007). *Supervised machine learning: A review of classification techniques*, Emerging Artificial Intelligence Applications in Computer  
860 Engineering, pages 3–24, IOS Press, 2007.

- [21] M. Wistuba, N. Schilling and L. Schmidt-Thieme, (2017). *Automatic Frankenstein: Creating Complex Ensembles Autonomously*, Proc. SIAM Int. Conf. on Data Mining, pages 741–749, SIAM, 2017.
- [22] J. Lévesque, C. Gagné and R. Sabourin, (2016). *Bayesian Hyper-parameter Optimization for Ensemble Learning*, Proc. 32nd Conference on Uncertainty in Artificial Intelligence (UAI), pages 437–446. Jersey City, New Jersey, USA, 2016.
- [23] A. Lacoste, H. Larochelle, F. Laviolette and M. Marchand, (2014). *Sequential Model-Based Ensemble Optimization*, Computing Research Repository (CoRR), 2014.
- [24] R. Olson, R. Urbanowicz, P. Andrews, N. Lavender, L. Kidd, and J.H. Moore, (2016). *Automating Biomedical Data Science Through Tree-Based Pipeline Optimization*, European Conference on the Applications of Evolutionary Computation, pages 123–137. Springer, 2016.
- [25] A.G.C. de Sá, G.L. Pappa and A.A. Freitas, (2018). *Automated selection and configuration of multi-label classification algorithms with grammar-based genetic programming*, To appear in Proc. of the 15th International Conf. on Parallel Problem Solving from Nature (PPSN-2018), to be held in Coimbra, Portugal, Sep. 2018.
- [26] A.G.C. de Sá, W.J.G.S. Pinto, L.O.V.B. Oliveira, G.L. Pappa, (2017). *RECIPE: A Grammar-based Framework for Automatically Evolving Classification Pipelines*, Proc. of the 20th European Conference on Genetic Programming (EuroGP’17), LNCS 10196, pages 246–261. Springer.
- [27] Demsar, Janez, (2006). *Statistical comparisons of classifiers over multiple data sets*, Journal of Machine learning research 7: 1-30. Jan. 2006.
- [28] J. C. Xavier-Junior, A. A. Freitas, A. Feitosa-Neto and T. B. Ludermir(2018). *A Novel Evolutionary Algorithm for Automated Machine Learn-*

*ing Focusing on Classifier Ensembles*, 7th Brazilian Conference on Intelligent Systems (BRACIS 2018 ), pages 462–467, 2018.

- <sup>890</sup> [29] L. Rokach, (2010). *Ensemble-based classifiers*. Artificial Intelligence Review 33(1-2), 1-39, 2010.
- [30] S.B. Kotsiantis, (2014). *Bagging and boosting variants for handling classification problems: a survey*. The Knowledge Engineering Review, 29(1), pages 78–100, 2014.
- <sup>895</sup> [31] O. Sagi and L. Rokach, (2018). *Ensemble learning: a survey*. Wiley Interdisciplinary Reviews: Data Mining and Knowledge Discovery, 8(4), e1249, 2018.