

# Benchmarking Feature-based Algorithm Selection Systems for Black-box Numerical Optimization

Ryoji Tanabe, *Member, IEEE*

**Abstract**—Feature-based algorithm selection aims to automatically find the best one from a portfolio of optimization algorithms on an unseen problem based on its landscape features. Feature-based algorithm selection has recently received attention in the research field of black-box numerical optimization. However, there is still room for analysis of algorithm selection for black-box optimization. Most previous studies have focused only on whether an algorithm selection system can outperform the single-best solver in a portfolio. In addition, a benchmarking methodology for algorithm selection systems has not been well investigated in the literature. In this context, this paper analyzes algorithm selection systems on the 24 noiseless black-box optimization benchmarking functions. First, we demonstrate that the first successful performance measure is more reliable than the expected runtime measure for benchmarking algorithm selection systems. Then, we examine the influence of randomness on the performance of algorithm selection systems. We also show that the performance of algorithm selection systems can be significantly improved by using sequential least squares programming as a pre-solver. We point out that the difficulty of outperforming the single-best solver depends on algorithm portfolios, cross-validation methods, and dimensions. Finally, we demonstrate that the effectiveness of algorithm portfolios depends on various factors. These findings provide fundamental insights for algorithm selection for black-box optimization.

**Index Terms**—Feature-based algorithm selection, black-box numerical optimization, benchmarking

## I. INTRODUCTION

**B**LACK-BOX numerical optimization aims to find a solution  $\mathbf{x} \in \mathbb{R}^n$  with an objective value  $f(\mathbf{x})$  as small as possible without any explicit knowledge of the objective function  $f: \mathbb{R}^n \rightarrow \mathbb{R}$ . Here,  $n$  is the dimension of the solution space. This paper considers only single-objective noiseless black-box optimization. A number of derivative-free black-box optimizers have been proposed in the literature, including mathematical optimization approaches, Bayesian optimization approaches, and evolutionary optimization approaches. In general, the best optimizer depends on the characteristics of a given problem [1]. It is also difficult for a user to select an appropriate optimizer for her/his problem through a trial-and-error process. Thus, automatic algorithm selection is essential for practical black-box optimization.

The algorithm selection problem [2], [3] involves selecting the best one from a portfolio of  $k$  algorithms  $\mathcal{A} = \{a_1, \dots, a_k\}$  on a set of problem instances  $\mathcal{I}$  in terms of a performance

measure  $m: \mathcal{A} \times \mathcal{I} \rightarrow \mathbb{R}$ . The algorithm selection problem is a fundamental research topic in the fields of artificial intelligence and evolutionary computation [4]–[6].

Feature-based offline algorithm selection is one of the most popular approaches for the algorithm selection problem [4]–[6]. First, a feature-based algorithm selection approach computes numerical features of a given problem. It is desirable that the features well capture the characteristics of the problem. Then, the approach predicts the most promising algorithm  $a^{\text{best}}$  from a pre-defined portfolio  $\mathcal{A}$  based on the features. Machine learning techniques are generally employed to build a selection model. Finally,  $a^{\text{best}}$  is applied to the problem. Feature-based algorithm selection has demonstrated its effectiveness on a wide range of problem domains, including the propositional satisfiability problem (SAT) [7], the traveling salesperson problem (TSP) [8], answer set programming (ASP) [9], and multi-objective optimization [10].

Six recent studies [11]–[16] have reported promising results of feature-based algorithm selection for black-box numerical optimization<sup>1</sup>. All of them performed algorithm selection on the 24 noiseless BBOB functions [17]. Throughout this paper, we denote the noiseless BBOB functions as the BBOB functions. Except for [12], these studies used exploratory landscape analysis (ELA) [18] for feature computation, where ELA computes a set of numerical features of a given problem based on a set of solutions. In addition to ELA, the study [13] proposed a feature computation method based on the tree constructed by simultaneous optimistic optimization (SOO) [19]. The study [13] also investigated the effectiveness of the SOO-based features. The results in the previous studies showed that feature-based algorithm selection systems could potentially outperform the single-best solver (SBS) on the BBOB function suite, where SBS is the best optimizer in a portfolio  $\mathcal{A}$  across all function instances. Here, we use the term “an algorithm selection system” to represent a whole system that includes, e.g., a feature computation method, an algorithm selection method, and an algorithm portfolio.

In other words, the previous studies have mainly focused only on whether an algorithm selection system can perform better than the SBS. Although an algorithm selection system consists of many elements, their influence has not been investigated in the literature. A better understanding of algorithm selection systems is needed for the next step. More importantly, a benchmarking methodology has not been well standardized

R. Tanabe is with Faculty of Environment and Information Sciences, Yokohama National University, Yokohama, Japan, and also with Center for Advanced Intelligence Project, RIKEN, Tokyo, Japan. e-mail: (rt.ryoji.tanabe@gmail.com).

<sup>1</sup>We say that a previous study relates to algorithm selection only when it *actually* performed algorithm selection. Although some previous studies identified “performance prediction” with “algorithm selection”, we strictly distinguish the two different tasks.

in the field of algorithm selection for black-box optimization. Table S.7 in the supplementary file shows the experimental settings in the five previous studies [11], [13]–[16], except for [12]. We do not explain Table S.7 due to the paper length limitation, but the experimental settings in the five previous studies are different, including a cross-validation method and an algorithm portfolio. Thus, none of the five previous studies [11], [13]–[16] adopted the same experimental setting.

### Contributions

In this context, this paper analyzes feature-based offline algorithm selection systems for black-box numerical optimization. Through a benchmarking study, this paper addresses the following six research questions.

**RQ1** *Is the expected runtime reliable for benchmarking algorithm selection systems?* The expected runtime (ERT) [20] is a general performance measure for benchmarking black-box optimizers in the fixed-target scenario [21]. Most previous studies also evaluated the performance of algorithm selection systems by using the ERT. However, as discussed in [13], the ERT is sensitive to the maximum number of function evaluations for an unsuccessful run. When different optimizers in a portfolio use different termination conditions, the ERT may incorrectly evaluate the performance of algorithm selection systems.

**RQ2** *How does the performance of algorithm selection systems depend on randomness?* Most operations in algorithm selection systems include randomness, e.g., the generation of the sample and the computation of features. However, the previous studies for black-box optimization [11]–[16] performed only a single run of an algorithm selection system. In [3], Lindauer et al. demonstrated that the performance of the winner of the algorithm selection competition in 2017 significantly depends on a random seed. Thus, it is necessary to understand the influence of randomness on the performance of algorithm selection systems for black-box optimization. The previous studies [16], [22] focused on randomness in the sampling phase. The study [23] pointed out that the VBS performance of a portfolio can be overestimated when the portfolio includes randomized solvers. The study [24] investigated the influence of the runtime variation of randomized SAT solvers on the accuracy of runtime predictors. In contrast, we are interested in randomness in the whole process of algorithm selection.

**RQ3** *How much can a pre-solver improve the performance of an algorithm selection system? Which optimizer is suitable for a pre-solver?* Some modern algorithm selection systems in the field of artificial intelligence (e.g., SATZILLA [7] and 3S [25]) adopt the concept of pre-solving. Here, pre-solving is an approach that aims to solve easy problem instances quickly before the algorithm selection process starts. In contrast, no previous study has used a pre-solver for black-box optimization. As investigated in [14], algorithm selection systems generally perform poorly on easy function instances, e.g.,  $f_1$  (the Sphere function) in the BBOB function set. This issue can potentially be addressed by using a pre-solver.

**RQ4** *Which algorithm selection method is the best?* As reviewed in [26], various algorithm selection methods have been proposed in the field of artificial intelligence. For example, the regression-based method [7] constructs a performance model for each algorithm and selects the best one from  $\mathcal{A}$  in terms of

the predicted performance. However, it is not clear which algorithm selection method is the best for black-box optimization. Although the study [14] evaluated the performance of three selection methods, it did not show details of the comparison results. Unlike [15], we are interested in the performance of algorithm selection methods rather than the performance of machine learning models.

**RQ5** *How difficult is it to outperform the SBS?* The studies [11]–[16] discussed the effectiveness of algorithm selection systems by comparing them with the SBS. They also compared algorithm selection systems with the virtual best solver (VBS), which is an oracle that always selects the best optimizer from  $\mathcal{A}$  on any given problem. A comparison with SBS and VBS allows understanding “how far” an algorithm selection system is from them using, e.g., a relative deviation. If an algorithm selection system  $S$  outperforms the SBS in  $\mathcal{A}$ , the previous studies concluded that  $S$  is effective. However, the difficulty of outperforming the SBS has not been well understood.

**RQ6** *How does the choice of algorithm portfolios influence the overall performance of algorithm selection systems?* The study [14] gave the following rule of thumb to construct a portfolio  $\mathcal{A}$ : “Ideally, the considered set should be as small and as complementary as possible and should include state-of-the-art optimizers”. Since it is difficult to select the best one from too many candidates, the portfolio size should be as small as possible. The performance of the VBS of  $\mathcal{A}$  should also be as good as possible. However, the influence of the VBS and the size of portfolios on the performance of algorithm selection systems is unclear for black-box optimization. It is also unclear which algorithm portfolio should be used in practice.

### Outline

Section II provides some preliminaries. Section III reviews previous studies. Section IV explains our approaches for benchmarking algorithm selection systems. Section V describes our experimental setting. Section VI shows analysis results. Section VII concludes this paper.

### Supplementary file

This paper refers to a figure and a table in the supplementary file as Figure S.\* and Table S.\*, respectively.

### Code availability

The source code used in this study is available at [https://github.com/ryojitnabe/as\\_bbo](https://github.com/ryojitnabe/as_bbo).

## II. PRELIMINARIES

First, Section II-A describes the BBOB function set [17] and the COCO data archive (<https://numbbo.github.io/data-archive>). Then, Section II-B explains the following three performance measures for black-box optimization: the expected runtime (ERT) [20], the relative ERT (relERT) [11], [14], and the successful performance 1 (SP1) [20]. Section II-B also explains other performance measures.

### A. The BBOB function set and the COCO data archive

The (noiseless) BBOB function set [17] consists of the 24 parameterized functions, which are grouped into the following five categories: separable functions ( $f_1, \dots, f_5$ ), functions with low or moderate conditioning ( $f_6, \dots, f_9$ ), functions with

high conditioning and unimodal ( $f_{10}, \dots, f_{14}$ ), multimodal functions with adequate global structure ( $f_{15}, \dots, f_{19}$ ), and multimodal functions with weak global structure ( $f_{20}, \dots, f_{24}$ ). Each BBOB function represents one or more difficulties in real-world black-box optimization. Each BBOB function is instantiated with different parameters.

The BBOB workshop is held at the GECCO conference almost every year. COCO [27] is a platform for benchmarking black-box optimizers. The COCO data archive provides the benchmarking results of almost all optimizers that participated in the BBOB workshop. Currently, except for incomplete results, the benchmarking results of 209 optimizers are available at the COCO data archive. The number of instances for each BBOB function is fixed to 15 for all years. However, as summarized in [14], only the first 5 out of 15 instances are commonly used in all years. For this reason, most previous studies on algorithm selection used only the first five instances whose instance IDs are 1, 2, 3, 4, and 5.

### B. Performance measures for black-box optimization

1) *ERT*: In the context of black-box numerical optimization, the runtime is generally measured in terms of the number of function evaluations rather than the computation time. The ERT [20] measures the expected number of function evaluations needed to reach a target value  $f_{\text{target}} = f(\mathbf{x}^*) + \epsilon$ , where  $\mathbf{x}^*$  is the optimal solution, and  $\epsilon$  is a precision level. See Section V for the  $\epsilon$  value used in this study. Note that most black-box optimizers (e.g., DE [28] and CMA-ES [29]) are invariant in terms of order-preserving transformations of the objective function value [30].

Suppose that an independent run of an optimizer  $a$  is performed for each of instances of a function  $f$ . In this case, the ERT value of  $a$  is calculated as follows:

$$\text{ERT} = \frac{\sum_{i=1}^{N^{\text{run}}} \text{FE}_i}{N^{\text{succ}}}, \quad (1)$$

where  $\text{FE}_i$  is the number of all function evaluations conducted in the  $i$ -th function instance until  $a$  terminates. Here,  $a$  immediately terminates when  $a$  reaches  $f_{\text{target}}$ .  $N^{\text{run}}$  in (1) is the number of runs, where  $N^{\text{run}}$  also represents the number of function instances in this study.  $N^{\text{succ}}$  is the number of successful runs. We say that a run of  $a$  on the  $i$ -th function instance is successful if  $a$  reaches  $f_{\text{target}}$ .

2) *relERT*: The relERT [11], [14] is a normalized ERT using the ERT value of the best optimizer (bestERT) in  $\mathcal{A}$  as follows:  $\text{relERT} = \text{ERT} / \text{bestERT}$ . Here, the best optimizer is determined based on its ERT value for all instances of the corresponding function. The ERT values significantly differ depending on the difficulty of a function. The relERT aims to evaluate the performance of optimizers on the same scale.

The ERT value in (1) and the relERT value are not computable when all runs of  $a$  are unsuccessful (i.e.,  $N^{\text{succ}} = 0$ ). The previous studies [11]–[14] imputed the missing relERT value using the penalized average runtime (PAR10) score [3]. Since PAR10 was used in many previous studies for algorithm selection (e.g., [3], [8], [11], [13], [14], [31]), we adopted PAR10. Similarly, we replaced the missing relERT value with ten times the worst relERT ( $\text{relERT}^{\text{worst}}$ )

value of all optimizers in  $\mathcal{A}$  for each  $n$ . Precisely, for each dimension  $n$ , we defined the  $\text{relERT}^{\text{worst}}$  value based on the relERT values of all algorithms in  $\mathcal{A}$  on all the 24 BBOB functions ( $f_1, \dots, f_{24}$ ) as follows:  $\text{relERT}^{\text{worst}} = \max_{a \in \mathcal{A}, f \in \{f_1, \dots, f_{24}\}} \{\text{relERT}(a, f)\}$ , where  $\text{relERT}(a, f)$  is the relERT value of  $a$  on the  $n$ -dimensional  $f$ .

3) *SP1*: Similar to the ERT, the SP1 [20] estimates the expected number of function evaluations to reach  $f_{\text{target}}$ . The SP1 assumes that the expected number of function evaluations for unsuccessful runs equals that for successful runs [32]. Unlike the ERT, the SP1 is not sensitive to the maximum number of function evaluations. The SP1 is defined as follows:

$$\text{SP1} = \frac{\text{FE}^{\text{avg}}}{p^{\text{succ}}}, \quad (2)$$

where  $\text{FE}^{\text{avg}}$  is the average number of function evaluations for successful runs. In (2),  $p^{\text{succ}}$  is the success probability, which is the number of successful runs  $N^{\text{succ}}$  divided by the number of runs  $N^{\text{run}}$  (i.e.,  $p^{\text{succ}} = N^{\text{succ}} / N^{\text{run}}$ ).

4) *Notes on ERT, relERT, and SP1*: Here, we describe some notes on ERT, relERT, and SP1. The three measures require  $f_{\text{target}}$ , which can depend on  $f(\mathbf{x}^*)$ . Since  $f(\mathbf{x}^*)$  is unknown in most real-world problems, the three measures are not always available. Each optimizer has one or more stopping conditions, e.g., the maximum budget of evaluations  $b^{\text{max}}$ . Since the ERT in (1) takes into account the number of function evaluations used in an unsuccessful run, the ERT is sensitive to a stopping condition of an optimizer. Section 3.3.2 in [13] shows how sensitive the ERT is to  $b^{\text{max}}$  using an intuitive example. The same is true for the relERT. An optimizer can possibly reach  $f_{\text{target}}$  when setting  $b^{\text{max}}$  to a sufficiently large number, and vice versa. In other words, the results of an optimizer depend on  $b^{\text{max}}$ . Thus, any performance measure (including the SP1) is influenced by  $b^{\text{max}}$  even when considering the same optimizer.

5) *Other performance measures*: The ERT, relERT, and SP1 measures are for the fixed-target scenario, which is representative in the BBOB community. The performance for the fixed-budget scenario is generally evaluated by the quality of the best-so-far solution [33]. The study [34] proposed a hybrid measure of the ERT and the error value.

Some previous studies (e.g., [35], [36]) proposed anytime performance measures. Roughly speaking, this kind of measure commonly aims to evaluate the anytime performance of an optimizer by calculating the “volume” of its performance profile. For example, the IOHprofiler platform [37] provides an anytime performance measure based on the area under the curve of the empirical cumulative distribution function.

## III. LITERATURE REVIEW

This section reviews previous studies on algorithm selection for black-box optimization problems as well as other problems. First, Section III-A describes features for algorithm selection only for black-box numerical optimization. Then, Section III-B introduces methods for constructing algorithm portfolios. Section III-C describes selection methods. Section III-D explains cross-validation methods. Section III-E describes pre-solvers. Finally, Section III-F describes related work in other domains.

## A. Features

Feature-based algorithm selection systems require a set of domain-dependent features, which represent the characteristics of a given problem. It is challenging to design helpful features for black-box optimization. As discussed in [12], unlike other problems (e.g., SAT), only scarce information about a problem is available from its definition. For this reason, features need to be computed based on a set of  $s$  solutions  $\mathcal{X} = \{\mathbf{x}_i\}_{i=1}^s$  and their objective values  $f(\mathcal{X})$ , where the  $s$  solutions can be generated by any method, e.g., random sampling and local search. However, the evaluation of solutions by the objective function  $f$  is generally computationally expensive. Thus, features should be computed based on a small-size  $\mathcal{X}$ .

The ELA approach [18] is generally used to compute features for black-box optimization. ELA computes a set of numerical features from  $\mathcal{X}$  and  $f(\mathcal{X})$ . Most previous studies used the R-package `flacco` [38] to compute ELA features. Table S.8 shows 17 feature classes currently provided by `flacco`. Each feature class consists of more than one feature. In total, 342 features are available in `flacco`. However, three feature classes (`ela_conv`, `ela_curv`, and `ela_local`) have not been used in most recent studies since they need additional function evaluations apart from  $\mathcal{X}$ . Since the five cell mapping feature classes (e.g., `cm_angle`) are computable only for  $n \leq 5$ , they have not been generally used.

Each feature class characterizes the properties of a problem in a different way. For example, the `ela_meta` feature class builds multiple regression models based on  $\mathcal{X}$  and  $f(\mathcal{X})$ . The model-fitting results are the `ela_meta` features, e.g., how well the regression models can fit the given data set. The `disp` features are computed by the dispersion metric [39], which was designed for quantifying the degree of the global structure.

## B. Construction methods of algorithm portfolios

Since no optimizer can perform the best on all function instances, it is necessary to select the most promising optimizer from a portfolio on a given instance. We explain three methods [6], [11], [40] for constructing a portfolio of  $k$  algorithms  $\mathcal{A} = \{a_1, \dots, a_k\}$  for algorithm selection. We do not describe methods for constructing *parallel algorithm portfolio approaches* (e.g., [41]). The three construction methods aim to construct a portfolio whose optimizer can solve at least one BBOB function. Even when a portfolio includes “a correct answer”, algorithm selection systems cannot know it due to the exclusive property of the cross-validation methods (see Section III-D). In real-world applications, an algorithm selection system selects the most promising algorithm from a portfolio on a given problem. HYDRA [42], [43] is an efficient automatic configurator for algorithm portfolios for the combinatorial domain. Although HYDRA requires a parameterized system, how to construct it in the continuous domain is unclear. An extension of HYDRA to the continuous domain is also beyond the scope of this paper. Section VI investigates the influence of  $k$  on the performance of algorithm selection systems later.

In [11], first, a set of candidates  $\mathcal{R}$  are selected so that each algorithm in  $\mathcal{R}$  performs the best on at least one BBOB function in terms of the ERT. Then,  $k$  algorithms in  $\mathcal{A}$  are

further selected from  $\mathcal{R}$  so that  $\mathcal{A}$  minimizes the worst-case performance of the VBS on the 24 BBOB functions. Here, the study [11] did not explain how to get  $\mathcal{A}$  from  $\mathcal{R}$ . A method proposed in [40] first constructs  $\mathcal{R}$  as in [11]. Then, the method iteratively selects an algorithm from  $\mathcal{R}$  based on a voting strategy. Here, the method uses the ERT as a performance measure. A method proposed in [6] constructs a portfolio on the BBOB functions with  $n \in \{2, 3, 5, 10\}$ . First, all candidates are ranked based on their ERT values. Then, four candidate sets  $\mathcal{R}_2$ ,  $\mathcal{R}_3$ ,  $\mathcal{R}_5$ , and  $\mathcal{R}_{10}$  are constructed for 2, 3, 5, and 10 dimensions, respectively.  $\mathcal{R}_n$  contains algorithms ranked within the top 3 of at least one  $n$ -dimensional BBOB function. Finally, the method selects algorithms that commonly belong to the four sets, i.e.,  $\mathcal{A} = \bigcap_{n \in \{2, 3, 5, 10\}} \mathcal{R}_n$ .

## C. Algorithm selection methods

An algorithm selection method aims to find a mapping from features to  $k$  algorithms in a portfolio  $\mathcal{A} = \{a_1, \dots, a_k\}$  by machine learning. As reviewed in [26], various algorithm selection methods have been proposed in the field of artificial intelligence. In [26], Lindauer et al. proposed AUTOFOLIO, which is a highly-parameterized algorithm selection framework for combinatorial optimization problems, including SAT and ASP. They generalized existing algorithm selection methods for AUTOFOLIO. Inspired by [26], this paper considers the following five general algorithm selection methods:

1) *Classification*: The classification-based method was used in LLAMA [44] and previous studies for black-box optimization [11], [14]. The classification-based method builds a classification model to directly predict the best algorithm from the  $k$  algorithms in  $\mathcal{A}$  based on a given feature set, where the best algorithm is determined for each function based on a given performance measure. As pointed out in [14], the classification-based method does not consider the performance rankings of the other  $k - 1$  algorithms.

2) *Regression*: The regression-based method was used in SATZILLA’09 [7] and recent studies for black-box optimization [13]–[15]. First, the regression-based method constructs  $k$  regression models for  $k$  algorithms in  $\mathcal{A}$ , respectively. Then, the regression-based method selects the best one from  $k$  algorithms based on their predicted performance.

3) *Pairwise classification*: The pairwise classification-based method was used in SATZILLA’11 [43]. The study [26] reported the promising performance of the pairwise classification-based method for combinatorial optimization. A similar selection method was also adopted by a study for black-box optimization [16]. In the training phase, the pairwise classification-based method builds a classification model for each pair of  $k$  algorithms in  $\mathcal{A}$ . In the testing phase, the method evaluates all  $\binom{k}{2}$  models. Then, the method selects the best one out of  $k$  algorithms in terms of the number of votes. In this study, ties are broken randomly.

4) *Pairwise regression*: Although the pairwise regression-based method was originally proposed for the TSP [8], it has been used for black-box optimization [14]. First, the method constructs a regression model for each pair of  $k$  algorithms in  $\mathcal{A}$ , where the model predicts the performance difference

between two algorithms for each pair. Then, the method selects the best one from  $k$  algorithms based on the sum of predicted performance differences.

5) *Clustering*: The clustering-based method was used in ISAC [45]. In the training phase, for each feature, feature values of function instances are normalized in the range  $[-1, 1]$ . Then, the  $g$ -means [46] clustering of function instances is performed based on their normalized feature values, where  $g$ -means automatically determines an appropriate number of clusters. In the testing phase, an unseen function instance is assigned to the nearest cluster based on its normalized feature values. Finally, the best algorithm in the nearest cluster is selected. This study determines the best algorithm in each cluster according to the average ranking based on a performance measure.

#### D. Cross-validation methods

Algorithm selection will be ultimately performed on a real-world problem. Thus, the performance of algorithm selection systems should be evaluated in an unbiased manner. For this purpose, most previous studies (e.g., [11], [13]–[16]) used cross-validation methods for benchmarking algorithm selection systems for black-box optimization.

Below, we explain the following four cross-validation methods used in the literature: leave-one-instance-out cross-validation (LOIO-CV) [11], [15], leave-one-problem-out cross-validation (LOPO-CV) [11], [13], leave-one-problem-out-across-dimensions cross-validation (LOPOAD-CV) [14], and 10-fold randomized-instance cross-validation (RI-CV).

As explained in Section II-A, this study considers only the first five instances for each BBOB function, i.e.,  $|\mathcal{I}_i| = 5$  for  $f_i$  ( $i \in \{1, \dots, 24\}$ ). As in [14], this paper sets  $n$  to 2, 3, 5, and 10. Let  $\mathcal{I}^{\text{train}}$  be a set of function instances used in the training phase. Let also  $\mathcal{I}^{\text{test}}$  be a set of function instances used in the testing phase. Note that setting described here depends on previous benchmarking studies on the BBOB function set. Note also that the LOIO-CV, the LOPO-CV, and the LOPOAD-CV cannot always be applied to any function set. For example, the LOIO-CV is inapplicable when the number of instances for each function is only one.

1) *LOIO-CV*: A 5-fold cross-validation is performed on the 24 BBOB functions for each dimension  $n$ . In the  $i$ -th fold,  $\mathcal{I}^{\text{test}}$  is the set of the 24  $i$ -th instances  $\mathcal{I}_i$ , where  $|\mathcal{I}^{\text{test}}| = 24 \times 1 = 24$ . Thus,  $\mathcal{I}^{\text{train}}$  is  $\mathcal{I}_1 \cup \dots \cup \mathcal{I}_5 \setminus \mathcal{I}_i$ , where  $|\mathcal{I}^{\text{train}}| = 24 \times 4 = 96$ . Since function instances used in the training and testing phases are relatively similar, the LOIO-CV is likely easier than the LOPO-CV explained later.

For the LOIO-CV, the calculation of a performance measure value needs special care to avoid “data leakage”. Let us consider the calculation of the ERT value on a function  $f$  in the  $i$ -th fold, where  $i \in \{1, \dots, 5\}$ . Generally, test datasets should not be available in the training phase. For this reason, the ERT value in the training phase should be calculated based only on the remaining four instances ( $\{1, \dots, 5\} \setminus i$ ).

2) *LOPO-CV*: A 24-fold cross-validation is performed on the 24 BBOB functions for each dimension  $n$ . In the  $i$ -th fold,  $\mathcal{I}^{\text{test}}$  is the set of the five instances of the  $i$ -th function  $f_i$ ,

where  $|\mathcal{I}^{\text{test}}| = 1 \times 5 = 5$ . Thus,  $\mathcal{I}^{\text{train}}$  is  $\mathcal{I}_1 \cup \dots \cup \mathcal{I}_{24} \setminus \mathcal{I}_i$ , where  $|\mathcal{I}^{\text{train}}| = 23 \times 5 = 115$ . Since the 24 BBOB functions have different properties, instances used in the training and testing phases can be quite dissimilar for the LOPO-CV. For this reason, it is expected that algorithm selection for the LOPO-CV is challenging.

3) *LOPOAD-CV*: While the LOPO-CV is performed for each  $n$ , the LOPOAD-CV is performed across all four dimensions ( $n \in \{2, 3, 5, 10\}$ ). In the LOPOAD-CV, a 96-fold cross-validation is conducted on the 24 BBOB functions with all 4 dimensions, where  $24 \times 4 = 96$ . Let  $\mathcal{I}_{j,l}$  be the set of the five instances of the  $j$ -th function  $f_j$  with the  $l$ -th dimension. In the  $(j \times l)$ -th fold,  $\mathcal{I}^{\text{test}}$  is  $\mathcal{I}_{j,l}$ , where  $|\mathcal{I}^{\text{test}}| = 1 \times 5 = 5$ . Thus,  $\mathcal{I}^{\text{train}}$  is  $\mathcal{I}_{1,1} \cup \dots \cup \mathcal{I}_{24,4} \setminus \mathcal{I}_{j,l}$ , where  $|\mathcal{I}^{\text{train}}| = 95 \times 5 = 475$ .

Unlike the LOPO-CV and the LOIO-CV, the LOPOAD-CV evaluates the performance of algorithm selection systems on multiple dimensions. If both the performance of algorithms in  $\mathcal{A}$  and features have good scalability with respect to dimension, the LOPOAD-CV may be easier than the LOPO-CV. This is because both  $\mathcal{I}^{\text{test}}$  and  $\mathcal{I}^{\text{train}}$  include instances of the same function.

4) *RI-CV*: A 10-fold random cross-validation is performed on the 24 BBOB functions for each dimension  $n$ . The 120 function instances ( $24 \times 5 = 120$ ) are randomly grouped into 10 subsets of size 12 as follows:  $\mathcal{I}_1, \dots, \mathcal{I}_{10}$ . In the  $i$ -th fold,  $\mathcal{I}^{\text{test}}$  is the  $i$ -th subset  $\mathcal{I}_i$ , where  $|\mathcal{I}^{\text{test}}| = 12 \times 1 = 12$ . Thus,  $\mathcal{I}^{\text{train}}$  is  $\mathcal{I}_1 \cup \dots \cup \mathcal{I}_{12} \setminus \mathcal{I}_i$ , where  $|\mathcal{I}^{\text{train}}| = 12 \times 9 = 108$ . A 10-fold random cross-validation method has been generally used for algorithm selection in the combinatorial domain [31]. In contrast, only the study [12] used the 10-fold random cross-validation method for black-box optimization, where it did not describe details of the cross-validation method. In the best case, function instances used in the training and testing phases can be similar as in the LOIO-CV.

#### E. Pre-solvers

The concept of pre-solving was first adopted in SATZILLA [7], which is a representative algorithm selection system for SAT. Pre-solving was also incorporated into some modern algorithm selection systems for combinatorial optimization, e.g., CLASPFOLIO 2 [9] and 3S [25].

Before starting the algorithm selection process, SATZILLA runs two pre-defined pre-solvers on a given SAT instance within a short amount of time. Only when the pre-solvers cannot solve the instance, SATZILLA performs algorithm selection. If the pre-solvers can quickly solve easy instances, an algorithm selection system can focus only on hard instances. For some easy instances, the computation time of the algorithm selection process (including feature computation) dominates the runtime of a selected algorithm. The use of the pre-solvers can avoid such unnecessary computation time for algorithm selection on easy instances.

#### F. Related work in other problem domains

Below, we explain difficulties in algorithm selection for black-box optimization. A review of all previous studies for algorithm selection is beyond the scope of this paper. Interested readers can refer to exhaustive survey papers [4]–[6].

It is difficult to compare the performance of algorithm selection systems on across different problem domains (e.g., SAT and ASP) in a common platform. A benchmark library for algorithm selection (ASlib) [31] addresses this issue by providing scenarios for various problem domains, which are represented in a standardized format. ASlib was used in the algorithm selection competitions in 2015 and 2017 [3].

In addition to the performance sensitivity of optimizers and features (see the beginning of Sections I and III-A, respectively), the similarity of problem instances in the training and testing phases may be the main difference between algorithm selection for black-box numerical optimization and that for combinatorial optimization. In most scenarios (e.g., for SAT and the TSP), at least some problem instances in the training and testing phases were taken from the same distribution of problem instances even when using heterogeneous instance distribution. In contrast, as in the LOPO-CV, function instances used in the training and testing phases can be quite dissimilar for algorithm selection for black-box numerical optimization. We discuss the rational reason to adopt the LOPO-CV in Section VI-E later.

#### IV. OUR APPROACHES

First, Section IV-A discussed the importance of considering the number of function evaluations in the sampling phase. Then, Section IV-B explains that the relERT can overestimate the performance of algorithm selection systems. Section IV-C presents a local search method for constructing an algorithm portfolio of any size  $k$ . Finally, Section IV-D describes pre-solvers for black-box optimization.

##### A. Necessity of considering the number of function evaluations in the sampling phase

As explained in Section III-A, the sample of  $s$  solutions  $\mathcal{X} = \{\mathbf{x}_i\}_{i=1}^s$  and their objective values  $f(\mathcal{X})$  are needed to compute landscape features. Thus, the sampling phase requires  $s$  function evaluations. As shown in Table S.7, the previous studies [11], [13]–[16] set  $s$  to different numbers. For example, the two previous studies [11] and [14] set  $s$  to  $500 \times n$  and  $50 \times n$ , respectively. The start of the run of a selected optimizer should be delayed by  $s$  function evaluations used in the sampling phase. When  $s$  is too large, the overall performance of an algorithm selection system can deteriorate. For this reason,  $s$  is an important factor for benchmarking of algorithm selection systems. However, except for [14], most previous studies did not include the number of  $s$  evaluations in the total number of function evaluations. Note that most previous studies for algorithm selection in the combinatorial domain (e.g., [7], [9]) included the computation time of features in the total computation time of algorithm selection.

On the one hand, as discussed in [14], the sample  $\mathcal{X}$  can be reused for the initial population of a population-based evolutionary algorithm (e.g., DE [28]). The number of  $s$  function evaluations may be negligibly small when a selected optimizer requires a much large number of function evaluations, e.g., on hard instances. On the other hand,  $\mathcal{X}$  cannot be reused when an algorithm selection system selects

a non-evolutionary algorithm or a model-based evolutionary algorithm (e.g., CMA-ES [29]). The number of  $s$  function evaluations is *not negligible* when a selected optimizer requires a small number of function evaluations, e.g., on easy instances. In any case, there is no rational reason to ignore  $s$  function evaluations in the sampling phase. Based on the above discussions, as in [14], this study includes the number of  $s$  function evaluations in the total number of function evaluations. In this case, the setting of  $s$  influences the overall performance of algorithm selection systems.

##### B. Undesirable property of the relERT

Except for [15], all the previous studies used the relERT for evaluating the performance of algorithm selection systems in the fixed-target scenario. It has been believed that the lower bound of the relERT value is 1. However, we point out that the true lower bound of the relERT value is  $1/\text{bestERT}$ . There are three cases where the relERT value is less than 1.

First, it has been assumed that  $f(\mathbf{x}) > f_{\text{target}}$  for any  $\mathbf{x}$  in the sample  $\mathcal{X}$ . In addition, none of the previous studies considered a pre-solver. When the sample  $\mathcal{X}$  or a pre-solver reaches  $f_{\text{target}}$  faster than the best optimizer in a portfolio  $\mathcal{A}$ , an algorithm selection system obtains a lower ERT value than the bestERT value. In this case, it is possible to achieve a relERT value of less than 1.

Second, it has been assumed that the same optimizer is selected from  $\mathcal{A}$  for all the five instances of a function  $f$ . In contrast, it is possible that different optimizers are selected from  $\mathcal{A}$  for the five instances, respectively. Note that the best optimizer mentioned here is determined based on the number of function evaluations to reach  $f_{\text{target}}$  on each instance. When an algorithm selection system selects the best optimizer for each function instance in terms of the number of function evaluations, it obtains a relERT value of less than 1.

Third, apart from the above two “nice mistakes”, a relERT value can accidentally be less than 1. Let us consider a portfolio  $\mathcal{A}$  of two algorithms  $a_1$  and  $a_2$ . While the maximum number of function evaluations in  $a_1$  is 50, that in  $a_2$  is 5. Suppose that  $a_1$  reached  $f_{\text{target}}$  within 30 and 20 on the first two out of the five instances of  $f$  respectively, and  $a_1$  failed to reach  $f_{\text{target}}$  on the other three instances. Suppose also that all five runs of  $a_2$  were unsuccessful. In this case, the ERT value of  $a_1$  is  $(30 + 20 + 50 \times 3)/2 = 100$ , and that of  $a_2$  is not computable. Note that the missing relERT value of  $a_2$  is imputed by the PAR10 score described in Section II-B. Thus, the bestERT value is 100. If an algorithm selection system selects  $a_1$  on the first instance and  $a_2$  on the other four instances, the ERT value of the system is  $(30 + 5 \times 4)/1 = 50$ . As a result, the system achieves the relERT value of 0.5 ( $= 50/100$ ), even though it actually failed to select the best algorithm. One may incorrectly conclude that the system is two times faster than  $a_1$  in solving  $f$ . Section VI-A shows a practical example later.

The third undesirable case is due to the sensitivity of the ERT to the maximum number of function evaluations for an unsuccessful run. When all algorithms in  $\mathcal{A}$  use the same maximum number of function evaluations, the third case

never occurs. However, it is not realistic to assume such a termination condition.

Based on the above discussion, this study uses the SP1, instead of the ERT. As explained in Section II-B, the SP1 is robust to the maximum number of function evaluations. We do not claim that the SP1 is a better measure than the ERT *for any purpose*. Instead, we claim that the SP1 is more appropriate than the ERT *for benchmarking algorithm selection systems*. Inspired by the relERT, this study uses the relative SP1 (relSP1), which is the SP1 value of an algorithm in  $\mathcal{A}$  divided by the SP1 value of the best algorithm in  $\mathcal{A}$  on a function. This is the same as the procedure for obtaining the best ERT described in Section II-B. This study also applies the PAR10 to a missing relSP1 value.

### C. Local search for constructing algorithm portfolios

To analyze the influence of the portfolio size on the performance of algorithm selection systems, this study requires algorithm portfolios of any size  $k$ . Although the three construction methods reviewed in Section III-B are available, they are not appropriate for this purpose. The two methods proposed in [14], [40] cannot control the size of  $\mathcal{A}$ . Unlike the method proposed in [11], our method aims to optimize the VBS performance of  $\mathcal{A}$ . Note that our method is only for an analysis of algorithm portfolios. Thus, we do not claim that our method is more effective than the three existing methods.

Let  $\mathcal{R}$  be a set of  $l$  optimizers, where  $l = 209$  in this study. Let also  $\mathcal{A}$  be a portfolio of  $k$  optimizers, where  $k < l$ . As pointed out in [40], constructing  $\mathcal{A}$  can be defined as a subset selection problem. The goal of the problem is to select  $\mathcal{A} \subset \mathcal{R}$  that minimizes a quality measure  $m : \mathcal{A} \times \mathcal{I} \rightarrow \mathbb{R}$  on a set of problem instances  $\mathcal{I}$ . Fortunately, a general local search method for the subset selection problem can be applied to this portfolio construction problem in a straightforward manner. This study uses a first-improvement local search method [47], which was proposed for the hypervolume subset selection problem. Algorithm S.1 shows the local search method. First, the local search method initializes  $\mathcal{A}$  by randomly selecting  $k$  optimizers from  $\mathcal{R}$ . Then, for each iteration, the local search method swaps a pair of optimizers  $a \in \mathcal{A}$  and  $a' \in \mathcal{R} \setminus \mathcal{A}$ . The swap operation is performed until there is no pair to improve the following ranking-based quality measure  $m$ :

$$m(\mathcal{A}) = \sum_{n \in \{2,3,5,10\}} \sum_{i \in \{1, \dots, 24\}} \text{score}(\mathcal{A}, f_i^n), \quad (3)$$

$$\text{score}(\mathcal{A}, f) = \begin{cases} c \times \min_{a \in \mathcal{A}} \{\text{rank}(a, f)\} & \text{if } \exists a \text{ solves } f \\ 1 & \text{otherwise} \end{cases}, \quad (4)$$

where  $f_i^n$  in equation (3) is the  $n$ -dimensional  $i$ -th BBOB function. We rank  $l$  optimizers in  $\mathcal{R}$  based on their SP1 values for each  $f_i^n$ . In equation (4),  $\text{rank}(a, f)$  is a ranking of  $a$  on  $f$  in  $\mathcal{R}$ . In equation (4),  $c$  is a coefficient value. We set  $c$  to  $1/(l \times 24 \times 4)$  for the 24 BBOB functions with  $n \in \{2, 3, 5, 10\}$ . If no  $a$  in  $\mathcal{A}$  reaches the target value, a penalty value of 1 is assigned to  $\text{score}(\mathcal{A}, f)$ . If  $m(\mathcal{A}) \geq 1$  in equation (3), it means that no optimizer in  $\mathcal{A}$  could reach the target value  $f_{\text{target}}$  for at least one function.

### D. Pre-solvers for black-box optimization

It is straightforward to incorporate the concept of pre-solving into an algorithm selection system for black-box optimization. In the pre-solving phase, a pre-solver is applied to a given problem with a small budget of function evaluations (e.g.,  $50 \times n$  evaluations) before the algorithm selection process starts. Note that it is possible to use more than one pre-solvers as in SATZILLA [7]. When the pre-solver reaches the target value  $f_{\text{target}}$  on the problem, algorithm selection is not performed. In this case, the algorithm selection system does not need to generate the sample  $\mathcal{X} = \{\mathbf{x}_i\}_{i=1}^s$  of size  $s$  to compute features. In other words,  $s$  function evaluations can be saved. We expect that using a pre-solver can improve the performance of algorithm selection systems on easy instances.

It is desirable that a pre-solver can reach  $f_{\text{target}}$  on an easy instance with a small number of function evaluations. This study uses SLSQP [48] and SMAC [49] as pre-solvers. Since evolutionary algorithms (e.g., DE and CMA-ES) generally require a relatively large number of function evaluations to find a good solution, they are not appropriate as pre-solvers. In [50], Hansen showed the excellent performance of SLSQP on the BBOB functions for a small number of function evaluations. SMAC is also a representative Bayesian optimization approach, which performs particularly well for computationally expensive optimization. Thus, SLSQP and SMAC are reasonable first choices.

## V. EXPERIMENTAL SETUP

This section explains the experimental setup. Unlike previous studies, we performed 31 independent runs of algorithm selection systems, including pre-solving, sampling, feature computation, and cross-validation. We conducted all experiments on a workstation with an Intel(R) 52-Core Xeon Platinum 8270 (26-Core $\times$ 2) 2.7GHz and 768GB RAM using Ubuntu 18.04. As in [14], we used the 24 noiseless BBOB functions [17] with dimensions  $n \in \{2, 3, 5, 10\}$ . We conducted our experiment by using the COCO platform [27]. We also used the benchmarking results of 209 optimizers provided by the COCO data archive. As in [13], [14], we set a precision level  $\epsilon$  to  $10^{-2}$  for the ERT and SP1 calculations. Note that the same  $\epsilon$  value is used for pre-solvers. The study [51] showed that the setting of  $\epsilon$  does not significantly influence the accuracy of performance predictors.

We used the following nine non-expensive and scalable flacco feature classes in Table S.8: `ela_distr`, `ela_level`, `ela_meta`, `nbc`, `disp`, `ic`, `basic`, `limo`, and `pca`. We employed `pflacco`, which is the Python interface of `flacco` (<https://github.com/Reiyan/pflacco>). Our preliminary results showed that feature selection deteriorates the performance of algorithm selection systems in some cases. This is consistent with the results reported in [52]. Except for the classification-based selection method, wrapper feature selection approaches require extremely high computational cost as the portfolio size increases. For these reasons, as in [11]–[13], [15], [16], we did not perform feature selection.

As in most previous studies, we used the improved Latin hypercube sampling [53] to generate the sample  $\mathcal{X} = \{\mathbf{x}_i\}_{i=1}^s$

for feature computation. We employed the `lhs` function in the `flacco` package. As in [14], we set  $s$  to  $50 \times n$  unless explicitly noted, where the study [14] selected this  $s$  value based on the results in [54]. For pre-solvers, we employed the SciPy implementation of SLSQP and the SMAC3 implementation [55] of SMAC. We used the default parameter settings, including the termination criteria for SLSQP. The maximum number of function evaluations was set to  $50 \times n$ , which is the same as the sampling phase.

We used random forest [56] for the four supervised learning-based algorithm selection methods explained in Section III-C. The random forest is a representative ensemble machine learning method that uses multiple decision trees. Each tree is fitted to a subset of randomly selected features. Then, the prediction is performed based on the number of votes by the trees. Typical hyperparameters in the random forest include the number of trees (`n_estimators`) and the feature subset size (`max_features`). As in ISAC [45], we used  $g$ -means [46] for the clustering-based selection method. We employed the scikit-learn implementation of random forest and the `pyclustering` implementation of  $g$ -means with the default parameters, e.g., `n_estimators=100` and `max_features=auto`.

Table S.9 shows 14 algorithm portfolios used in Section VI. The first five portfolios ( $\mathcal{A}_{kt}$ , ...,  $\mathcal{A}_{mk}$ ) were used in the previous studies [11], [13]–[16], respectively. We named each portfolio by taking the initial letters of the authors of the corresponding paper. For each of the nine portfolios ( $\mathcal{A}_{ls2}$ , ...,  $\mathcal{A}_{ls18}$ ), we performed 31 runs of the local search method explained in Section IV-C. Then, the best portfolio is selected in terms of  $m$  in equation (3). We confirmed that all the nine portfolios achieve  $m$  values less than 1.

Sections VI-D and VI-F use the performance score [57] to rank multiple algorithm selection systems. For each dimension  $n$ , let us consider a comparison of  $l$  algorithm selection systems  $S_1, \dots, S_l$  based on their 31 mean relSP1 values from 31 independent runs. For  $i \in \{1, \dots, l\}$  and  $j \in \{1, \dots, l\} \setminus \{i\}$ , if  $S_j$  performs significantly better than  $S_i$  using the Wilcoxon rank-sum test with  $p < 0.05$ , then  $\delta_{i,j} = 1$ ; otherwise,  $\delta_{i,j} = 0$ . The score  $P(S_i)$  is defined as follows:  $P(S_i) = \sum_{j \in \{1, \dots, l\} \setminus \{i\}} \delta_{i,j}$ . The score  $P(S_i)$  is the number of systems that outperform  $S_i$  for each  $n$ . A small  $P(S_i)$  means that  $S_i$  has a better performance among the  $l$  systems.

In addition, we calculated the average rankings of algorithm selection systems by the Friedman test [58]. We used the `CONTROLTEST` package (<https://sci2s.ugr.es/sicidm>) to calculate the rankings. However, the rankings by the Friedman test were generally consistent with those by the performance score. For this reason, this paper shows only the latter results.

## VI. RESULTS

This section analyzes algorithm selection systems to answer the six research questions **RQ1–RQ6** discussed in Section I. First, Section VI-A shows a comparison of the ERT and the SP1 (**RQ1**). Then, Section VI-B examines the influence of randomness on results of algorithm selection, including pre-solving, sampling, feature computation, and cross-validation (**RQ2**). Section VI-C investigates the effectiveness of SLSQP

TABLE I: Comparison of HCMA, HMLSL, and the regression-based algorithm selection system (AS) using  $\mathcal{A}_{kt}$  on the 5-dimensional  $f_{24}$  for the LOPO-CV. If the corresponding run was successful, the number in parentheses is 1. Otherwise, it is 0.

	HCMA	HMLSL	AS #1	AS #2
FES on $i_1$	3 097 698 (1)	100 021 (0)	3 097 948 (1)	3 097 948 (1)
FES on $i_2$	2 254 446 (1)	100 002 (0)	2 254 696 (1)	100 252 (0)
FES on $i_3$	5 001 954 (0)	100 206 (0)	100 456 (0)	100 456 (0)
FES on $i_4$	5 001 015 (0)	100 008 (0)	100 258 (0)	100 258 (0)
FES on $i_5$	5 001 872 (0)	100 006 (0)	100 256 (0)	100 256 (0)
ERT	10 178 492.50	Na	2 826 807	3 499 170
SP1	6 690 180	Na	6 690 805	15 489 740
relERT	1	Na	0.28	0.34
relSP1	1	Na	1.00009	2.32

and SMAC as a pre-solver (**RQ3**). Based on results shown in Section VI-C, Sections VI-D, VI-E, and VI-F use SLSQP as a pre-solver. Section VI-D shows a comparison of the five algorithm selection methods (**RQ4**). Section VI-E analyzes the difficulty of outperforming the SBS (**RQ5**). Finally, Section VI-F examines the performance of algorithm selection systems with the 14 algorithm portfolios shown in Table S.9 (**RQ6**).

The influence of cross-validation methods on the results of algorithm selection systems is unclear. For this reason, we here investigate it by using the four cross-validation methods. We essentially aim to analyze the influence of the similarity of function instances used in the training and testing phases. However, it would be misleading to make conclusions based on the results achieved by multiple cross-validation methods. Thus, we answer each question based only on the representative results for the LOPO-CV, which is the most practical (see the discussion in Section VI-E).

### A. Comparison of the ERT and the SP1 (RQ1)

Table I shows a comparison of HCMA, HMLSL, and the regression-based algorithm selection system using  $\mathcal{A}_{kt}$  on the 5-dimensional  $f_{24}$  for the LOPO-CV. Note that only this section uses the relERT, instead of the relSP1. Table I reports the number of function evaluations (FES) used in the search on each of the five instances. Table I also reports the ERT, SP1, relERT, and relSP1 values. HCMA performs the best on  $f_{24}$  with  $n = 5$ . All five runs of HMLSL are unsuccessful on  $f_{24}$  with  $n = 5$ . Table I shows results of the best run (AS #1) and the second best run (AS #2) out of the 31 runs of the system in terms of the ERT (not the SP1). For AS #1 and AS #2, the number of function evaluations includes  $50 \times 5 = 250$  function evaluations in the sampling phase.

As shown in Table I, AS #1 selected HCMA on the first two instances and HMLSL on the other instances. In contrast, AS #2 selected HCMA on the first instance and HMLSL on the other instances. Thus, both AS #1 and AS #2 failed to select HCMA on all five instances. Nevertheless, AS #1 and AS #2 perform significantly better than HCMA in terms of the ERT. AS #1 and AS #2 also achieve relERT values of less than 1. As seen from Table I, the maximum number of function evaluations of HCMA is approximately 50 times larger than that of HMLSL. For this reason, when HMLSL is selected for unsuccessful runs, a better ERT value can be obtained. This is exactly the third case explained in Section IV-B.

We observed the undesirable ERT (and relERT) results due to the third case only on  $f_{24}$ , where no optimizer in  $\mathcal{A}_{kt}$  can solve all the five instances of  $f_{24}$ . However, as demonstrated here, the ERT and the relERT can possibly overestimate the performance of algorithm selection systems. In contrast, as shown in Table I, the SP1 and the relSP1 do not overestimate the performance of AS #1 and AS #2. AS #1 performs quite slightly worse than HCMA in terms of the SP1 and the relSP1. The relSP1 of AS #2 is 2.32 due to the unsuccessful run on the second instance.

**Answers to RQ1** We suggest using the SP1 and the relSP1 for benchmarking algorithm selection systems, instead of the ERT and the relERT. It should be highly noted that the above-discussed issue of the ERT never occurs for benchmarking a single optimizer and multiple optimizers with exactly the same termination conditions.

### B. Influence of randomness (RQ2)

The previous studies [11]–[14] discussed the performance of algorithm selection systems based on the mean of the performance measure values (e.g., the relERT value) on the 24 BBOB functions for a single run. In contrast, Fig. 1 shows the distributions of 31 “mean relSP1” values, which were obtained by 31 runs of the pairwise classification-based system with  $\mathcal{A}_{kt}$  for the four cross-validation methods. Figs. S.4–S.8 show results of the five systems with the five portfolios, respectively. We do not explain the results in Figs. S.4–S.8 here, but they are almost consistent with Fig. 1. Figs. S.9–S.13 also show the non-log scale versions of Figs. S.4–S.8, respectively.

As shown in Fig. 1, the distribution of 31 “mean relSP1” values depends on the cross-validation method. While the dispersion in the distribution of the “mean relSP1” values is relatively small for the LOIO-CV (except for the results on  $n = 3$ ), that is relatively large for the LOPO-CV, the LOPOAD-CV, and the RI-CV (except for the results on  $n = 5$ ). For example, the minimum “mean relSP1” value on  $n = 10$  for the LOPO-CV is 7.10, but the maximum one is 542.62. Thus, the best-case performance of the algorithm selection system is approximately 76 times better than the worst-case one. Figs. S.14–S.17 show the distribution of “relSP1” (not “mean relSP1”) values on the 24 BBOB functions for the four cross-validation methods, respectively. Figs. S.14–S.17 indicate that large variations in relSP1 values are due to unsuccessful selection on multimodal functions.

Note that we cannot generalize our observation in Fig. 1 for any portfolio. For example, as seen from Fig. S.7(c), when using  $\mathcal{A}_{bmtp}$ , the dispersion in the distribution of the “mean relSP1” values is relatively large for the LOIO-CV. Thus, our observations in Figs. 1 and S.7(c) are inconsistent with each other. However, we can say that the influence of randomness is not negligible in any case.

**Answers to RQ2** We demonstrated that randomness influences results of algorithm selection systems in most cases. This observation is consistent with the results for combinatorial optimization investigated in [3]. Because the best-case and worst-case results can be significantly different, we suggest performing multiple runs of algorithm selection systems.

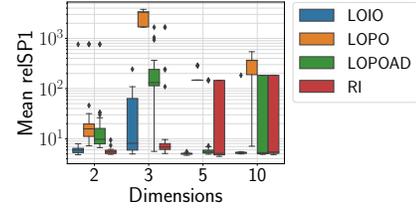


Fig. 1: Distributions of 31 mean relSP1 values of the pairwise classification-based algorithm selection system with  $\mathcal{A}_{kt}$ .

### C. Effectiveness of a pre-solver (RQ3)

Table II shows results of the pairwise classification-based system with and without a pre-solver. In Table II, AS ( $50 \times n$ ) and AS ( $100 \times n$ ) are the systems with the sample size  $50 \times n$  and  $100 \times n$ , respectively. Recall that the default sample size is  $50 \times n$ . SLSQP-AS and SMAC-AS are the system using SLSQP and SMAC as pre-solvers, respectively. Tables II(a) and II(b) show the median of 31 “mean relSP1” values when using  $\mathcal{A}_{kt}$  and  $\mathcal{A}_{bmtp}$  as algorithm portfolios, respectively. Table II does not show the results for the RI-CV, but they are similar to the results for the LOIO-CV. Tables S.10–S.14 show results of the five systems with the five portfolios ( $\mathcal{A}_{kt}$ , ...,  $\mathcal{A}_{mk}$ ), respectively. The symbols + and – indicate that a given system performs significantly better (+) and significantly worse (–) than AS ( $50 \times n$ ) according to the Wilcoxon rank-sum test with  $p < 0.05$ . Apart from the comparison with AS ( $50 \times n$ ), results that are better than those of the SBS are highlighted in **dark gray**. The comparison with the SBS is discussed later in Section VI-E.

As seen from Table II, SLSQP-AS performs significantly better than AS ( $50 \times n$ ) for the three cross-validation methods in most cases. Tables S.10–S.14 also show that SLSQP-AS does not perform significantly worse than AS ( $50 \times n$ ) in this study. These results suggest that the performance of algorithm selection systems can be significantly improved by using SLSQP as a pre-solver.

Fig. 2 shows a comparison of AS ( $50 \times n$ ) and SLSQP-AS on the 24 BBOB functions with  $n = 5$ . Fig. 2 shows the results of a single run with a median “mean relSP1” value among 31 runs. Fig. S.18 also shows the distribution of the number of function evaluations used by SLSQP in the pre-solving phase. Recall that the maximum number of function evaluations in the pre-solving phase is  $50 \times n$ . As shown in Fig. 2, SLSQP-AS obtains much smaller values of the relSP1 than AS ( $50 \times n$ ) on  $f_1$  (the Sphere function),  $f_5$  (the Linear Slope function), and  $f_{14}$  (the different powers function). These results demonstrate that the use of SLSQP is highly effective on some unimodal functions. In other words,  $s$  function evaluations can be saved by using SLSQP as the pre-solver. As noted in [50], the poor performance of SLSQP on  $f_6$  (the Attractive Sector function) is mainly due to a small number of function evaluations. In contrast, SLSQP does not reach the target value within  $50 \times n$  function evaluations on the 12 multimodal functions  $f_3, f_4, f_{15}, \dots, f_{24}$ , except for some instances of  $f_{21}$  (the Gallagher’s Gaussian 101-me Peaks function). Since the initial  $50 \times n$  function evaluations in the

TABLE II: Results of the pairwise classification-based algorithm selection system with and without the pre-solvers.

(a) LOIO-CV ( $\mathcal{A}_{kt}$ )					(b) LOPO-CV ( $\mathcal{A}_{kt}$ )					(c) LOPOAD-CV ( $\mathcal{A}_{kt}$ )				
System	$n = 2$	$n = 3$	$n = 5$	$n = 10$	System	$n = 2$	$n = 3$	$n = 5$	$n = 10$	System	$n = 2$	$n = 3$	$n = 5$	$n = 10$
AS ( $50 \times n$ )	5.94	8.13	5.01	5.28	AS ( $50 \times n$ )	15.75	3337.17	146.83	363.62	AS ( $50 \times n$ )	9.76	131.42	5.42	5.29
AS ( $100 \times n$ )	8.75-	11.33-	7.88-	8.59-	AS ( $100 \times n$ )	16.66	3340.90	149.85-	366.78	AS ( $100 \times n$ )	11.22-	124.18	8.01-	8.60-
SLSQP-AS	3.53+	4.85+	2.54+	2.49+	SLSQP-AS	11.58+	3334.35	144.06+	360.35	SLSQP-AS	7.36+	128.10	2.92+	2.50+
SMAC-AS	6.18	7.82	7.07-	8.60-	SMAC-AS	15.59	3338.37	148.80-	366.99-	SMAC-AS	9.68	115.13	7.54-	8.64-
SBS	5.81	9.82	4.49	6.76	SBS	5.81	9.82	4.49	6.76	SBS	5.81	9.82	4.49	6.76

(d) LOIO-CV ( $\mathcal{A}_{bmtP}$ )					(e) LOPO-CV ( $\mathcal{A}_{bmtP}$ )					(f) LOPOAD-CV ( $\mathcal{A}_{bmtP}$ )				
System	$n = 2$	$n = 3$	$n = 5$	$n = 10$	System	$n = 2$	$n = 3$	$n = 5$	$n = 10$	System	$n = 2$	$n = 3$	$n = 5$	$n = 10$
AS ( $50 \times n$ )	7.41	36.15	96.86	158.07	AS ( $50 \times n$ )	12.30	947.32	109.69	323.98	AS ( $50 \times n$ )	12.93	13.28	7.19	158.48
AS ( $100 \times n$ )	7.82	18.40+	7.07	160.06-	AS ( $100 \times n$ )	14.23-	944.64	106.29	324.18	AS ( $100 \times n$ )	14.82-	12.54	7.41	160.06-
SLSQP-AS	5.55	34.76	95.28+	156.52+	SLSQP-AS	10.20+	743.19	103.99	314.43+	SLSQP-AS	11.41+	11.73	5.29+	156.84+
SMAC-AS	6.77	36.83	97.93-	159.64-	SMAC-AS	11.51	716.73	111.16	325.62	SMAC-AS	10.76+	13.79	8.27	160.10-
SBS	11.53	35.92	102.10	315.32	SBS	11.53	35.92	102.10	315.32	SBS	11.53	35.92	102.10	315.32

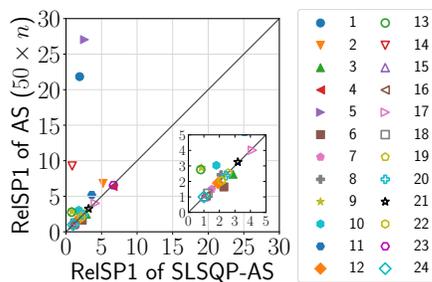


Fig. 2: Distribution of 31 relSP1 values of the pairwise classification-based algorithm selection system with  $\mathcal{A}_{kt}$  on the 24 BBOB functions with  $n = 5$  for the LOIO-CV.

pre-solving phase are wasted in this case, SLSQP-AS performs slightly worse than AS ( $50 \times n$ ) in terms of the relSP1. However, any optimizer requires a large number of function evaluations for hard multimodal functions to reach the target value. For this reason, as shown in Fig. 2, a small number of additional function evaluations in the pre-solving phase does not drastically deteriorate the performance of the system.

In contrast to SLSQP-AS, Table II shows that the performance of SMAC-AS is not better than that of AS ( $50 \times n$ ) in most cases. Our results suggest that SMAC is not appropriate for a pre-solver. AS ( $100 \times n$ ) performs significantly worse than AS ( $50 \times n$ ) in most cases, except for the result for  $n = 2$  in Table S.10(h) and the result for  $n = 3$  in Table S.12(j). Based on these results, it would be better to allocate half of the budget ( $50 \times n$  function evaluations) to the pre-solving and sampling phases, instead of allocating all the budget ( $100 \times n$  function evaluations) to the sampling phase.

One may be interested in the influence of the sample size  $s$  on the performance of algorithm selection systems. Table S.15 further shows the comparison of the five algorithm selection system with  $s \in \{50 \times n, 25 \times n, 100 \times n, 200 \times n\}$ .  $\mathcal{A}_{kt}$  was used in this comparison. As shown in Table S.15,  $s \geq 100 \times n$  is less effective than  $s = 50 \times n$  in most cases. In contrast,  $s = 25 \times n$  is more effective than  $s = 50 \times n$  in some cases. These results suggest that  $s$  can possibly be set to

$s \leq 50 \times n$ , which is the general setting. An in-depth analysis of  $s$  is another future work.

Let  $\mathcal{Y}$  be a set of all solutions found so far by a pre-solver. As in [52], [59], it is possible to compute features based on the union of the sample  $\mathcal{X}$  and  $\mathcal{Y}$ . A similar idea was also discussed in [13]. Table S.16 shows a comparison of two SLSQP-AS variants that compute features based on  $\mathcal{X}$  and  $\mathcal{X} \cup \mathcal{Y}$ , respectively.  $\mathcal{A}_{kt}$  was used in this comparison. Although SLSQP-AS with  $\mathcal{X} \cup \mathcal{Y}$  outperforms SLSQP-AS with  $\mathcal{X}$  in some cases, we cannot say that SLSQP-AS with  $\mathcal{X} \cup \mathcal{Y}$  generally performs better than SLSQP-AS with  $\mathcal{X}$ . Table S.16 also shows that features computed based on  $\mathcal{X} \cup \mathcal{Y}$  are less effective than those based on  $\mathcal{X}$  alone in some cases. The unpromising results may be due to the extremely biased distribution of solutions in  $\mathcal{Y}$ . Since SLSQP terminates early on some functions, the size of  $\mathcal{X} \cup \mathcal{Y}$  is not constant for all the 24 BBOB functions. However, as revealed in [22], the sample size should be the same for all the functions to obtain effective features. Designing a method for “cleansing”  $\mathcal{Y}$  is a future research topic.

**Answers to RQ3** Our results demonstrated that the overall performance of algorithm selection systems can be significantly improved by using SLSQP as a pre-solver in most cases, especially for easy function instances. Although the pre-solving approach has not been considered for black-box optimization, it would be better to incorporate a pre-solver into algorithm selection systems. We believe that our promising findings here facilitate the use of pre-solvers to researchers in the field of black-box numerical optimization.

#### D. Comparison of the five algorithm selection methods (RQ4)

Based on the results reported in Section VI-C, we use SLSQP-AS in the rest of this paper. Tables III and IV show performance score values of the five algorithm selection systems using  $\mathcal{A}_{kt}$  and  $\mathcal{A}_{bmtP}$ , respectively. Tables III and IV do not show the results for the RI-CV, which are similar to the results for the LOIO-CV. The best and second-best results are highlighted in **dark gray** and **gray**, respectively. Tables S.17–S.19 show results when using the 14 portfolios. Tables S.20–S.22 also show results of the Friedman test.

TABLE III: Results of the five algorithm selection systems. Tables (a)–(c) show performance score values of the five systems using  $\mathcal{A}_{kt}$  for the three cross-validation methods, respectively.

(a) LOIO-CV					(b) LOPO-CV					(c) LOPOAD-CV				
$n$					$n$					$n$				
2 3 5 10					2 3 5 10					2 3 5 10				
Cla.	1	1	0	0	Cla.	3	4	3	3	Cla.	2	4	3	3
Reg.	2	3	2	2	Reg.	1	0	0	0	Reg.	0	0	1	0
P-Cla.	0	0	1	1	P-Cla.	0	2	1	2	P-Cla.	0	0	0	0
P-Reg.	3	4	3	3	P-Reg.	3	0	2	1	P-Reg.	4	1	1	1
Clu.	3	2	4	4	Clu.	0	2	3	4	Clu.	0	0	3	4

TABLE IV: Results of the five systems using  $\mathcal{A}_{bmtmp}$ .

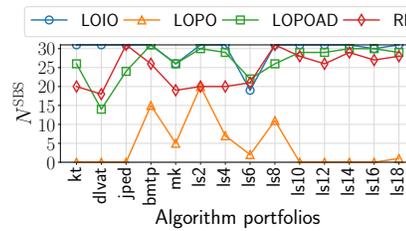
(a) LOIO-CV					(b) LOPO-CV					(c) LOPOAD-CV				
$n$					$n$					$n$				
2 3 5 10					2 3 5 10					2 3 5 10				
Cla.	0	3	0	2	Cla.	0	2	3	3	Cla.	3	0	0	0
Reg.	3	1	2	0	Reg.	0	0	0	0	Reg.	0	3	3	3
P-Cla.	0	0	1	0	P-Cla.	0	2	0	2	P-Cla.	0	0	2	1
P-Reg.	0	0	1	0	P-Reg.	1	1	0	0	P-Reg.	0	2	1	2
Clu.	4	4	4	4	Clu.	4	2	4	4	Clu.	4	4	4	4

On the one hand, as seen from Table III, for the LOIO-CV and the LOPOAD-CV, the pairwise classification-based system performs the best in most cases when using  $\mathcal{A}_{kt}$ . This observation is consistent with the results for combinatorial optimization shown in [26]. For the LOIO-CV, the classification-based system also performs well. For the LOPO-CV, the regression-based system is the best performer. The results of the systems with  $\mathcal{A}_{dlvat}$  and  $\mathcal{A}_{jped}$  are similar to those with  $\mathcal{A}_{kt}$ .

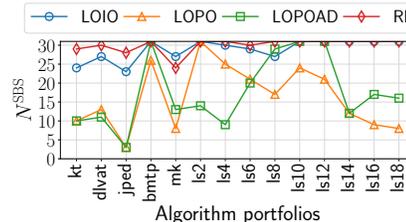
On the other hand, as seen from Table IV, for the LOIO-CV, the pairwise regression-based and pairwise classification-based systems have the same performance when using  $\mathcal{A}_{bmtmp}$ . For the LOPO-CV, the pairwise regression-based system is competitive with the regression-based system for  $n \in \{5, 10\}$ . These results indicate that the best algorithm selection method depends on algorithm portfolios and cross-validation methods.

Since  $|\mathcal{A}_{kt}| = 12$  and  $|\mathcal{A}_{bmtmp}| = 4$ , one may think that the portfolio size determines the best algorithm selection method. To investigate the scalability of the five systems with respect to the portfolio size, Tables S.18 and S.19 show results of the five systems with  $\mathcal{A}_{ls2}, \dots, \mathcal{A}_{ls18}$ , where  $|\mathcal{A}_{ls2}| = 2, \dots, |\mathcal{A}_{ls18}| = 18$ . As seen from Tables S.18 and S.19, there is no clear correlation between the portfolio size and the performance rankings of the five systems. For example, for the LOIO-CV, the pairwise regression-based system does not obtain the best performance on any dimension when using  $\mathcal{A}_{ls4}$  (Table S.18(d)). In contrast, for the LOPO-CV, the pairwise regression-based system performs the best for all dimensions when using  $\mathcal{A}_{ls12}$  (Table S.18(q)). Note that  $|\mathcal{A}_{bmtmp}| = |\mathcal{A}_{ls4}|$  and  $|\mathcal{A}_{kt}| = |\mathcal{A}_{ls12}|$ . These results suggest that components of portfolios are more important than the portfolio size to determine the best algorithm selection method. Such an observation has not been reported even in combinatorial optimization.

One exception is the results of the classification-based system for the LOPOAD-CV. Tables III, IV, S.17–S.19 show that the classification-based system is the best performer for some



(a) Pairwise classification



(b) Regression

Fig. 3:  $N^{\text{SBS}}$  in which the pairwise classification-based and regression-based systems outperform the SBS for  $n = 5$ .

dimensions when using a small-size portfolio (e.g.,  $\mathcal{A}_{bmtmp}$  and  $\mathcal{A}_{ls6}$ ). In contrast, the classification-based system cannot perform the best when using a large-size portfolio (e.g.,  $\mathcal{A}_{kt}$ ,  $\mathcal{A}_{ls10}, \dots, \mathcal{A}_{ls18}$ ). The poor scalability of the classification-based system is simply because multi-class classification with many classes is difficult. Here, the classification-based system uses a random forest model (see Section III-C).

**Answers to RQ4** Our results showed that the best algorithm selection method mainly depends on the components of algorithm portfolios. For the challenging LOPO-CV, the regression-based method is the best performer, followed by the pairwise regression-based method. This observation suggests that the regression-based method is likely to perform the best when the training and testing instances are quite dissimilar as in the LOPO-CV.

#### E. On the difficulty of outperforming the SBS (RQ5)

Fig. 3 shows a comparison of the SBS and the pairwise classification-based and regression-based systems according to their 31 mean relSPI values for  $n = 5$ . Fig. 3 shows the number of times ( $N^{\text{SBS}} \in [0, 31]$ ) in which a system outperforms the SBS over 31 runs. We use the 14 portfolios. We measure the difficulty of outperforming the SBS based on  $N^{\text{SBS}}$ . If a system with a portfolio  $\mathcal{A}$  achieves a small  $N^{\text{SBS}}$  value, we say that it is difficult for the system to outperform the SBS in  $\mathcal{A}$ . Figs. S.19–S.23 show results of all the five systems for  $n \in \{2, 3, 5, 10\}$ .

As shown in Fig. 3, for the LOIO-CV and the LOPOAD-CV, the pairwise classification-based system outperforms the regression-based system for most portfolios in terms of  $N^{\text{SBS}}$ . In contrast, for the LOPO-CV and the RI-CV, the regression-based system achieves a better  $N^{\text{SBS}}$  value than the pairwise classification-based system for all the 14 portfolios. These results based on  $N^{\text{SBS}}$  are consistent with the results shown in Section VI-D. This is the first study to show that a system

outperforms the SBS for the LOIO-CV and the LOPO-CV for black-box numerical optimization *when considering the number of function evaluations used in the sampling phase*.

As seen from Fig. 3,  $N^{\text{SBS}}$  depends on portfolios and cross-validation methods. However, we cannot find the correlation between the portfolio size and  $N^{\text{SBS}}$ . For the LOIO-CV, Fig. 3(a) shows that the pairwise classification-based system performs better than the SBS for all 31 runs when using 11 out of the 14 portfolios. The results for the LOPOAD-CV, the RI-CV, and the LOIO-CV are similar. For the LOPO-CV, when using 7 out of the 14 portfolios, the pairwise classification-based system obtains  $N^{\text{SBS}} = 0$ . Our results show that it is generally the most difficult to outperform the SBS for the LOPO-CV. Of course, this is not always true. For example, Fig. 3(b) shows that  $N^{\text{SBS}}$  obtained by the regression-based system for the LOPO-CV is larger than that for the LOPOAD-CV when using  $\mathcal{A}_{\text{dlvat}}$ ,  $\mathcal{A}_{\text{ls2}}$ ,  $\mathcal{A}_{\text{ls4}}$ , and  $\mathcal{A}_{\text{ls6}}$ .

Although the above discussions are based on the results for  $n = 5$ , Figs. S.19–S.23 show that  $N^{\text{SBS}}$  depends on the dimension  $n$ . As shown in Figs. S.19–S.23,  $N^{\text{SBS}}$  for  $n \in \{2, 3\}$  is generally smaller than  $N^{\text{SBS}}$  for  $n \in \{5, 10\}$ . This is because the SBS can reach the target value for a relatively small number of function evaluations on low-dimensional function instances, even including hard multimodal instances. In this case, the budget of function evaluations used in the sampling phase is a critical disadvantage for algorithm selection.

Below, we discuss which cross-validation method should be used for benchmarking algorithm selection systems. In [60], van Stein et al. analyzed the fitness landscape of a neural architecture search (NAS) problem with  $n = 23$ . Their results based on the ELA approach revealed that the fitness landscape of the NAS problem is different from that of any BBOB function. We do not intend to generalize their conclusion, but it is practical to assume that problem instances in the testing and training phases are different. The LOPO-CV is appropriate for this purpose. In contrast to the LOPO-CV, it would *not* be better to use the LOIO-CV. In the LOIO-CV, function instances used in the training and testing phases are always similar. If a researcher uses the LOIO-CV, she/he can overestimate the performance of an algorithm selection system that does not actually work well for any real-world setting. For a similar reason, we do not suggest using the RI-CV.

Let  $f^{\text{real}}$  with  $n^{\text{real}}$  be an  $n^{\text{real}}$ -dimensional real-world problem. It is rare that  $f^{\text{real}}$  instances with  $n \neq n^{\text{real}}$  are available in the training phase. It is also practical to use the same  $n$  in the training and testing phases. Our results also showed that the LOPOAD-CV is less challenging than the LOPO-CV. It may be better not to use the LOPOAD-CV without a particular reason.

**Answers to RQ5** We demonstrated that the difficulty of outperforming the SBS depends on algorithm portfolios and dimensions. For example, even for the LOPO-CV, our results showed that the regression-based system can often outperform the SBS when using  $\mathcal{A}_{\text{bmtp}}$ . Since a result obtained using a single portfolio can be misleading, it would be better to use multiple portfolios (e.g.,  $\mathcal{A}_{\text{bmtp}}$  and  $\mathcal{A}_{\text{kt}}$ ) for benchmarking algorithm selection systems. Based on the discussion, we suggest using the LOPO-CV.

TABLE V: Mean relSP1 values of the VBS and the SBS in the 14 algorithm portfolios.

AP	VBS				SBS			
	$n = 2$	$n = 3$	$n = 5$	$n = 10$	$n = 2$	$n = 3$	$n = 5$	$n = 10$
$\mathcal{A}_{\text{kt}}$	2.59	1.76	1.62	1.56	15.71	14.36	6.63	9.53
$\mathcal{A}_{\text{dlvat}}$	3.19	1.84	1.64	1.60	15.71	14.36	6.63	9.53
$\mathcal{A}_{\text{jped}}$	2.52	1.83	1.52	1.58	15.71	14.36	6.63	9.53
$\mathcal{A}_{\text{bmtp}}$	75.22	56.16	5.20	5.55	115.88	174.38	40173.96	46560.09
$\mathcal{A}_{\text{mk}}$	14.96	54.71	6.77	8.63	33.06	174.38	40173.96	46560.09
$\mathcal{A}_{\text{ls2}}$	9.58	7.04	3.37	2.37	15.71	14.36	6.63	9.53
$\mathcal{A}_{\text{ls4}}$	2.84	5.75	2.60	1.71	15.71	14.36	6.63	9.53
$\mathcal{A}_{\text{ls6}}$	2.43	4.65	1.54	1.43	15.71	14.36	6.63	9.53
$\mathcal{A}_{\text{ls8}}$	2.17	4.61	1.24	1.67	15.71	14.36	40205.16	23745.28
$\mathcal{A}_{\text{ls10}}$	2.12	4.49	1.20	1.12	15.71	14.36	40205.16	532.13
$\mathcal{A}_{\text{ls12}}$	2.10	4.20	1.12	1.09	15.71	14.36	40205.16	532.13
$\mathcal{A}_{\text{ls14}}$	2.00	1.21	1.13	1.11	15.71	14.36	40156.68	532.13
$\mathcal{A}_{\text{ls16}}$	1.04	1.20	1.03	1.10	27.63	134.12	40156.68	532.13
$\mathcal{A}_{\text{ls18}}$	1.02	1.18	1.02	1.06	27.63	134.12	40156.68	532.13

### F. Comparison of algorithm portfolios (RQ6)

Table V shows the mean relSP1 values of the VBS and the SBS in the 14 portfolios. While the other sections calculate the relSP1 value based on each portfolio, only this section calculates the relSP1 value based on the union of all the 14 portfolios  $\mathcal{A}_{\text{kt}} \cup \dots \cup \mathcal{A}_{\text{ls18}}$ . For each function, the best SP1 value for the relSP1 calculation is obtained from the results of all optimizers in the 14 portfolios (see Table S.9). For this reason, the relSP1 value of even the VBS is not 1.

As shown in Table V, a larger-size portfolio achieves better VBS performance. While the VBS performance of  $\mathcal{A}_{\text{ls18}}$  is the best for any  $n$ , that in  $\mathcal{A}_{\text{bmtp}}$  and  $\mathcal{A}_{\text{mk}}$  is the worst for  $n \in \{2, 3\}$  and  $n \in \{5, 10\}$ , respectively. For  $n \in \{5, 10\}$ , HCMA is the SBS in  $\mathcal{A}_{\text{kt}}$ ,  $\mathcal{A}_{\text{dlvat}}$ ,  $\mathcal{A}_{\text{jped}}$ ,  $\mathcal{A}_{\text{ls2}}$ ,  $\mathcal{A}_{\text{ls4}}$ , and  $\mathcal{A}_{\text{ls6}}$ . HCMA is also the best optimizer in the union of the 14 portfolios. Since  $\mathcal{A}_{\text{bmtp}}$ ,  $\mathcal{A}_{\text{mk}}$ ,  $\mathcal{A}_{\text{ls8}}$ , ...,  $\mathcal{A}_{\text{ls18}}$  do not include HCMA, their SBS performance is poor for  $n \in \{5, 10\}$ . These results indicate that optimizing the VBS performance of  $\mathcal{A}$  does not always mean optimizing the SBS performance of  $\mathcal{A}$ .

Table VI shows performance score values of the pairwise classification-based system with the 14 portfolios for the three cross-validation methods. We do not show the results for the RI-CV, which are similar to the results for the LOIO-CV. As in Tables III and IV, the best and second-best results are highlighted in **dark gray** and **gray**, respectively. Tables S.23–S.27 show the results of the five systems, respectively. With some exceptions, the results of the other four systems are relatively similar to the results in Table VI. Tables S.28–S.32 also show results of the Friedman test.

As seen from Table VI, the effectiveness of portfolios depends on cross-validation methods. For the LOIO-CV, the system with the three variants of  $\mathcal{A}_{\text{kt}}$  ( $\mathcal{A}_{\text{kt}}$ ,  $\mathcal{A}_{\text{dlvat}}$ , and  $\mathcal{A}_{\text{jped}}$ ) performs the best for any  $n$ . For the LOPO-CV,  $\mathcal{A}_{\text{ls2}}$  and  $\mathcal{A}_{\text{ls4}}$  are the most effective for the system. For the LOPOAD-CV, the system performs well when using the three variants of  $\mathcal{A}_{\text{kt}}$ . In addition,  $\mathcal{A}_{\text{ls2}}$  and  $\mathcal{A}_{\text{ls4}}$  are still effective for  $n \in \{5, 10\}$ . These results suggest that a small-size portfolio including HCMA is effective when function instances in the training and testing phases are significantly different. Otherwise, the

TABLE VI: Results of the pairwise classification-based algorithm selection systems with the 14 algorithm portfolios for  $n \in \{2, 3, 5, 10\}$ . Tables (a)–(c) show the performance score values for the three cross-validation methods, respectively.

(a) LOIO-CV					(b) LOPO-CV					(c) LOPOAD-CV				
$n$					$n$					$n$				
	2	3	5	10		2	3	5	10		2	3	5	10
$\mathcal{A}_{kt}$	0	0	0	0	$\mathcal{A}_{kt}$	0	5	1	3	$\mathcal{A}_{kt}$	0	1	0	1
$\mathcal{A}_{divat}$	0	0	0	0	$\mathcal{A}_{divat}$	1	0	5	3	$\mathcal{A}_{divat}$	1	1	6	1
$\mathcal{A}_{jped}$	0	0	0	0	$\mathcal{A}_{jped}$	0	3	1	3	$\mathcal{A}_{jped}$	0	1	0	1
$\mathcal{A}_{bmtmp}$	8	7	12	6	$\mathcal{A}_{bmtmp}$	8	11	5	6	$\mathcal{A}_{bmtmp}$	10	0	7	6
$\mathcal{A}_{mk}$	6	6	12	6	$\mathcal{A}_{mk}$	1	4	9	7	$\mathcal{A}_{mk}$	8	5	11	7
$\mathcal{A}_{is2}$	9	5	3	4	$\mathcal{A}_{is2}$	0	0	0	0	$\mathcal{A}_{is2}$	2	5	0	0
$\mathcal{A}_{is4}$	9	7	4	0	$\mathcal{A}_{is4}$	0	1	0	1	$\mathcal{A}_{is4}$	3	8	0	1
$\mathcal{A}_{is6}$	11	9	5	4	$\mathcal{A}_{is6}$	0	5	1	1	$\mathcal{A}_{is6}$	0	7	0	5
$\mathcal{A}_{is8}$	4	4	11	6	$\mathcal{A}_{is8}$	5	5	5	8	$\mathcal{A}_{is8}$	2	1	6	13
$\mathcal{A}_{is10}$	4	4	6	8	$\mathcal{A}_{is10}$	10	10	12	7	$\mathcal{A}_{is10}$	2	1	5	8
$\mathcal{A}_{is12}$	6	4	6	8	$\mathcal{A}_{is12}$	10	8	12	7	$\mathcal{A}_{is12}$	2	1	5	8
$\mathcal{A}_{is14}$	0	0	6	8	$\mathcal{A}_{is14}$	5	9	6	7	$\mathcal{A}_{is14}$	1	7	5	8
$\mathcal{A}_{is16}$	12	11	6	8	$\mathcal{A}_{is16}$	10	10	10	7	$\mathcal{A}_{is16}$	10	12	7	8
$\mathcal{A}_{is18}$	12	11	6	8	$\mathcal{A}_{is18}$	9	10	6	7	$\mathcal{A}_{is18}$	10	11	7	8

variants of  $\mathcal{A}_{kt}$  are appropriate for algorithm selection systems. Based on the poor performance of the system with  $\mathcal{A}_{is16}$  and  $\mathcal{A}_{is18}$ , we can say that a portfolio of size less than 16 may be effective. Although a general rule of thumb “set the portfolio size as small as possible” has been accepted in the literature, this is the first study to show the correctness of the rule.

Although the VBS performance of  $\mathcal{A}_{bmtmp}$  and  $\mathcal{A}_{mk}$  is worst, Table VI shows that the system with  $\mathcal{A}_{bmtmp}$  and  $\mathcal{A}_{mk}$  achieves better performance score values than that with  $\mathcal{A}_{is16}$  and  $\mathcal{A}_{is18}$  in some cases. Our results also show that the system with  $\mathcal{A}_{bmtmp}$  and  $\mathcal{A}_{mk}$  performs the best in a few cases, e.g., the result for  $n = 3$  in Table VI(c). Since the VBS performance of  $\mathcal{A}_{is12}$  is better than that of  $\mathcal{A}_{kt}$  except for  $n = 3$ ,  $\mathcal{A}_{is12}$  is likely more effective than  $\mathcal{A}_{kt}$ . Unexpectedly, Table VI and Tables S.23–S.27 show that  $\mathcal{A}_{kt}$  is more effective than  $\mathcal{A}_{is12}$  except for a few cases. These results suggest that the VBS performance of a portfolio  $\mathcal{A}$  does not always represent the effectiveness of  $\mathcal{A}$ . This is due to the difficulty of selecting the best algorithm, especially for the LOPO-CV. Our observation can be a useful clue to construct effective algorithm portfolios.

**Answers to RQ6** We found that a small-size portfolio (i.e.,  $\mathcal{A}_{is2}$  and  $\mathcal{A}_{is4}$ ) is generally the most effective for the LOPO-CV. Our results showed that the effectiveness of an algorithm portfolio  $\mathcal{A}$  depends not only on its size  $|\mathcal{A}|$  but also on its components. We also demonstrated that a portfolio  $\mathcal{A}$  with high VBS performance is not always effective.

## VII. CONCLUSION

We have investigated the performance of algorithm selection systems for black-box numerical optimization. Through a benchmarking study, we have answered the six research questions (RQ1–RQ6). Our findings can contribute to the design of more efficient algorithm selection systems. For example, we showed that using SLSQP as a pre-solver can significantly improve the performance of algorithm selection systems. We found that the regression-based selection method performs well for the practical LOPO-CV. We also demonstrated that a small-size portfolio is generally effective for the LOPO-CV.

As in previous studies [11]–[16], we fixed the target value  $f_{\text{target}}$  and the number of functions and instances. We focused

only on functions with up to  $n = 10$ . We also focused only on the fixed-target scenario. In addition, there is room for analysis of algorithm selection for another type of black-box optimization, e.g., constrained black-box optimization. There is much room for investigation of the algorithm portfolio construction for real-world black-box numerical optimization. An investigation of these factors is need in future research. An analysis of the performance of algorithm selection systems on real-world applications is also another topic for future work.

We believe that our findings contribute to the standardization of a benchmarking methodology for black-box optimization. Our results showed that the best algorithm selection system depends on various factors. Since hand-tuning is difficult in practice, automatic configuration of algorithm selection systems as in AUTOFOLIO [26] is promising. It is also interesting to compare feature-based offline algorithm selection systems with online ones (e.g., [61]) and rule-based ones (e.g., [62]) in the same platform.

## ACKNOWLEDGMENT

This work was supported by Leading Initiative for Excellent Young Researchers, MEXT, Japan.

## REFERENCES

- [1] N. Hansen, A. Auger, R. Ros, S. Finck, and P. Posik, “Comparing results of 31 algorithms from the black-box optimization benchmarking BBOB-2009,” in *Genetic and Evolutionary Computation Conference (GECCO, Companion)*, 2010, pp. 1689–1696.
- [2] J. R. Rice, “The Algorithm Selection Problem,” *Adv. Comput.*, vol. 15, pp. 65–118, 1976.
- [3] M. Lindauer, J. N. van Rijn, and L. Kotthoff, “The algorithm selection competitions 2015 and 2017,” *Artif. Intell.*, vol. 272, pp. 86–100, 2019.
- [4] K. Smith-Miles, “Cross-disciplinary perspectives on meta-learning for algorithm selection,” *ACM Comput. Surv.*, vol. 41, no. 1, pp. 6:1–6:25, 2008.
- [5] L. Kotthoff, “Algorithm Selection for Combinatorial Search Problems: A Survey,” in *Data Mining and Constraint Programming - Foundations of a Cross-Disciplinary Approach*, 2016, vol. 10101, pp. 149–190.
- [6] P. Kerschke, H. H. Hoos, F. Neumann, and H. Trautmann, “Automated Algorithm Selection: Survey and Perspectives,” *Evol. Comput.*, vol. 27, no. 1, pp. 3–45, 2019.
- [7] L. Xu, F. Hutter, H. H. Hoos, and K. Leyton-Brown, “SATzilla: Portfolio-based Algorithm Selection for SAT,” *J. Artif. Intell. Res.*, vol. 32, pp. 565–606, 2008.
- [8] P. Kerschke, L. Kotthoff, J. Bossek, H. H. Hoos, and H. Trautmann, “Leveraging TSP Solver Complementarity through Machine Learning,” *Evol. Comput.*, vol. 26, no. 4, pp. 597–620, 2018.
- [9] H. H. Hoos, M. Lindauer, and T. Schaub, “claspfolio 2: Advances in algorithm selection for answer set programming,” *Theory Pract. Log. Program.*, vol. 14, no. 4-5, pp. 569–585, 2014.
- [10] A. Liefoghe, S. Vérel, B. Lacroix, A. Zavoianu, and J. A. W. McCall, “Landscape features and automated algorithm selection for multi-objective interpolated continuous optimisation problems,” in *Genetic and Evolutionary Computation Conference (GECCO)*, 2021, pp. 421–429.
- [11] B. Bischl, O. Mersmann, H. Trautmann, and M. Preuß, “Algorithm selection based on exploratory landscape analysis and cost-sensitive learning,” in *Genetic and Evolutionary Computation Conference (GECCO)*, 2012, pp. 313–320.
- [12] T. Abell, Y. Malitsky, and K. Tierney, “Features for Exploiting Black-Box Optimization Problem Structure,” in *Learning and Intelligent Optimization (LION)*, 2013, pp. 30–36.
- [13] B. Derbel, A. Liefoghe, S. Vérel, H. E. Aguirre, and K. Tanaka, “New features for continuous exploratory landscape analysis based on the SOO tree,” in *Foundations of Genetic Algorithms (FOGA)*, 2019, pp. 72–86.
- [14] P. Kerschke and H. Trautmann, “Automated Algorithm Selection on Continuous Black-Box Problems by Combining Exploratory Landscape Analysis and Machine Learning,” *Evol. Comput.*, vol. 27, no. 1, pp. 99–127, 2019.

- [15] A. Jankovic, G. Popovski, T. Eftimov, and C. Doerr, "The impact of hyper-parameter tuning for landscape-aware performance regression and algorithm selection," in *Genetic and Evolutionary Computation Conference (GECCO)*, 2021, pp. 687–696.
- [16] M. A. Muñoz and M. Kirley, "Sampling Effects on Algorithm Selection for Continuous Black-Box Optimization," *Algorithms*, vol. 14, no. 1, p. 19, 2021.
- [17] N. Hansen, S. Finck, R. Ros, and A. Auger, "Real-parameter black-box optimization benchmarking 2009: Noiseless functions definitions," INRIA, Tech. Rep., 2009.
- [18] O. Mersmann, B. Bischl, H. Trautmann, M. Preuss, C. Weihs, and G. Rudolph, "Exploratory Landscape Analysis," in *Genetic and Evolutionary Computation Conference (GECCO)*, 2011, pp. 829–836.
- [19] R. Munos, "Optimistic Optimization of a Deterministic Function without the Knowledge of its Smoothness," in *Advances in Neural Information Processing Systems (NIPS)*, 2011, pp. 783–791.
- [20] A. Auger and N. Hansen, "Performance evaluation of an advanced local search evolutionary algorithm," in *CEC*, 2005, pp. 1777–1784.
- [21] N. Hansen, A. Auger, D. Brockhoff, D. Tusar, and T. Tusar, "COCO: performance assessment," *arXiv*, vol. abs/1605.03560, 2016.
- [22] Q. Renau, C. Doerr, J. Dréo, and B. Doerr, "Exploratory Landscape Analysis is Strongly Sensitive to the Sampling Strategy," in *Parallel Problem Solving from Nature (PPSN)*, vol. 12270, 2020, pp. 139–153.
- [23] C. Cameron, H. H. Hoos, and K. Leyton-Brown, "Bias in algorithm portfolio performance evaluation," in *Proceedings of the Twenty-Fifth International Joint Conference on Artificial Intelligence (IJCAI)*. IJCAI/AAAI Press, 2016, pp. 712–719.
- [24] B. Hurlley and B. O'Sullivan, "Statistical regimes and runtime prediction," in *Proceedings of the Twenty-Fourth International Joint Conference on Artificial Intelligence (IJCAI)*. AAAI Press, 2015, pp. 318–324.
- [25] S. Kadioglu, Y. Malitsky, A. Sabharwal, H. Samulowitz, and M. Sellmann, "Algorithm Selection and Scheduling," in *Principles and Practice of Constraint Programming (CP)*, 2011, pp. 454–469.
- [26] M. Lindauer, H. H. Hoos, F. Hutter, and T. Schaub, "AutoFolio: An Automatically Configured Algorithm Selector," *J. Artif. Intell. Res.*, vol. 53, pp. 745–778, 2015.
- [27] N. Hansen, A. Auger, R. Ros, O. Mersmann, T. Tušar, and D. Brockhoff, "COCO: a platform for comparing continuous optimizers in a black-box setting," *Optim. Methods Softw.*, vol. 36, no. 1, pp. 114–144, 2021.
- [28] R. Storm and K. Price, "Differential Evolution - A Simple and Efficient Heuristic for Global Optimization over Continuous Spaces," *J. Glo. Opt.*, vol. 11, no. 4, pp. 341–359, 1997.
- [29] N. Hansen, "The CMA Evolution Strategy: A Tutorial," *arXiv*, vol. abs/1604.00772, 2016.
- [30] —, "Invariance, self-adaptation and correlated mutations and evolution strategies," in *Parallel Problem Solving from Nature (PPSN)*, vol. 1917. Springer, 2000, pp. 355–364.
- [31] B. Bischl, P. Kerschke, L. Kotthoff, M. Lindauer, Y. Malitsky, A. Fréchet, H. H. Hoos, F. Hutter, K. Leyton-Brown, K. Tierney, and J. Vanschoren, "Aslib: A benchmark library for algorithm selection," *Artif. Intell.*, vol. 237, pp. 41–58, 2016.
- [32] A. Auger, N. Hansen, J. M. P. Zepa, R. Ros, and M. Schoenauer, "Experimental comparisons of derivative free optimization algorithms," in *Symposium on Experimental Algorithms (SEA)*, 2009, pp. 3–15.
- [33] A. Jankovic and C. Doerr, "Landscape-aware fixed-budget performance regression and algorithm selection for modular CMA-ES variants," in *Genetic and Evolutionary Computation Conference (GECCO)*, 2020, pp. 841–849.
- [34] S. van Rijn, H. Wang, B. van Stein, and T. Bäck, "Algorithm configuration data mining for CMA evolution strategies," in *Proceedings of the Genetic and Evolutionary Computation Conference, GECCO 2017, Berlin, Germany, July 15-19, 2017*, P. A. N. Bosman, Ed. ACM, 2017, pp. 737–744.
- [35] M. López-Ibáñez and T. Stützle, "Automatically improving the anytime behaviour of optimisation algorithms," *Eur. J. Oper. Res.*, vol. 235, no. 3, pp. 569–582, 2014.
- [36] A. D. Jesus, A. Liefvooghe, B. Derbel, and L. Paquete, "Algorithm selection of anytime algorithms," in *GECCO '20: Genetic and Evolutionary Computation Conference, Cancún Mexico, July 8-12, 2020*, C. A. C. Coello, Ed. ACM, 2020, pp. 850–858.
- [37] C. Doerr, F. Ye, N. Horesh, H. Wang, O. M. Shir, and T. Bäck, "Benchmarking discrete optimization heuristics with IOHprofiler," *Appl. Soft Comput.*, vol. 88, p. 106027, 2020.
- [38] P. Kerschke and H. Trautmann, "Comprehensive Feature-Based Landscape Analysis of Continuous and Constrained Optimization Problems Using the R-package flacco," in *Applications in Statistical Computing – From Music Data Analysis to Industrial Quality Improvement*. Springer, 2019, pp. 93–123.
- [39] M. Lunacek and D. Whitley, "The dispersion metric and the CMA evolution strategy," in *Genetic and Evolutionary Computation Conference (GECCO)*, 2006, pp. 477–484.
- [40] M. A. Muñoz and M. Kirley, "ICARUS: Identification of complementary algorithms by uncovered sets," in *IEEE Congress on Evolutionary Computation (CEC)*, 2016, pp. 2427–2432.
- [41] K. Tang, F. Peng, G. Chen, and X. Yao, "Population-based Algorithm Portfolios with automated constituent algorithms selection," *Inf. Sci.*, vol. 279, pp. 94–104, 2014.
- [42] L. Xu, H. H. Hoos, and K. Leyton-Brown, "Hydra: Automatically configuring algorithms for portfolio-based selection," in *Proceedings of the Twenty-Fourth AAAI Conference on Artificial Intelligence, AAAI 2010, Atlanta, Georgia, USA, July 11-15, 2010*, M. Fox and D. Poole, Eds. AAAI Press, 2010.
- [43] L. Xu, F. Hutter, H. H. Hoos, and K. Leyton-Brown, "Hydra-MIP: Automated Algorithm Configuration and Selection for Mixed Integer Programming," in *Proceedings of the RCRA workshop at IJCAI*, 2011.
- [44] L. Kotthoff, "LLAMA: leveraging learning to automatically manage algorithms," *arXiv*, vol. abs/1306.1031, 2013.
- [45] S. Kadioglu, Y. Malitsky, M. Sellmann, and K. Tierney, "ISAC - Instance-Specific Algorithm Configuration," in *European Conference on Artificial Intelligence (ECAI)*, 2010, pp. 751–756.
- [46] G. Hamerly and C. Elkan, "Learning the  $k$  in  $k$ -means," in *Advances in Neural Information Processing Systems (NIPS)*, 2003, pp. 281–288.
- [47] M. Basseur, B. Derbel, A. Goëffon, and A. Liefvooghe, "Experiments on Greedy and Local Search Heuristics for  $d$ -dimensional Hypervolume Subset Selection," in *Genetic and Evolutionary Computation Conference (GECCO)*, 2016, pp. 541–548.
- [48] D. Kraft, "A Software Package for Sequential Quadratic Programming," Deutsche Forschungs- und Versuchsanstalt für Luft- und Raumfahrt, Tech. Rep. DFVLR-FB 88-28, 1988.
- [49] F. Hutter, F. H. Hoos, and K. Leyton-Brown, "Sequential Model-Based Optimization for General Algorithm Configuration," in *Conference on Learning and Intelligent Optimization (LION)*, 2011, pp. 507–523.
- [50] N. Hansen, "A global surrogate assisted CMA-ES," in *Genetic and Evolutionary Computation Conference (GECCO)*, 2019, pp. 664–672.
- [51] M. A. Muñoz, M. Kirley, and S. K. Halgamuge, "A Meta-learning Prediction Model of Algorithm Performance for Continuous Optimization Problems," in *Parallel Problem Solving from Nature (PPSN)*, 2012, pp. 226–235.
- [52] A. Jankovic, T. Eftimov, and C. Doerr, "Towards Feature-Based Performance Regression Using Trajectory Data," in *Applications of Evolutionary Computation (EvoApplications 2021)*, 2021, pp. 601–617.
- [53] B. Beachkofski and R. Grandhi, "Improved Distributed Hypercube Sampling," in *AIAA/ASME/ASCE/AHS/ASC Structures, Structural Dynamics, and Materials Conference*, 2002, p. 1274.
- [54] P. Kerschke, M. Preuss, S. Wessing, and H. Trautmann, "Detecting Funnel Structures by Means of Exploratory Landscape Analysis," in *Genetic and Evolutionary Computation Conference (GECCO)*, 2015, pp. 265–272.
- [55] M. Lindauer, K. Eggensperger, M. Feurer, A. Biedenkapp, D. Deng, C. Benjamins, T. Ruhkopf, R. Sass, and F. Hutter, "SMAC3: A Versatile Bayesian Optimization Package for Hyperparameter Optimization," *J. Mach. Learn. Res.*, vol. 23, pp. 54:1–54:9, 2022.
- [56] L. Breiman, "Random forests," *Mach. Learn.*, vol. 45, no. 1, pp. 5–32, 2001.
- [57] J. Bader and E. Zitzler, "HypE: An Algorithm for Fast Hypervolume-Based Many-Objective Optimization," *Evol. Comput.*, vol. 19, no. 1, pp. 45–76, 2011.
- [58] J. Derrac, S. García, D. Molina, and F. Herrera, "A practical tutorial on the use of nonparametric statistical tests as a methodology for comparing evolutionary and swarm intelligence algorithms," *Swarm Evol. Comput.*, vol. 1, no. 1, pp. 3–18, 2011.
- [59] Y. He, S. Y. Yuen, and Y. Lou, "Exploratory landscape analysis using algorithm based sampling," in *Genetic and Evolutionary Computation Conference (GECCO, Companion)*, 2018, pp. 211–212.
- [60] B. van Stein, H. Wang, and T. Bäck, "Neural Network Design: Learning from Neural Architecture Search," in *IEEE Symposium Series on Computational Intelligence*, 2020, pp. 1341–1349.
- [61] P. Baudiš and P. Pošík, "Online Black-Box Algorithm Portfolios for Continuous Optimization," in *Parallel Problem Solving from Nature (PPSN)*, 2014, pp. 40–49.
- [62] J. Liu, A. Moreau, M. Preuss, J. Rapin, B. Rozière, F. Teytaud, and O. Teytaud, "Versatile black-box optimization," in *Genetic and Evolutionary Computation Conference (GECCO)*, 2020, pp. 620–628.
- [63] M. A. Muñoz, M. Kirley, and S. K. Halgamuge, "Exploratory Landscape

- Analysis of Continuous Space Optimization Problems Using Information Content,” *IEEE Trans. Evol. Comput.*, vol. 19, no. 1, pp. 74–87, 2015.
- [64] P. Kerschke, M. Preuss, C. Hernández, O. Schütze, J. Sun, C. Grimme, G. Rudolph, B. Bischl, and H. Trautmann, “Cell mapping techniques for exploratory landscape analysis,” in *EVOLVE-A*, 2014, pp. 115–131.
- [65] P. Posik and P. Baudis, “Dimension Selection in Axis-Parallel Brent-STEP Method for Black-Box Optimization of Separable Continuous Functions,” in *Genetic and Evolutionary Computation Conference (GECCO, Companion)*, 2015, pp. 1151–1158.
- [66] L. Pál, “Comparison of multistart global optimization algorithms on the BBOB noiseless testbed,” in *Genetic and Evolutionary Computation Conference (GECCO, Companion)*, 2013, pp. 1153–1160.
- [67] —, “Benchmarking a hybrid multi level single linkage algorithm on the bbob noiseless testbed,” in *Genetic and Evolutionary Computation Conference (GECCO, Companion)*, 2013, pp. 1145–1152.
- [68] W. Huyer and A. Neumaier, “Benchmarking of MCS on the Noiseless Function Testbed,” Tech. Rep., 2009.
- [69] A. Auger, D. Brockhoff, and N. Hansen, “Benchmarking the local meta-model CMA-ES on the noiseless BBOB’2013 test bed,” in *Genetic and Evolutionary Computation Conference (GECCO, Companion)*, 2013, pp. 1225–1232.
- [70] I. Loshchilov, M. Schoenauer, and M. Sebag, “Bi-population CMA-ES algorithms with surrogate models and line searches,” in *Genetic and Evolutionary Computation Conference (GECCO, Companion)*, 2013, pp. 1177–1184.
- [71] A. Atamna, “Benchmarking IPOP-CMA-ES-TPA and IPOP-CMA-ES-MSR on the BBOB Noiseless Testbed,” in *Genetic and Evolutionary Computation Conference (GECCO, Companion)*, 2015, pp. 1135–1142.
- [72] F. Hutter, H. H. Hoos, and K. Leyton-Brown, “An evaluation of sequential model-based optimization for expensive blackbox functions,” in *Genetic and Evolutionary Computation Conference (GECCO, Companion)*, 2013, pp. 1209–1216.
- [73] N. Hansen, “Benchmarking a BI-population CMA-ES on the BBOB-2009 function testbed,” in *Genetic and Evolutionary Computation Conference (GECCO, Companion)*, 2009, pp. 2389–2396.
- [74] R. Ros, “Benchmarking the BFGS algorithm on the BBOB-2009 function testbed,” in *Genetic and Evolutionary Computation Conference (GECCO, Companion)*, 2009, pp. 2409–2414.
- [75] P. Posik, “Bbob-benchmarking two variants of the line-search algorithm,” in *Genetic and Evolutionary Computation Conference (GECCO, Companion)*, 2009, pp. 2329–2336.
- [76] A. Auger, D. Brockhoff, and N. Hansen, “Comparing the (1+1)-CMA-ES with a mirrored (1+2)-CMA-ES with sequential selection on the noiseless BBOB-2010 testbed,” in *Genetic and Evolutionary Computation Conference (GECCO, Companion)*, 2010, pp. 1543–1550.
- [77] B. Doerr, M. Fouz, M. Schmidt, and M. Wahlström, “BBOB: Nelder-Mead with resize and halfruns,” in *Genetic and Evolutionary Computation Conference (GECCO, Companion)*, 2009, pp. 2239–2246.
- [78] T. Glasmachers and O. Krause, “The Hessian Estimation Evolution Strategy,” in *Parallel Problem Solving from Nature (PPSN)*, 2020, pp. 597–609.
- [79] D. M. Nguyen, “Benchmarking a variant of the CMAES-APOP on the BBOB noiseless testbed,” in *Genetic and Evolutionary Computation Conference (GECCO, Companion)*, 2018, pp. 1521–1528.
- [80] L. Bajer, Z. Pitra, J. Repický, and M. Holena, “Gaussian Process Surrogate Models for the CMA Evolution Strategy,” *Evol. Comput.*, vol. 27, no. 4, pp. 665–697, 2019.
- [81] Z. Pitra, L. Bajer, J. Repický, and M. Holena, “Comparison of ordinal and metric gaussian process regression as surrogate models for CMA evolution strategy,” in *Genetic and Evolutionary Computation Conference (GECCO, Companion)*, 2017, pp. 1764–1771.
- [82] A. LaTorre, S. Muelas, and J. M. Peña, “Benchmarking a MOS-based algorithm on the BBOB-2010 noiseless function testbed,” in *Genetic and Evolutionary Computation Conference (GECCO, Companion)*, 2010, pp. 1649–1656.
- [83] K. Nishida and Y. Akimoto, “Benchmarking the PSA-CMA-ES on the BBOB noiseless testbed,” in *Genetic and Evolutionary Computation Conference (GECCO, Companion)*, 2018, pp. 1529–1536.
- [84] R. Ros, “Comparison of NEWUOA with different numbers of interpolation points on the BBOB noisy testbed,” in *Genetic and Evolutionary Computation Conference (GECCO, Companion)*, 2010, pp. 1495–1502.
- [85] T. Tran, D. Brockhoff, and B. Derbel, “Multiobjectivization with NSGA-II on the noiseless BBOB testbed,” in *Genetic and Evolutionary Computation Conference (GECCO, Companion)*, 2013, pp. 1217–1224.



**Ryoji Tanabe** is an Assistant Professor at Yokohama National University, Japan (2019–). Previously, he was a Research Assistant Professor at Southern University of Science and Technology, China (2017–2019). He was also a Post-Doctoral Researcher at Japan Aerospace Exploration Agency, Japan (2016–2017). He received the Ph.D. degree in Science from The University of Tokyo, Japan, in 2016. His research interests include single- and multi-objective black-box optimization, analysis of evolutionary algorithms, and automatic algorithm configuration.

## SUPPLEMENT

TABLE S.7: Experimental settings in the five previous studies [S.11]–[15]. This table also shows configurations of each algorithm selection system. Since the experimental setting in [S.12] is not very clear, we do not describe it here.

Ref.	Dimension $n$	S. size	Feature classes	Selector	Algorithm portfolio	Cross-validation
[S.11]	10	$500 \times n$	ela_distr, ela_level, ela_meta, ela_conv, ela_local	Classif.	$\mathcal{A}_{\text{bmtp}}$	LOIO-CV, LOPO-CV
[S.14]	2, 3, 5, 10	$50 \times n$	Features in flacco	Classif. Reg. Pairwise reg.	$\mathcal{A}_{\text{kt}}$	LOPOAD-CV
[S.13]	2, 3, 5	$100 \times n$	13 feature classes in flacco, SOO-based features	Reg.	$\mathcal{A}_{\text{dlvat}}$	LOPO-CV
[S.15]	5	$400 \times n$	ela_distr, ela_level ela_meta, nbc, disp, ic	Reg.	$\mathcal{A}_{\text{jped}}$	LOIO-CV
[S.16]	2, 5, 10, 20	$100 \times n$ – $1000 \times n$	ela_distr, ela_level ela_meta, nbc, disp, ic	Pairwise Classif.	$\mathcal{A}_{\text{mk}}$	LOIO-CV, LOPO-CV

TABLE S.8: 17 feature classes provided by flacco.

Feature class	Name	Num. features
ela_conv [S.18]	convexity	6
ela_curv [S.18]	curvature	26
ela_local [S.18]	local search	16
ela_distr [S.18]	$y$ -distribution	5
ela_level [S.18]	levelset	20
ela_meta [S.18]	meta-model	11
nbc [S.54]	nearest better clustering (NBC)	7
disp [S.54]	dispersion	18
ic [S.63]	information content	7
basic [S.38]	basic	15
limo [S.38]	linear model	14
pca [S.38]	principal component analysis	10
cm_angle [S.64]	cell mapping angle	10
cm_conv [S.64]	cell mapping convexity	6
cm_grad [S.64]	cell mapping gradient homog.	6
gcm [S.64]	generalized cell mapping	75
bt [S.38]	barrier tree	90

---

**Algorithm S.1:** First-improvement local search method [S.47]

---

```
1 Initialize a subset  $\mathcal{A} \subset \mathcal{R}$  of size  $k$ ;  
2 while there is a pair to improve a quality measure  $m$  do  
3   for  $a' \in \mathcal{R} \setminus \mathcal{A}$  do  
4     for  $a \in \mathcal{A}$  do  
5        $\mathcal{A}^{\text{new}} \leftarrow \mathcal{A} \setminus \{a\} \cup \{a'\}$ ;  
6       if  $m(\mathcal{A}^{\text{new}}) < m(\mathcal{A})$  then  
7          $\mathcal{A} \leftarrow \mathcal{A}^{\text{new}}$ ;
```

---

TABLE S.9: 14 algorithm portfolios used in this study. Note that “DTS-CMA-ES18” is “DTS-CMA-ES\_005-2pop\_v26\_1model”.

Portfolio	Algorithms
$\mathcal{A}_{kt}$ [S.14]	BrentSTEPrr [S.65], BrentSTEPqi [S.65], fmincon [S.66], fminunc [S.66], MLSL [S.67], HMLS [S.67], MCS [S.68], IPOP400D [S.69], HCMA [S.70], CMA-CSA [S.71], SMAC-BBOB [S.72], OQNLP [S.66]
$\mathcal{A}_{dlvat}$ [S.13]	$\mathcal{A}_{kt} \setminus \{\text{BrentSTEPqi [S.65], SMAC-BBOB [S.72]}\}$
$\mathcal{A}_{jped}$ [S.15]	$\mathcal{A}_{kt} \cup \{\text{BIPOP-CMA-ES [S.73]}\} \setminus \{\text{SMAC-BBOB [S.72]}\}$
$\mathcal{A}_{bmtp}$ [S.11]	BFGS [S.74], BIPOP-CMA-ES [S.73], LSfminbd [S.75], LSstep [S.75]
$\mathcal{A}_{mk}$ [S.16]	BIPOP-CMA-ES [S.73], lplus2mirser [S.76], LSstep [S.75], NELDERDOERR [S.77]
$\mathcal{A}_{is2}$	HCMA [S.70], HMLS [S.67]
$\mathcal{A}_{is4}$	HCMA [S.70], HMLS [S.67], BrentSTEPqi [S.65], HE-ES [S.78]
$\mathcal{A}_{is6}$	HCMA [S.70], HMLS [S.67], BrentSTEPqi [S.65], CMAES-APOP-Var1 [S.79], DTS-CMA-ES18 [S.80], HE-ES [S.78]
$\mathcal{A}_{is8}$	BIPOP-saACM-k [S.70], HMLS [S.67], BrentSTEPqi [S.65], DTS-CMA-ES [S.81], CMAES-APOP-Var1 [S.79], DTS-CMA-ES18 [S.80], HE-ES [S.78], SLSQP-11-scipy [S.50]
$\mathcal{A}_{is10}$	MOS [S.82], BIPOP-saACM-k [S.70], HMLS [S.67], SMAC-BBOB [S.72], BrentSTEPqi [S.65], DTS-CMA-ES [S.81], PSA-CMA-ES [S.83], DTS-CMA-ES18 [S.80], HE-ES [S.78], SLSQP-11-scipy [S.50]
$\mathcal{A}_{is12}$	LSstep [S.75], MOS [S.82], BIPOP-saACM-k [S.70], fmincon [S.66], HMLS [S.67], Imm-CMA-ES [S.69], SMAC-BBOB [S.72], BrentSTEPqi [S.65], PSA-CMA-ES [S.83], DTS-CMA-ES18 [S.80], HE-ES [S.78], SLSQP-11-scipy [S.50]
$\mathcal{A}_{is14}$	LSstep [S.75], MOS [S.82], BIPOP-saACM-k [S.70], fmincon [S.66], HMLS [S.67], OQNLP [S.66], SMAC-BBOB [S.72], BrentSTEPqi [S.65], DTS-CMA-ES [S.81], PSA-CMA-ES [S.83], DTS-CMA-ES18 [S.80], HE-ES [S.78], lq-CMA-ES [S.50], SLSQP-11-scipy [S.50]
$\mathcal{A}_{is16}$	LSstep [S.75], MCS [S.68], AVGNEUOA [S.84], MOS [S.82], BIPOP-saACM-k [S.70], fmincon [S.66], MLSL [S.67], OQNLP [S.66], P-DCN [S.85], BrentSTEPqi [S.65], DTS-CMA-ES [S.81], PSA-CMA-ES [S.83], DTS-CMA-ES18 [S.80], HE-ES [S.78], lq-CMA-ES [S.50], SLSQP-11-scipy [S.50]
$\mathcal{A}_{is18}$	LSstep [S.75], MCS [S.68], AVGNEUOA [S.84], MOS [S.82], BIPOP-saACM-k [S.70], fmincon [S.66], Imm-CMA-ES [S.69], MLSL [S.67], OQNLP [S.66], P-DCN [S.85], BrentSTEPqi [S.65], BrentSTEPrr [S.65], DTS-CMA-ES [S.81], PSA-CMA-ES [S.83], DTS-CMA-ES18 [S.80], HE-ES [S.78], lq-CMA-ES [S.50], SLSQP-11-scipy [S.50]

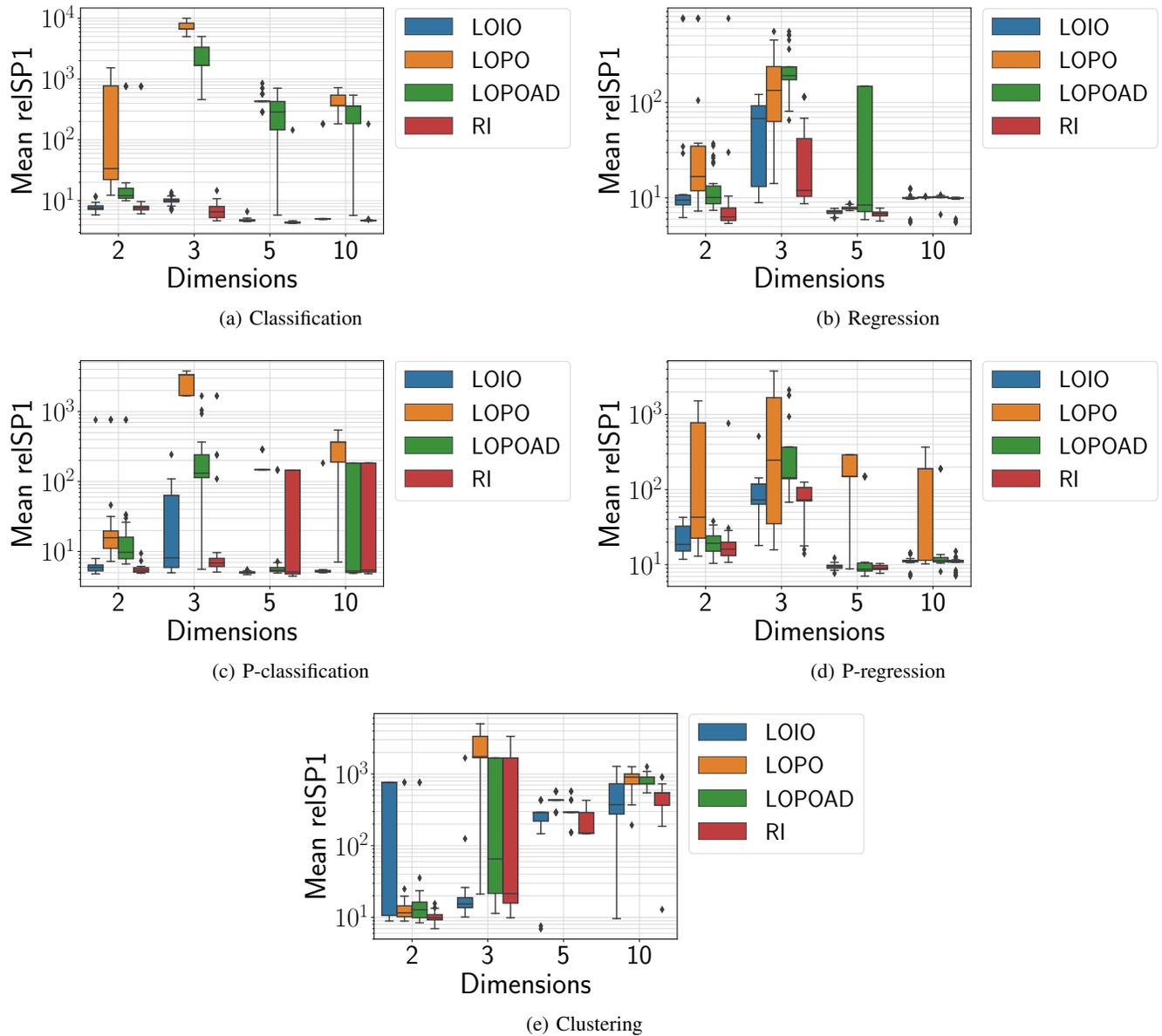


Fig. S.4: Distribution of 31 mean relSP1 values of the five algorithm selection systems with  $\mathcal{A}_{kt}$  for the LOIO-CV, the LOPO-CV, and the LOIO-CV.

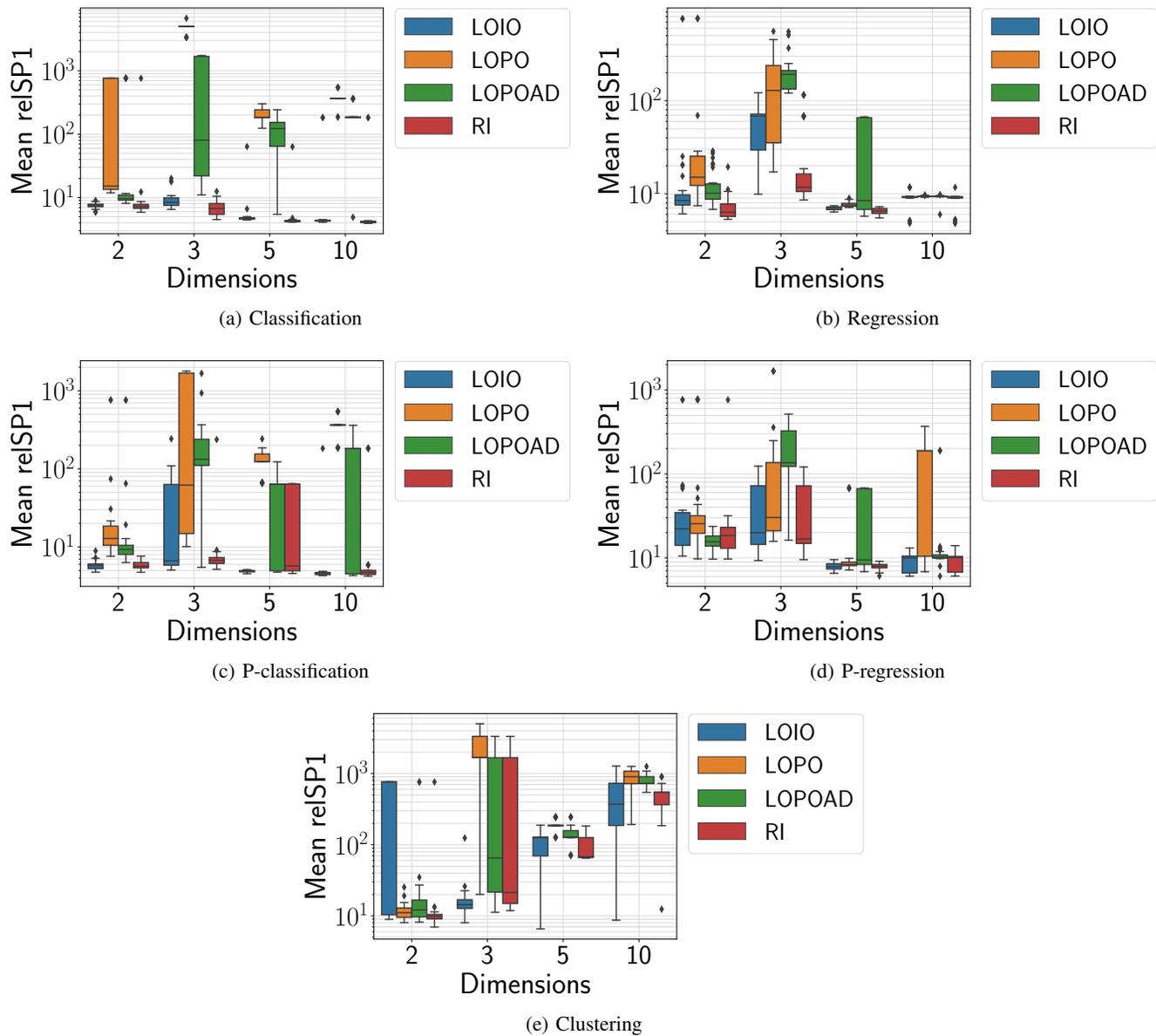


Fig. S.5: Distribution of 31 mean relSP1 values of the five algorithm selection systems with  $\mathcal{A}_{\text{divat}}$  for the LOIO-CV, the LOPO-CV, and the LOIO-CV.

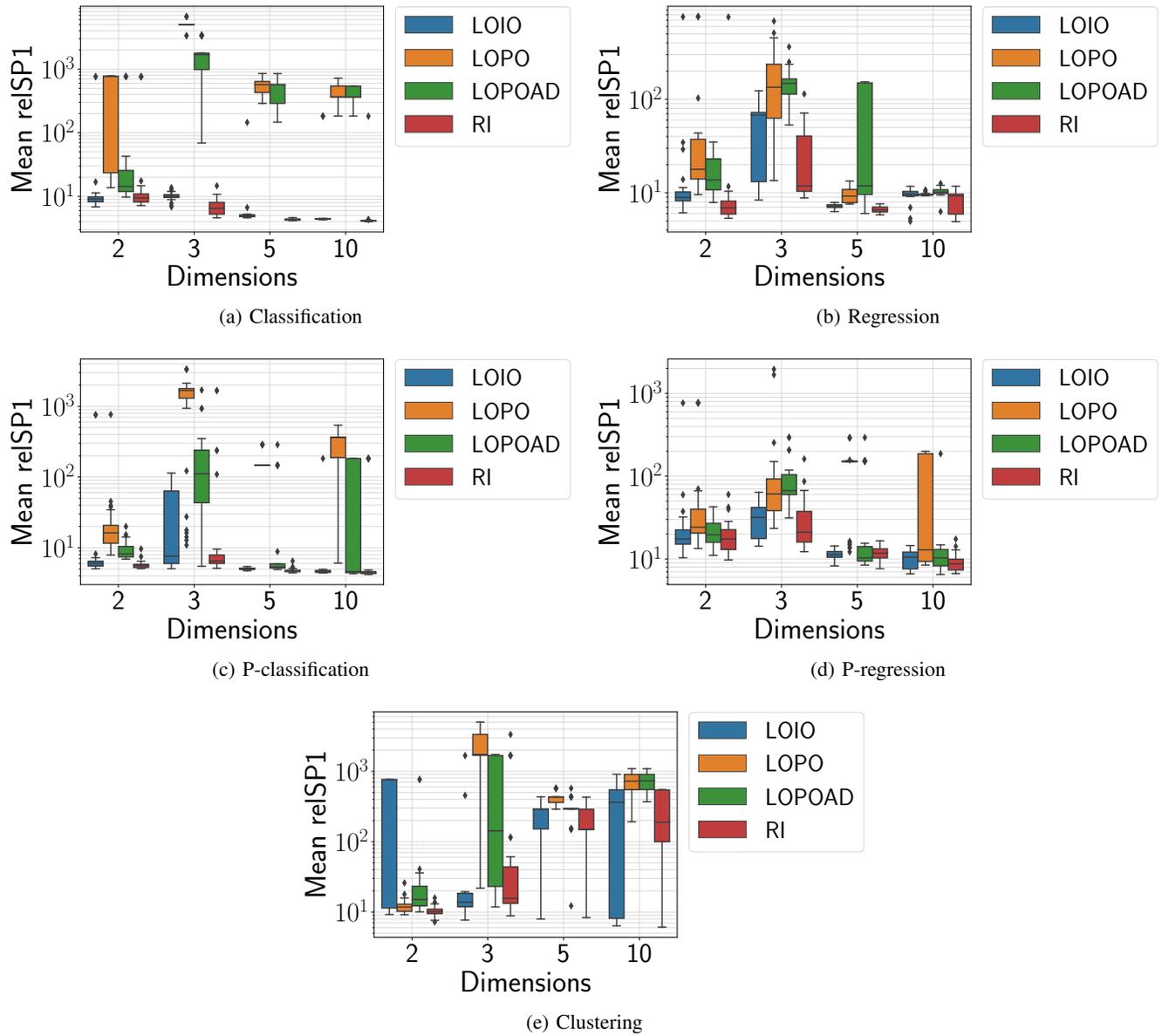


Fig. S.6: Distribution of 31 mean relSP1 values of the five algorithm selection systems with  $\mathcal{A}_{jped}$  for the LOIO-CV, the LOPO-CV, and the LOIO-CV.

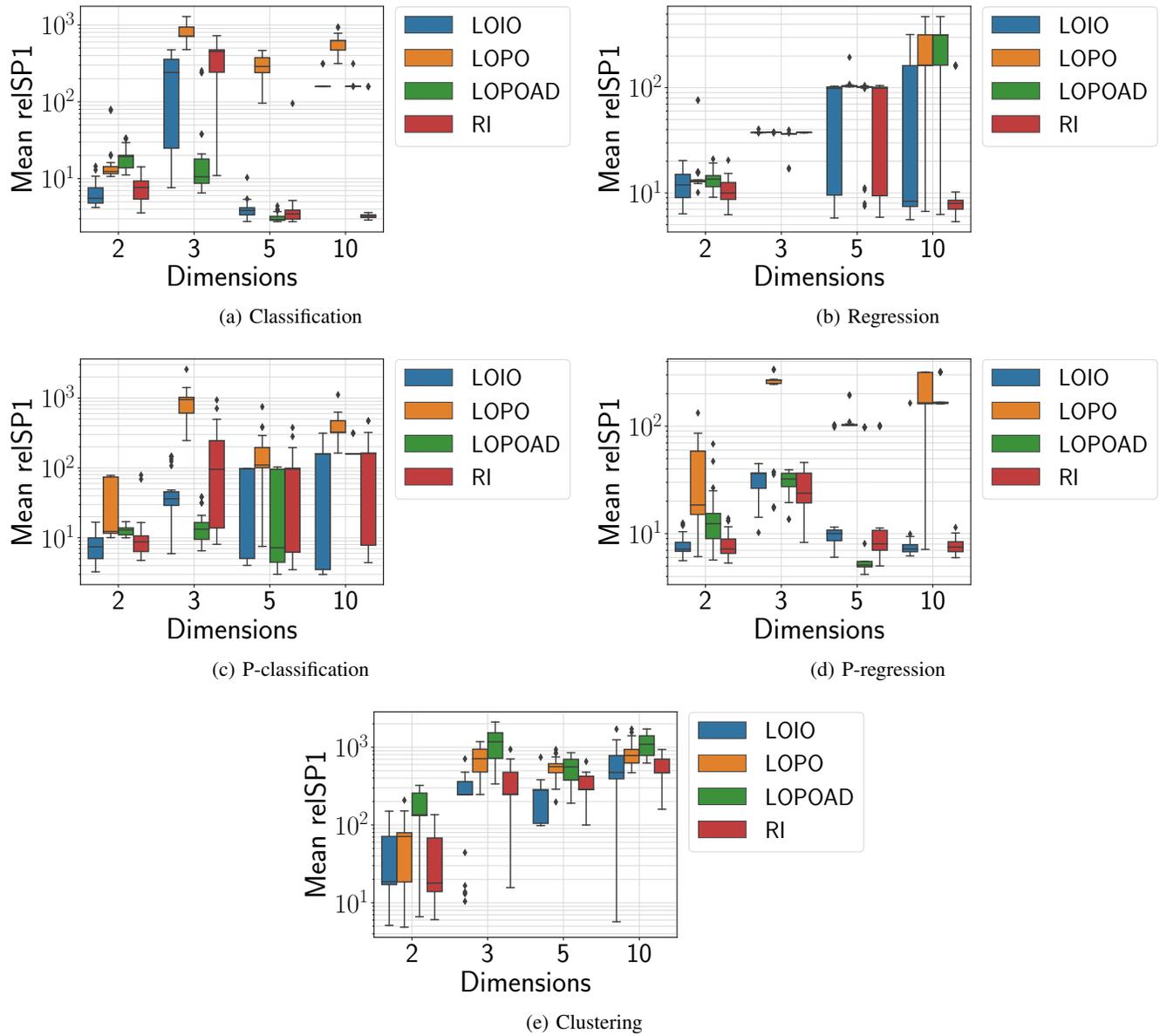


Fig. S.7: Distribution of 31 mean relSP1 values of the five algorithm selection systems with  $\mathcal{A}_{\text{bmtP}}$  for the LOIO-CV, the LOPO-CV, and the LOIO-CV.

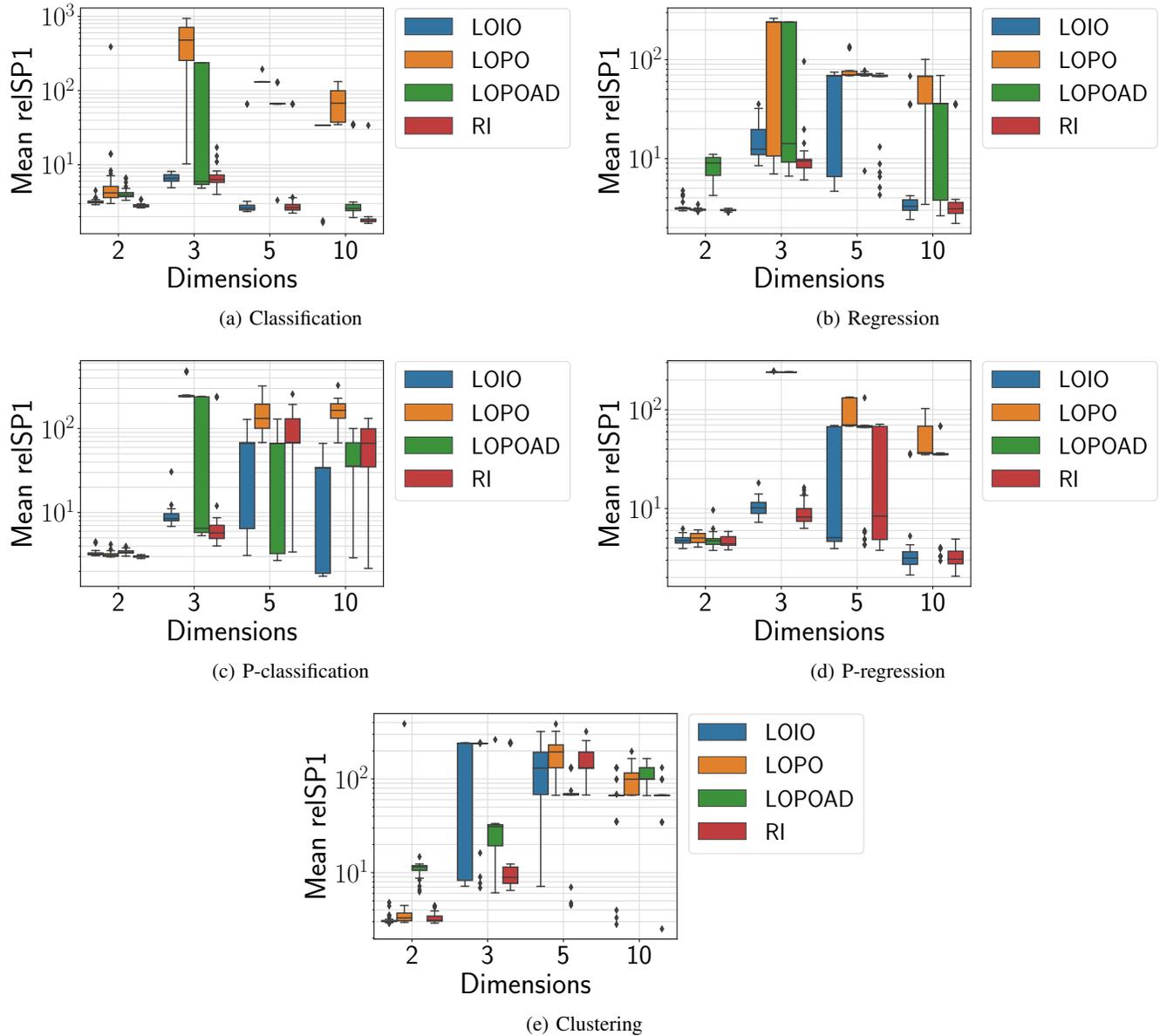


Fig. S.8: Distribution of 31 mean relSP1 values of the five algorithm selection systems with  $\mathcal{A}_{mk}$  for the LOIO-CV, the LOPO-CV, and the LOIO-CV.

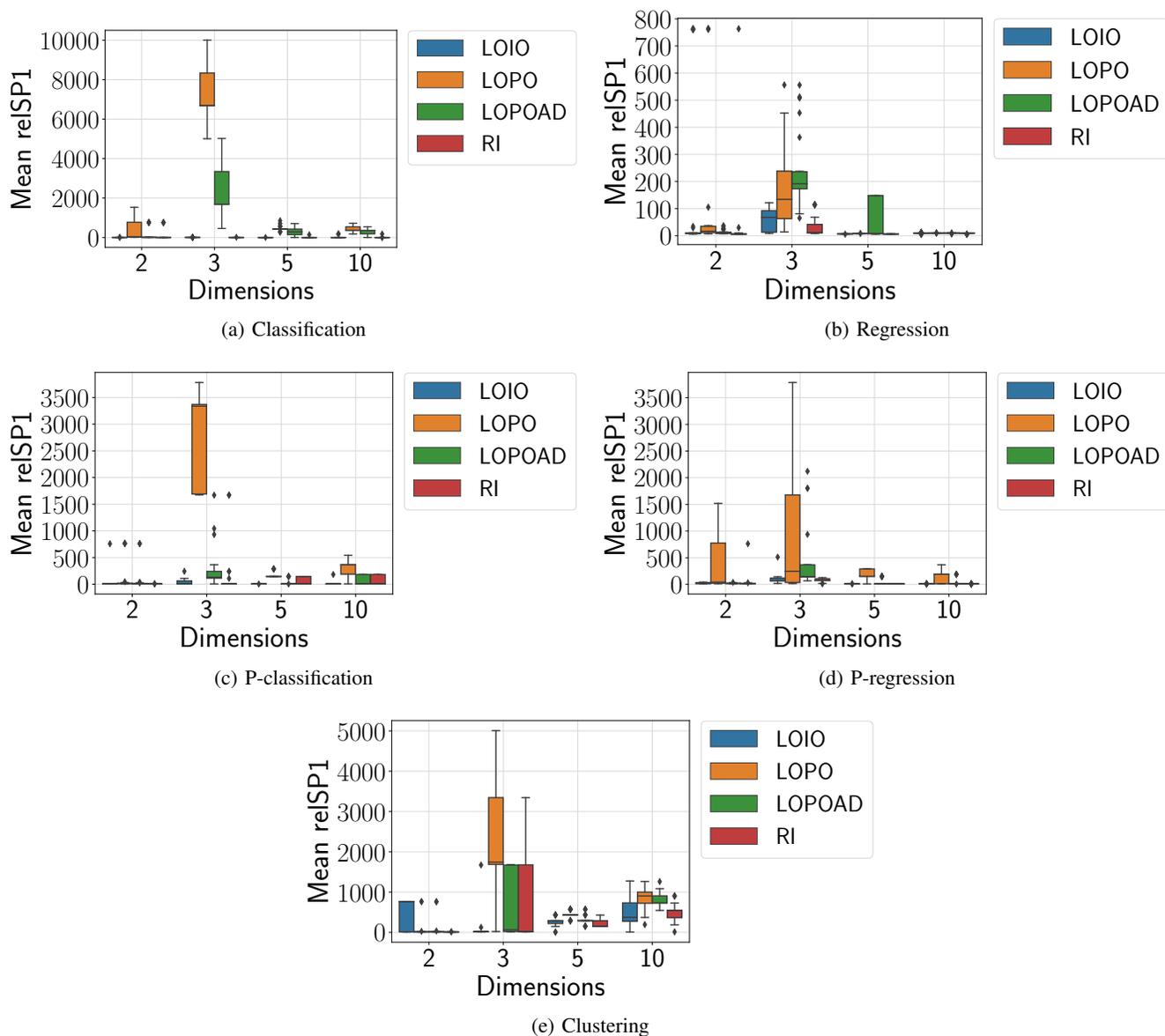


Fig. S.9: Distribution of 31 mean relSP1 values of the five algorithm selection systems with  $\mathcal{A}_{kt}$  for the LOIO-CV, the LOPO-CV, and the LOIO-CV.

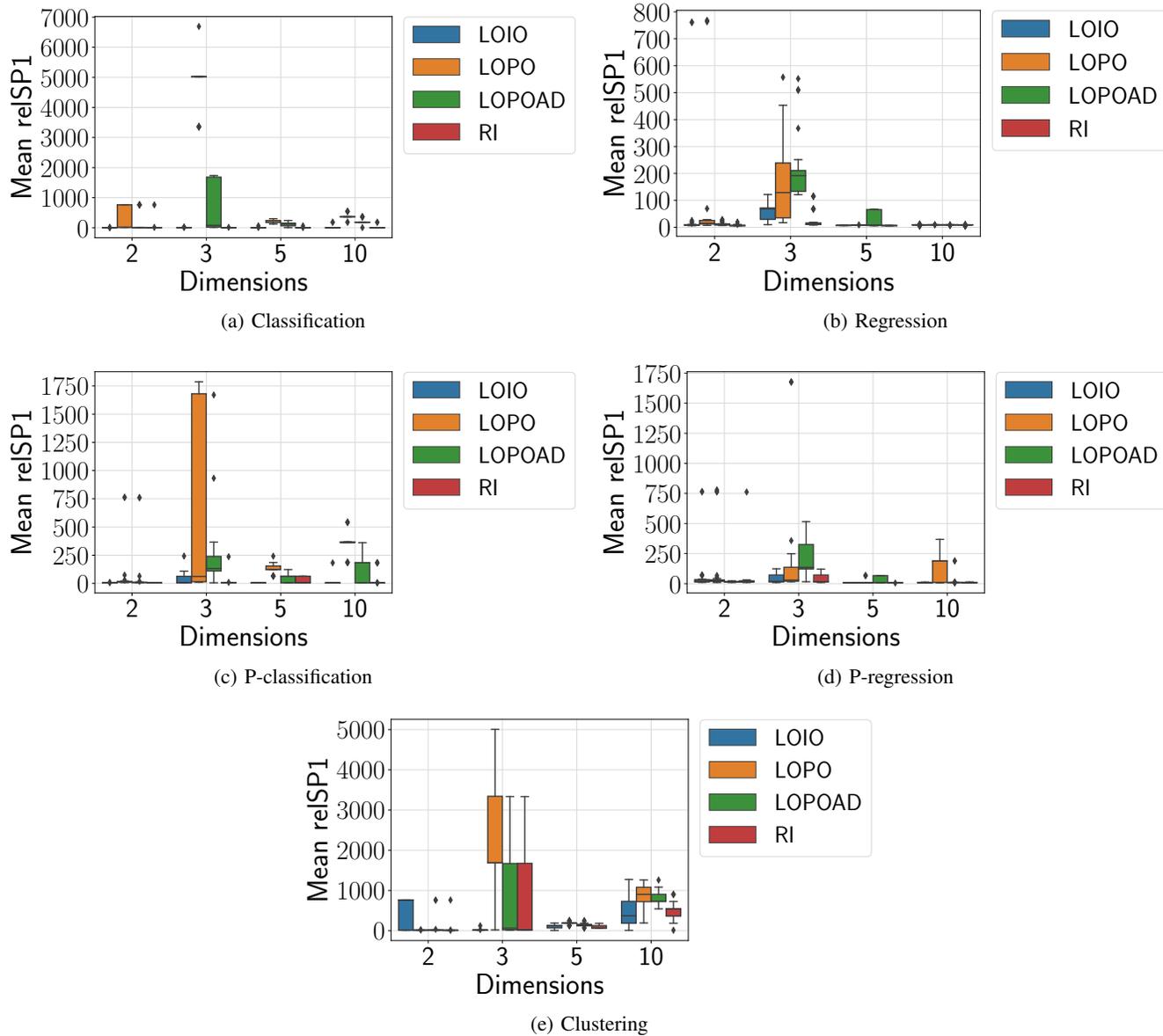


Fig. S.10: Distribution of 31 mean relSP1 values of the five algorithm selection systems with  $\mathcal{A}_{dlvat}$  for the LOIO-CV, the LOPO-CV, and the LOIO-CV.

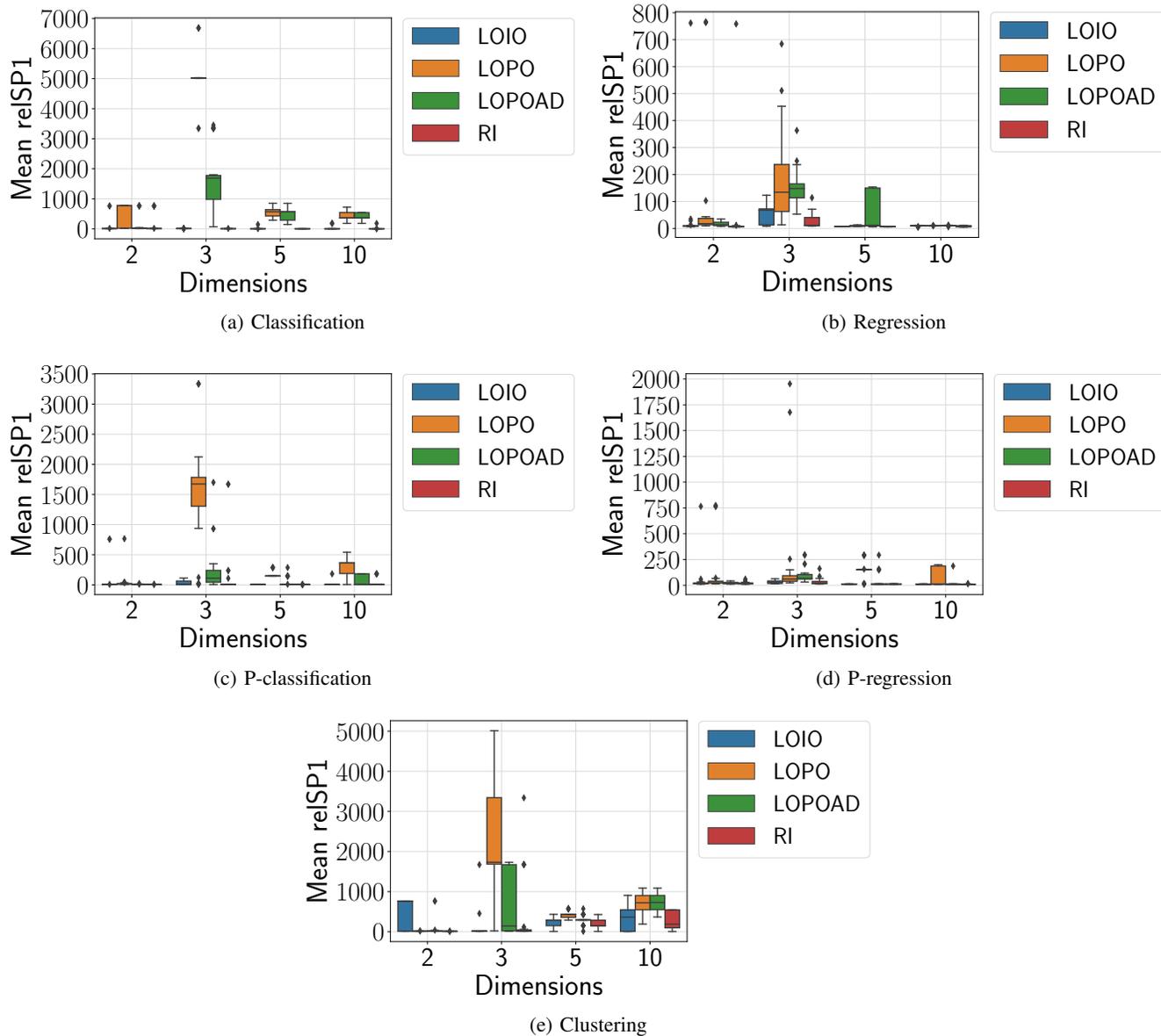


Fig. S.11: Distribution of 31 mean relSP1 values of the five algorithm selection systems with  $\mathcal{A}_{jped}$  for the LOIO-CV, the LOPO-CV, and the LOIO-CV.

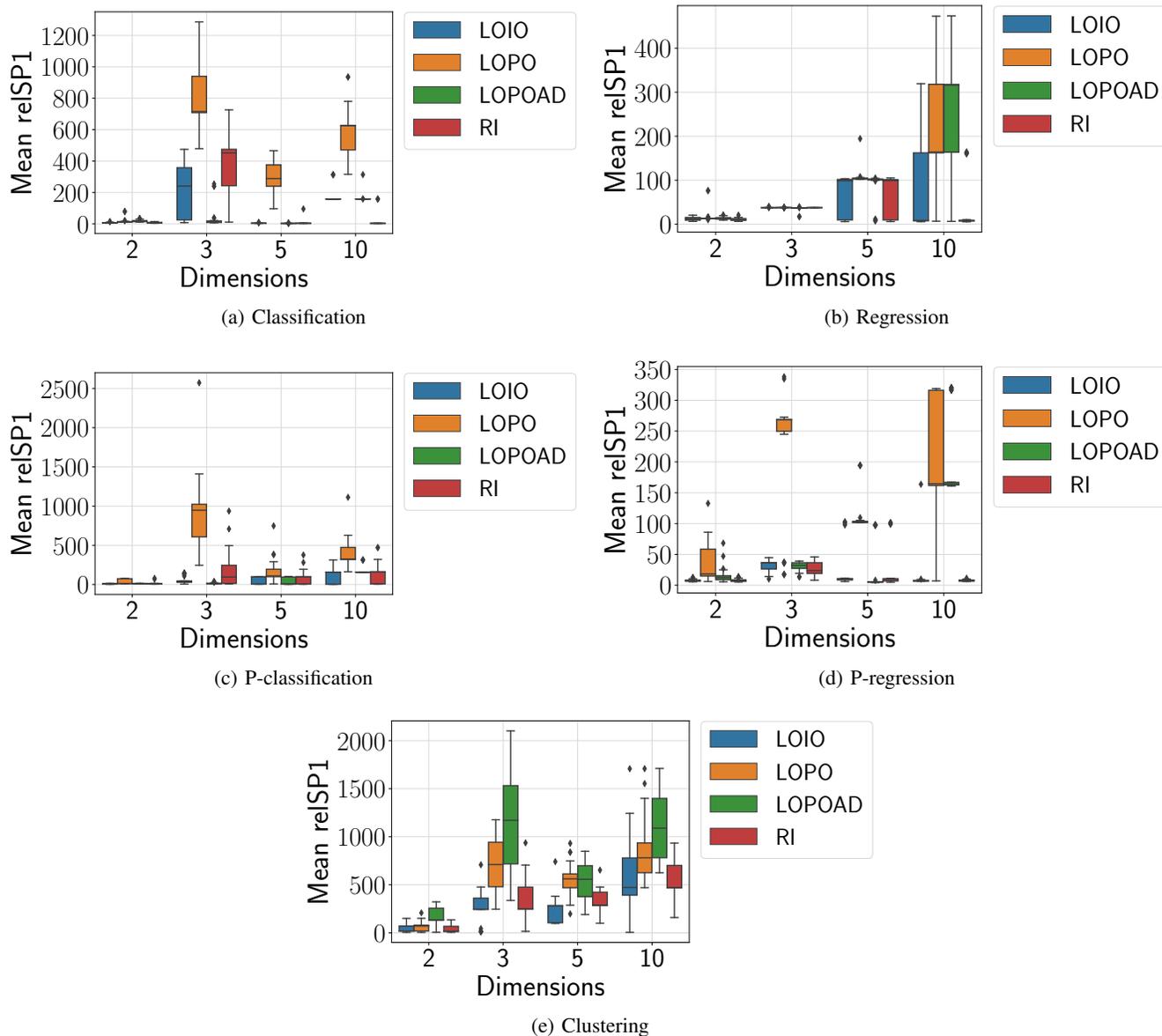


Fig. S.12: Distribution of 31 mean relSP1 values of the five algorithm selection systems with  $\mathcal{A}_{\text{bmtP}}$  for the LOIO-CV, the LOPO-CV, and the LOIO-CV.

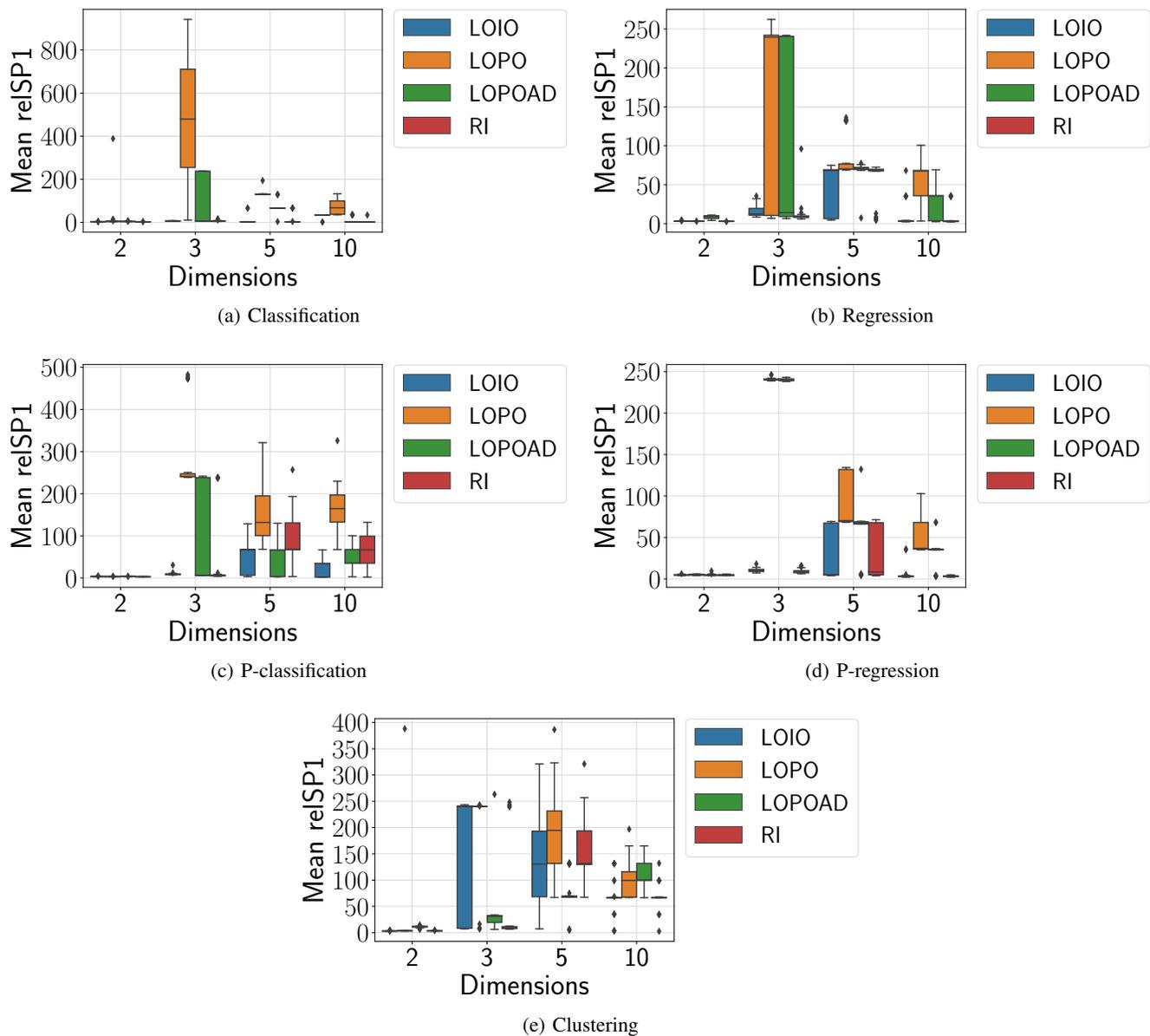


Fig. S.13: Distribution of 31 mean relSP1 values of the five algorithm selection systems with  $\mathcal{A}_{mk}$  for the LOIO-CV, the LOPO-CV, and the LOIO-CV.

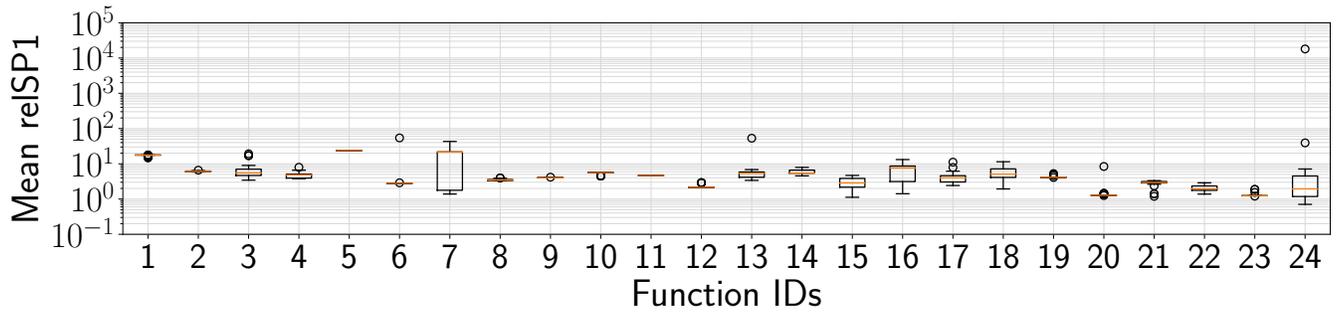
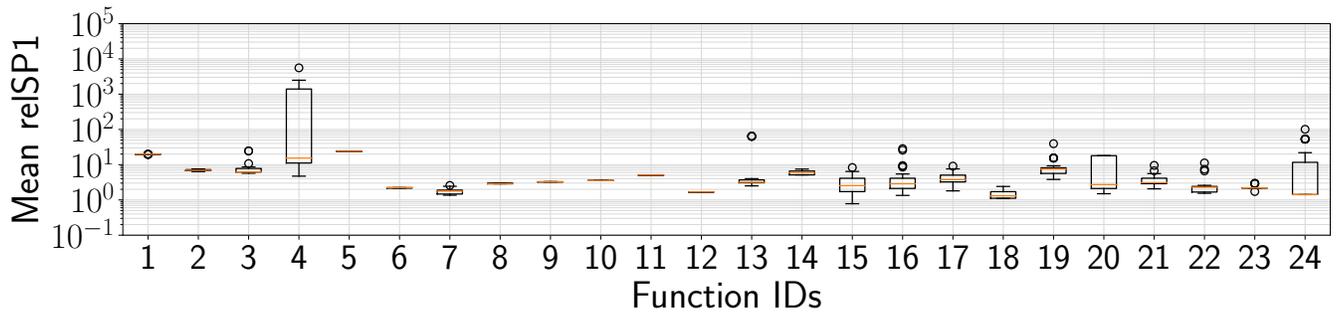
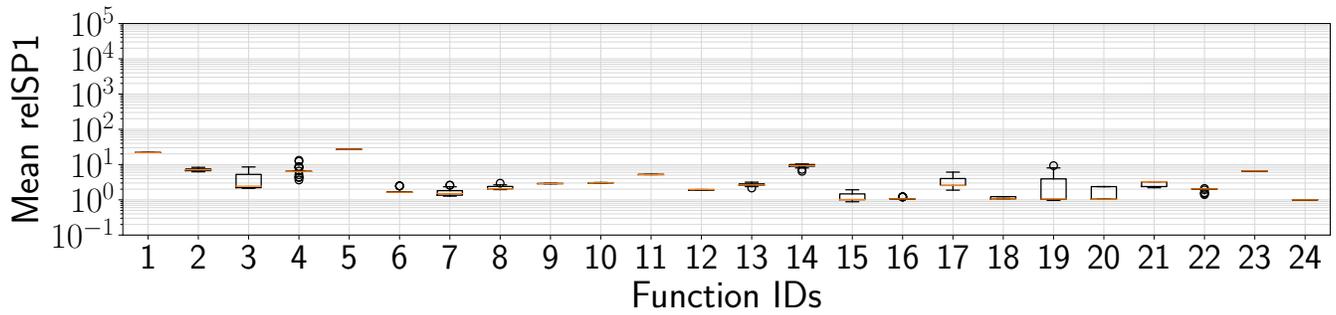
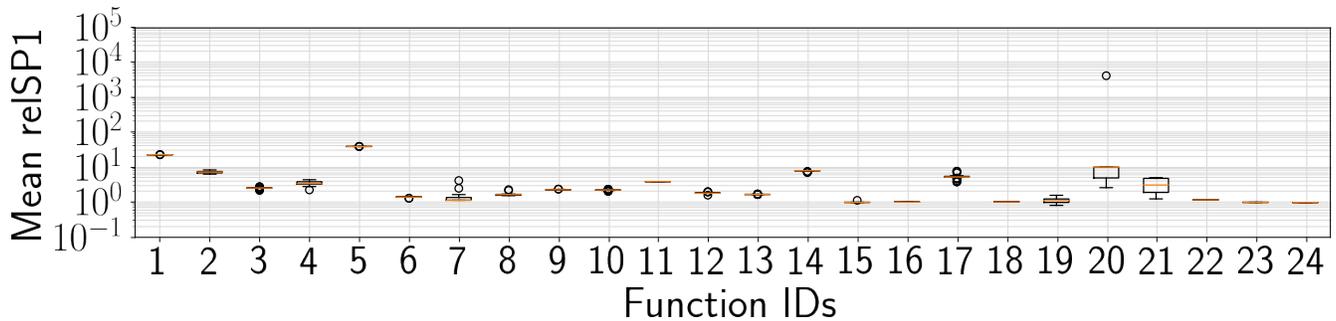
(a)  $n = 2$ (b)  $n = 3$ (c)  $n = 5$ (d)  $n = 10$ 

Fig. S.14: Distribution of 31 mean relSP1 values of the pairwise classification-based algorithm selection system with  $\mathcal{A}_{kt}$  for the LOIO-CV.

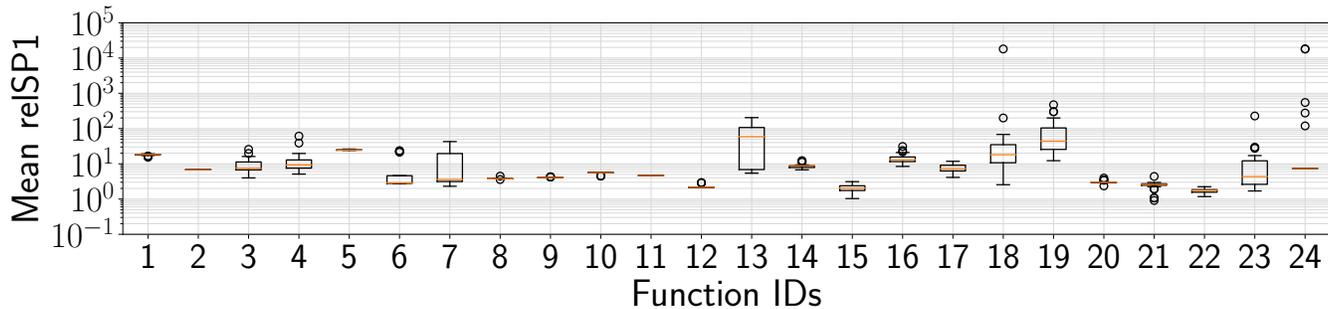
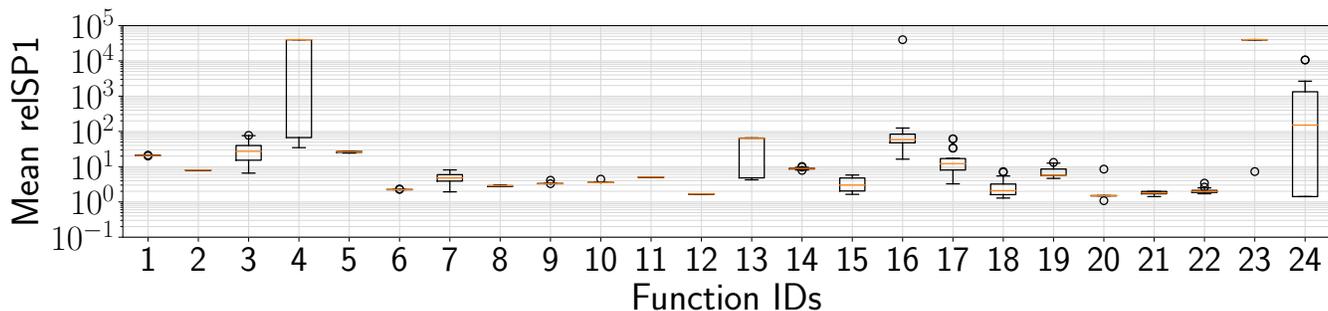
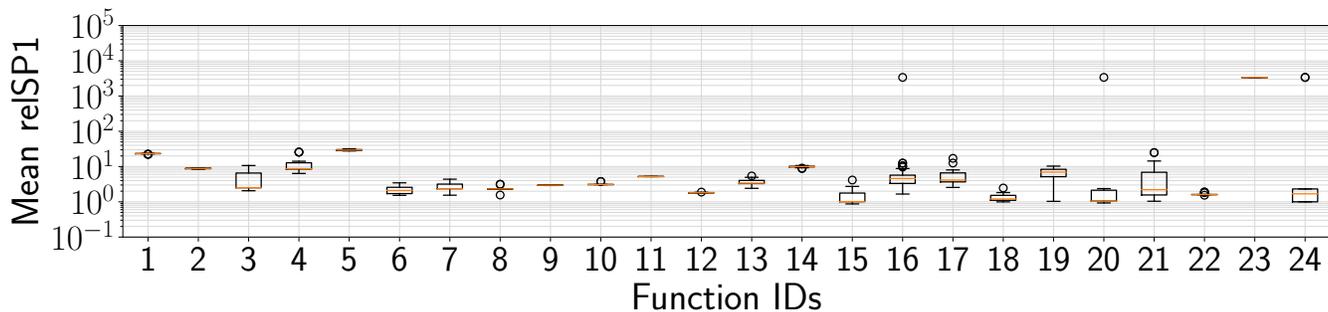
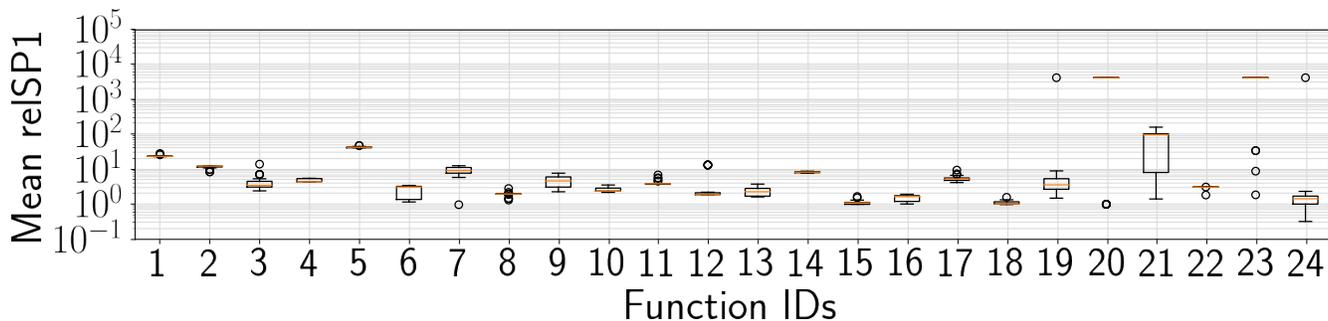
(a)  $n = 2$ (b)  $n = 3$ (c)  $n = 5$ (d)  $n = 10$ 

Fig. S.15: Distribution of 31 mean relSP1 values of the pairwise classification-based algorithm selection system with  $\mathcal{A}_{kt}$  for the LOPO-CV.

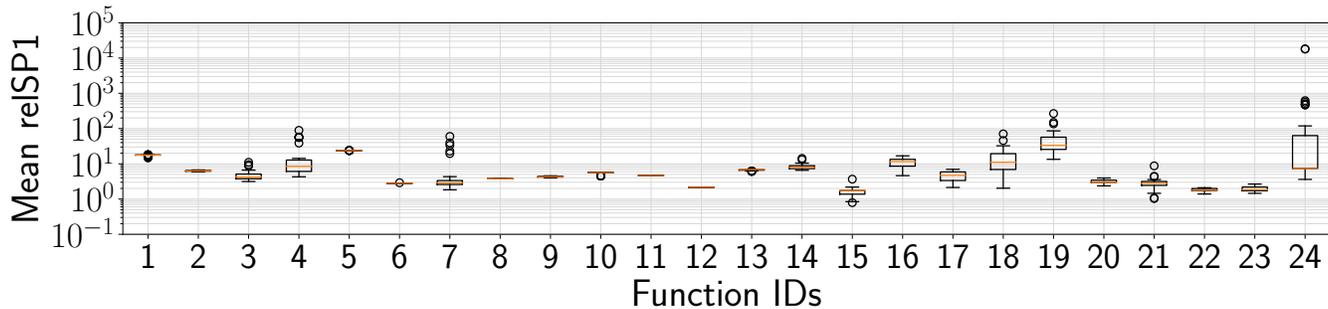
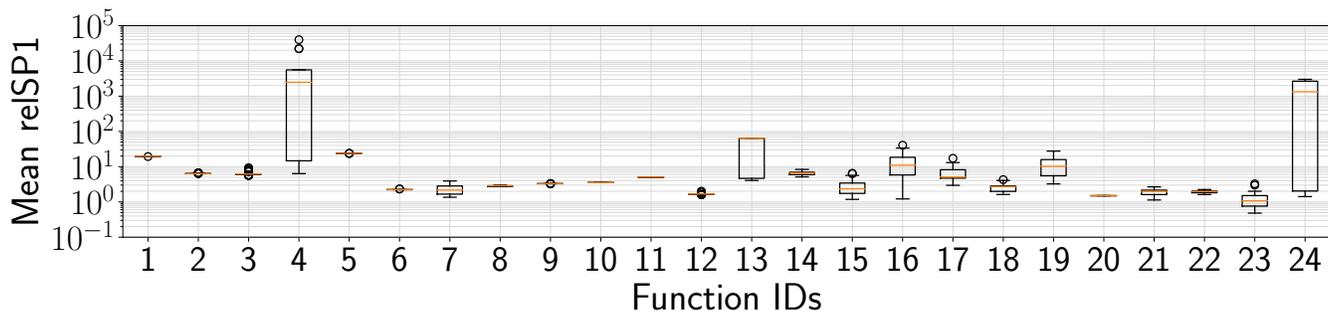
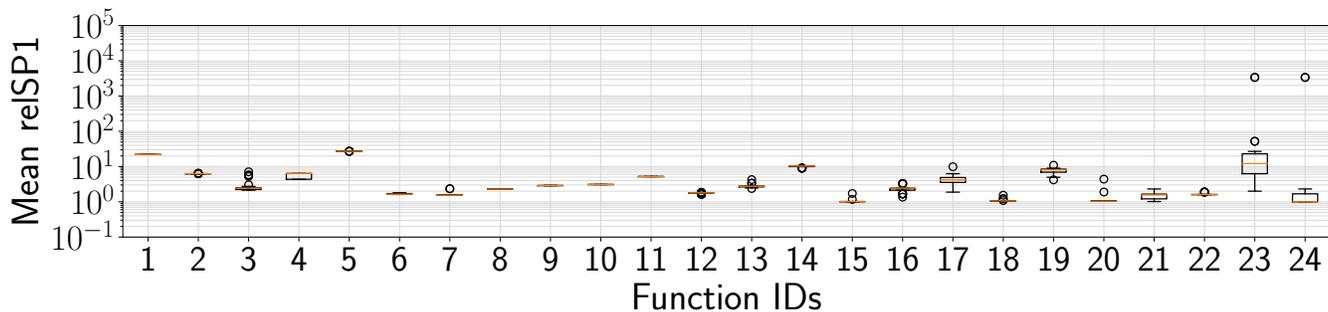
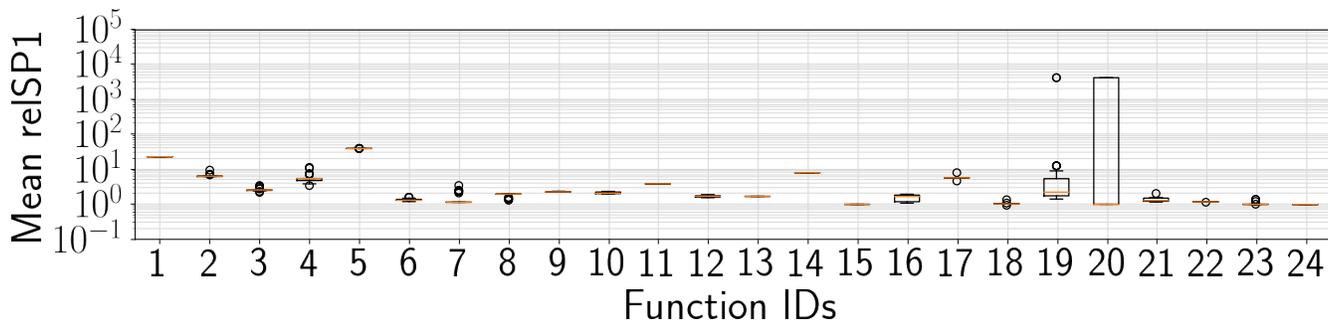
(a)  $n = 2$ (b)  $n = 3$ (c)  $n = 5$ (d)  $n = 10$ 

Fig. S.16: Distribution of 31 mean relSP1 values of the pairwise classification-based algorithm selection system with  $\mathcal{A}_{kt}$  for the LOPOAD-CV.

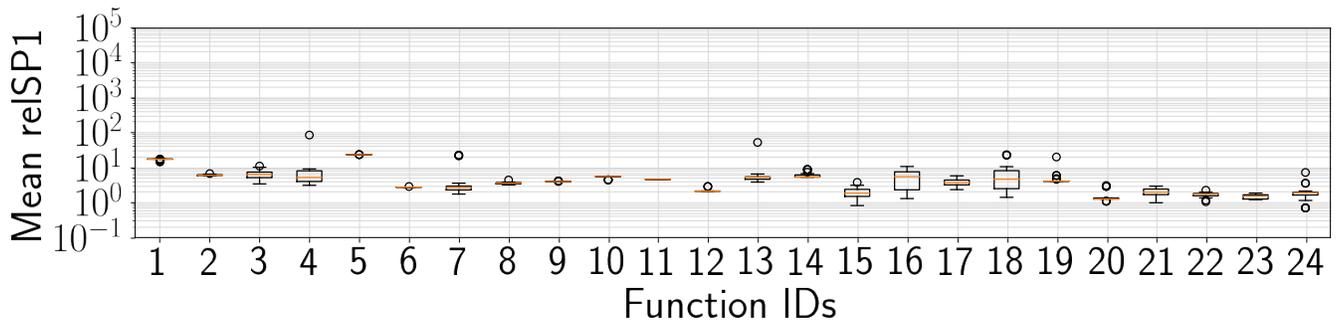
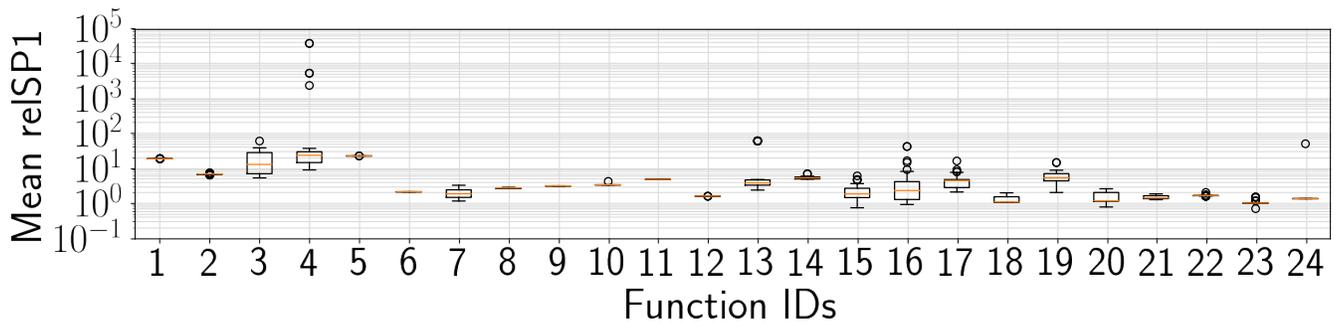
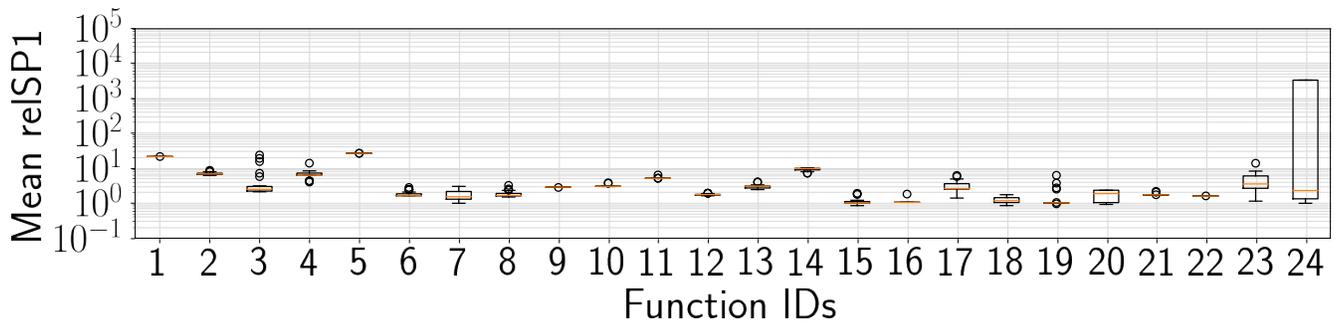
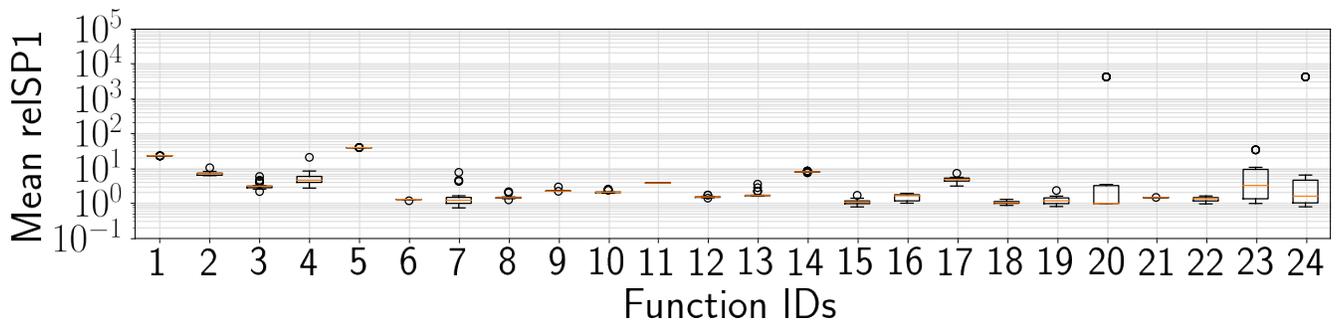
(a)  $n = 2$ (b)  $n = 3$ (c)  $n = 5$ (d)  $n = 10$ 

Fig. S.17: Distribution of 31 mean relSP1 values of the pairwise classification-based algorithm selection system with  $\mathcal{A}_{kt}$  for the RI-CV.

TABLE S.10: Median of 31 “mean relSP1” values of the classification-based algorithm selection system with and without the pre-solvers. “AS-50n” is the algorithm selection system with the sample size  $50 \times n$ . “AS-100n” is the algorithm selection system with the sample size  $100 \times n$ . “SLSQP-AS” is the algorithm selection system using SLSQP as a pre-solver. “SMAC-AS” is the algorithm selection system using SMAC as a pre-solver.

(a) LOIO-CV ( $\mathcal{A}_{kt}$ )					(b) LOPO-CV ( $\mathcal{A}_{kt}$ )					(c) LOPOAD-CV ( $\mathcal{A}_{kt}$ )					(d) RI-CV ( $\mathcal{A}_{kt}$ )				
System	$n = 2$	$n = 3$	$n = 5$	$n = 10$	System	$n = 2$	$n = 3$	$n = 5$	$n = 10$	System	$n = 2$	$n = 3$	$n = 5$	$n = 10$	System	$n = 2$	$n = 3$	$n = 5$	$n = 10$
AS-50n	7.54	10.10	4.75	4.99	AS-50n	33.77	6686.64	429.01	367.30	AS-50n	12.07	1695.37	288.15	361.93	AS-50n	7.63	6.57	4.38	4.73
AS-100n	10.84	12.86	7.79	8.39	AS-100n	33.53	6683.59	572.14	370.59	AS-100n	13.92	3351.36	430.87	365.25	AS-100n	10.15	10.11	7.34	8.04
SLSQP-AS	5.15+	7.80+	2.32+	2.27+	SLSQP-AS	27.21	6682.92	426.11+	363.88	SLSQP-AS	9.08+	1689.36	285.64	359.00	SLSQP-AS	5.30+	3.41+	1.95+	1.98+
SMAC-AS	7.86	10.97	6.84	8.32	SMAC-AS	33.24	5024.38+	430.95	370.63	SMAC-AS	11.82	129.39+	290.42	365.25	SMAC-AS	7.82	7.69	6.42	8.05
SBS	5.81	9.82	4.49	6.76	SBS	5.81	9.82	4.49	6.76	SBS	5.81	9.82	4.49	6.76	SBS	5.81	9.82	4.49	6.76

(e) LOIO-CV ( $\mathcal{A}_{dlvat}$ )					(f) LOPO-CV ( $\mathcal{A}_{dlvat}$ )					(g) LOPOAD-CV ( $\mathcal{A}_{dlvat}$ )					(h) RI-CV ( $\mathcal{A}_{dlvat}$ )				
System	$n = 2$	$n = 3$	$n = 5$	$n = 10$	System	$n = 2$	$n = 3$	$n = 5$	$n = 10$	System	$n = 2$	$n = 3$	$n = 5$	$n = 10$	System	$n = 2$	$n = 3$	$n = 5$	$n = 10$
AS-50n	7.56	8.44	4.70	4.31	AS-50n	15.20	5018.89	183.43	366.29	AS-50n	9.67	80.59	123.04	183.63	AS-50n	7.25	6.76	4.25	4.15
AS-100n	10.29	10.78	7.60	6.99	AS-100n	15.19	5017.77	186.99	368.81	AS-100n	12.68	1687.85	126.19	185.97	AS-100n	9.98	10.17	7.12	6.83
SLSQP-AS	5.17+	5.93+	2.34+	2.13+	SLSQP-AS	11.17+	5015.52	180.67+	363.78+	SLSQP-AS	7.18+	74.54	120.59	181.31+	SLSQP-AS	4.92+	3.52+	1.88+	1.92+
SMAC-AS	7.74	9.56	6.57	6.98	SMAC-AS	14.47	3355.69+	184.96	368.95	SMAC-AS	9.95	30.18	124.88	186.29	SMAC-AS	7.39	7.75	6.15	6.80
SBS	5.58	9.79	4.44	6.67	SBS	5.58	9.79	4.44	6.67	SBS	5.58	9.79	4.44	6.67	SBS	5.58	9.79	4.44	6.67

(i) LOIO-CV ( $\mathcal{A}_{jped}$ )					(j) LOPO-CV ( $\mathcal{A}_{jped}$ )					(k) LOPOAD-CV ( $\mathcal{A}_{jped}$ )					(l) RI-CV ( $\mathcal{A}_{jped}$ )				
System	$n = 2$	$n = 3$	$n = 5$	$n = 10$	System	$n = 2$	$n = 3$	$n = 5$	$n = 10$	System	$n = 2$	$n = 3$	$n = 5$	$n = 10$	System	$n = 2$	$n = 3$	$n = 5$	$n = 10$
AS-50n	9.06	10.12	4.92	4.40	AS-50n	768.17	5020.02	568.50	366.58	AS-50n	14.34	1688.45	568.31	365.54	AS-50n	9.46	6.50	4.32	4.12
AS-100n	11.31	12.69	7.86	7.18	AS-100n	21.43+	5019.51	572.53	369.31	AS-100n	16.41	1764.36	571.96	542.52	AS-100n	10.93	10.13	7.29	6.85
SLSQP-AS	6.59+	7.99+	2.47+	2.22+	SLSQP-AS	763.43+	5015.93	565.78+	363.81	SLSQP-AS	11.38+	1683.78	565.89	363.18	SLSQP-AS	7.14+	3.37+	1.93+	1.91+
SMAC-AS	9.20	11.15	6.87	7.14	SMAC-AS	767.13	3356.91+	570.32	369.31	SMAC-AS	14.22	50.73+	430.71	368.28	SMAC-AS	9.57	7.88	6.24	6.86
SBS	6.05	9.81	4.50	6.72	SBS	6.05	9.81	4.50	6.72	SBS	6.05	9.81	4.50	6.72	SBS	6.05	9.81	4.50	6.72

(m) LOIO-CV ( $\mathcal{A}_{bmtip}$ )					(n) LOPO-CV ( $\mathcal{A}_{bmtip}$ )					(o) LOPOAD-CV ( $\mathcal{A}_{bmtip}$ )					(p) RI-CV ( $\mathcal{A}_{bmtip}$ )				
System	$n = 2$	$n = 3$	$n = 5$	$n = 10$	System	$n = 2$	$n = 3$	$n = 5$	$n = 10$	System	$n = 2$	$n = 3$	$n = 5$	$n = 10$	System	$n = 2$	$n = 3$	$n = 5$	$n = 10$
AS-50n	5.56	240.63	3.85	158.07	AS-50n	12.35	715.12	288.17	625.37	AS-50n	19.31	10.58	2.93	158.26	AS-50n	7.67	451.54	3.47	3.23
AS-100n	6.49	243.44	4.95	159.63	AS-100n	14.84	714.82	291.53	626.50	AS-100n	20.78	10.86	4.46	159.82	AS-100n	7.81	246.19	4.92	4.85
SLSQP-AS	4.03+	239.16	2.42+	156.51+	SLSQP-AS	9.92+	710.52	286.30	622.69+	SLSQP-AS	17.81+	9.22	1.40+	156.55+	SLSQP-AS	6.20+	450.16	1.96+	1.66+
SMAC-AS	5.28	241.27	4.99	159.67	SMAC-AS	11.26+	715.54	289.43	626.95	SMAC-AS	18.72	11.33	4.08	159.84	SMAC-AS	7.54	358.71	4.53	4.83
SBS	11.53	35.92	102.10	315.32	SBS	11.53	35.92	102.10	315.32	SBS	11.53	35.92	102.10	315.32	SBS	11.53	35.92	102.10	315.32

(q) LOIO-CV ( $\mathcal{A}_{mk}$ )					(r) LOPO-CV ( $\mathcal{A}_{mk}$ )					(s) LOPOAD-CV ( $\mathcal{A}_{mk}$ )					(t) RI-CV ( $\mathcal{A}_{mk}$ )				
System	$n = 2$	$n = 3$	$n = 5$	$n = 10$	System	$n = 2$	$n = 3$	$n = 5$	$n = 10$	System	$n = 2$	$n = 3$	$n = 5$	$n = 10$	System	$n = 2$	$n = 3$	$n = 5$	$n = 10$
AS-50n	3.14	6.56	2.60	34.10	AS-50n	4.18	479.53	130.76	67.63	AS-50n	3.90	5.97	66.30	2.60	AS-50n	2.80	6.30	2.64	1.82
AS-100n	4.52	8.57	3.27	34.62	AS-100n	4.93	479.90	131.60	68.44	AS-100n	5.51	6.51	67.32	3.08	AS-100n	4.11	7.14	3.28	2.33
SLSQP-AS	2.00+	5.75+	1.69+	33.43+	SLSQP-AS	2.16+	469.43+	129.52+	66.69+	SLSQP-AS	2.30+	4.57+	65.28+	1.82+	SLSQP-AS	1.67+	5.17+	1.71+	1.17+
SMAC-AS	3.21	6.84	3.04	34.63	SMAC-AS	3.64	474.62	131.16	68.18	SMAC-AS	3.86	5.64	66.67	3.16	SMAC-AS	2.89	6.27	3.04	2.36
SBS	1.64	31.99	68.74	67.44	SBS	1.64	31.99	68.74	67.44	SBS	1.64	31.99	68.74	67.44	SBS	1.64	31.99	68.74	67.44

TABLE S.11: Median of 31 “mean relSP1” values of the regression-based algorithm selection system with and without the pre-solvers. “AS-50n” is the algorithm selection system with the sample size  $50 \times n$ . “AS-100n” is the algorithm selection system with the sample size  $100 \times n$ . “SLSQP-AS” is the algorithm selection system using SLSQP as a pre-solver. “SMAC-AS” is the algorithm selection system using SMAC as a pre-solver.

(a) LOIO-CV ( $\mathcal{A}_{kt}$ )					(b) LOPO-CV ( $\mathcal{A}_{kt}$ )					(c) LOPOAD-CV ( $\mathcal{A}_{kt}$ )					(d) RI-CV ( $\mathcal{A}_{kt}$ )				
System	$n = 2$	$n = 3$	$n = 5$	$n = 10$	System	$n = 2$	$n = 3$	$n = 5$	$n = 10$	System	$n = 2$	$n = 3$	$n = 5$	$n = 10$	System	$n = 2$	$n = 3$	$n = 5$	$n = 10$
AS-50n	9.48	67.95	7.18	9.98	AS-50n	16.73	134.31	7.74	10.08	AS-50n	10.11	191.74	8.43	10.20	AS-50n	6.31	11.96	6.82	9.98
AS-100n	10.70	70.79	10.23	13.21	AS-100n	24.55	136.81	10.81	13.44	AS-100n	15.66	169.12	10.17	13.52	AS-100n	8.99	13.53	9.78	13.10
SLSQP-AS	7.02+	65.50	4.23+	6.88+	SLSQP-AS	10.83+	131.85	4.58+	6.90+	SLSQP-AS	7.08+	189.23	5.44+	6.99+	SLSQP-AS	3.54+	9.58+	3.97+	6.86+
SMAC-AS	9.40	68.07	9.08	13.30	SMAC-AS	14.49	119.50	9.75	13.40	SMAC-AS	9.85	123.26	10.33	13.52	SMAC-AS	6.40	11.43	8.75	13.30
SBS	5.81	9.82	4.49	6.76	SBS	5.81	9.82	4.49	6.76	SBS	5.81	9.82	4.49	6.76	SBS	5.81	9.82	4.49	6.76
(e) LOIO-CV ( $\mathcal{A}_{dlvat}$ )					(f) LOPO-CV ( $\mathcal{A}_{dlvat}$ )					(g) LOPOAD-CV ( $\mathcal{A}_{dlvat}$ )					(h) RI-CV ( $\mathcal{A}_{dlvat}$ )				
System	$n = 2$	$n = 3$	$n = 5$	$n = 10$	System	$n = 2$	$n = 3$	$n = 5$	$n = 10$	System	$n = 2$	$n = 3$	$n = 5$	$n = 10$	System	$n = 2$	$n = 3$	$n = 5$	$n = 10$
AS-50n	8.47	68.43	7.08	9.28	AS-50n	15.06	129.05	7.56	9.34	AS-50n	10.14	191.80	8.48	9.39	AS-50n	6.33	11.76	6.55	9.24
AS-100n	10.61	70.79	9.81	11.86	AS-100n	19.94	135.92	10.63	12.04	AS-100n	15.00	144.31	9.83	12.09	AS-100n	8.64	13.41	9.44	11.76
SLSQP-AS	6.02+	65.97	4.24+	6.70+	SLSQP-AS	11.02+	126.60	4.52+	6.75+	SLSQP-AS	7.45+	189.32	5.63+	6.78+	SLSQP-AS	3.78+	9.34+	3.90+	6.68+
SMAC-AS	8.58	68.32	8.84	11.93	SMAC-AS	14.26	120.71	9.64	12.00	SMAC-AS	10.05	123.35	10.76	12.04	SMAC-AS	6.18	11.09	8.45	11.89
SBS	5.58	9.79	4.44	6.67	SBS	5.58	9.79	4.44	6.67	SBS	5.58	9.79	4.44	6.67	SBS	5.58	9.79	4.44	6.67
(i) LOIO-CV ( $\mathcal{A}_{jped}$ )					(j) LOPO-CV ( $\mathcal{A}_{jped}$ )					(k) LOPOAD-CV ( $\mathcal{A}_{jped}$ )					(l) RI-CV ( $\mathcal{A}_{jped}$ )				
System	$n = 2$	$n = 3$	$n = 5$	$n = 10$	System	$n = 2$	$n = 3$	$n = 5$	$n = 10$	System	$n = 2$	$n = 3$	$n = 5$	$n = 10$	System	$n = 2$	$n = 3$	$n = 5$	$n = 10$
AS-50n	8.93	67.72	7.20	9.70	AS-50n	17.85	134.54	9.23	9.50	AS-50n	13.77	148.09	11.87	10.17	AS-50n	6.89	11.85	6.57	9.27
AS-100n	10.75	70.79	10.19	12.21	AS-100n	31.97	134.18	13.29	12.25	AS-100n	18.46	124.02	12.39	12.44	AS-100n	9.04	13.49	9.61	9.60
SLSQP-AS	6.49+	65.28	4.37+	7.12+	SLSQP-AS	15.20	132.10	6.16+	6.88+	SLSQP-AS	10.64	145.55	8.79+	7.24+	SLSQP-AS	3.85+	9.46+	3.86+	6.68+
SMAC-AS	9.14	67.78	9.09	12.43	SMAC-AS	17.70	115.25	11.16	12.22	SMAC-AS	14.09	103.12	14.04	12.90	SMAC-AS	6.84	11.56	8.62	11.99
SBS	6.05	9.81	4.50	6.72	SBS	6.05	9.81	4.50	6.72	SBS	6.05	9.81	4.50	6.72	SBS	6.05	9.81	4.50	6.72
(m) LOIO-CV ( $\mathcal{A}_{bmtip}$ )					(n) LOPO-CV ( $\mathcal{A}_{bmtip}$ )					(o) LOPOAD-CV ( $\mathcal{A}_{bmtip}$ )					(p) RI-CV ( $\mathcal{A}_{bmtip}$ )				
System	$n = 2$	$n = 3$	$n = 5$	$n = 10$	System	$n = 2$	$n = 3$	$n = 5$	$n = 10$	System	$n = 2$	$n = 3$	$n = 5$	$n = 10$	System	$n = 2$	$n = 3$	$n = 5$	$n = 10$
AS-50n	11.94	37.47	98.73	8.34	AS-50n	13.07	37.64	103.38	163.95	AS-50n	13.53	36.45	101.67	316.23	AS-50n	9.98	37.59	99.02	7.96
AS-100n	12.67	39.25	16.04	10.26	AS-100n	15.01	39.36	104.37	316.58	AS-100n	14.49	38.20	101.26	317.89	AS-100n	12.63	39.33	12.35	9.16
SLSQP-AS	9.23+	34.58+	95.48+	3.26+	SLSQP-AS	10.25+	34.72+	100.02+	159.64+	SLSQP-AS	11.42+	34.04+	99.07+	312.87+	SLSQP-AS	7.33+	34.71+	96.28+	3.21+
SMAC-AS	10.84	37.79	99.81	9.93	SMAC-AS	11.99+	37.96	104.34	165.54	SMAC-AS	11.86+	36.74	102.76	317.81	SMAC-AS	9.01	37.84	100.11	9.57
SBS	11.53	35.92	102.10	315.32	SBS	11.53	35.92	102.10	315.32	SBS	11.53	35.92	102.10	315.32	SBS	11.53	35.92	102.10	315.32
(q) LOIO-CV ( $\mathcal{A}_{mk}$ )					(r) LOPO-CV ( $\mathcal{A}_{mk}$ )					(s) LOPOAD-CV ( $\mathcal{A}_{mk}$ )					(t) RI-CV ( $\mathcal{A}_{mk}$ )				
System	$n = 2$	$n = 3$	$n = 5$	$n = 10$	System	$n = 2$	$n = 3$	$n = 5$	$n = 10$	System	$n = 2$	$n = 3$	$n = 5$	$n = 10$	System	$n = 2$	$n = 3$	$n = 5$	$n = 10$
AS-50n	3.13	12.46	68.31	3.29	AS-50n	3.05	239.51	70.37	67.65	AS-50n	9.02	14.18	70.64	35.68	AS-50n	3.02	9.46	69.09	3.10
AS-100n	4.57	11.31	68.69	3.68	AS-100n	4.42	241.55	70.76	37.42	AS-100n	7.20	10.52	72.36	36.10	AS-100n	4.37	9.71	68.87	3.52
SLSQP-AS	2.01+	11.30	67.15+	2.18+	SLSQP-AS	1.93+	238.45	69.07+	66.83+	SLSQP-AS	7.47+	13.08	69.47+	34.77+	SLSQP-AS	1.88+	8.38+	67.62+	2.17+
SMAC-AS	3.23	12.65	68.47	3.85	SMAC-AS	3.12	10.00+	70.74	68.17	SMAC-AS	8.76	8.99+	70.89	36.25	SMAC-AS	3.11	9.54	69.16	3.65
SBS	1.64	31.99	68.74	67.44	SBS	1.64	31.99	68.74	67.44	SBS	1.64	31.99	68.74	67.44	SBS	1.64	31.99	68.74	67.44

TABLE S.12: Median of 31 “mean relSP1” values of the pairwise classification-based algorithm selection system with and without the pre-solvers. “AS-50n” is the algorithm selection system with the sample size  $50 \times n$ . “AS-100n” is the algorithm selection system with the sample size  $100 \times n$ . “SLSQP-AS” is the algorithm selection system using SLSQP as a pre-solver. “SMAC-AS” is the algorithm selection system using SMAC as a pre-solver.

(a) LOIO-CV ( $\mathcal{A}_{kt}$ )					(b) LOPO-CV ( $\mathcal{A}_{kt}$ )					(c) LOPOAD-CV ( $\mathcal{A}_{kt}$ )					(d) RI-CV ( $\mathcal{A}_{kt}$ )				
System	$n = 2$	$n = 3$	$n = 5$	$n = 10$	System	$n = 2$	$n = 3$	$n = 5$	$n = 10$	System	$n = 2$	$n = 3$	$n = 5$	$n = 10$	System	$n = 2$	$n = 3$	$n = 5$	$n = 10$
AS-50n	5.94	8.13	5.01	5.28	AS-50n	15.75	3337.17	146.83	363.62	AS-50n	9.76	131.42	5.42	5.29	AS-50n	5.43	6.86	5.10	5.44
AS-100n	8.75-	11.33-	7.88-	8.59-	AS-100n	16.66	3340.90	149.85-	366.78	AS-100n	11.22-	124.18	8.01-	8.60-	AS-100n	8.37-	10.98-	7.85-	8.64-
SLSQP-AS	3.53+	4.85+	2.54+	2.49+	SLSQP-AS	11.58+	3334.35	144.06+	360.35	SLSQP-AS	7.36+	128.10	2.92+	2.50+	SLSQP-AS	2.98+	4.25+	2.65+	2.66+
SMAC-AS	6.18	7.82	7.07-	8.60-	SMAC-AS	15.59	3338.37	148.80-	366.99-	SMAC-AS	9.68	115.13	7.54-	8.64-	SMAC-AS	5.62-	7.89-	7.29-	8.79-
SBS	5.81	9.82	4.49	6.76	SBS	5.81	9.82	4.49	6.76	SBS	5.81	9.82	4.49	6.76	SBS	5.81	9.82	4.49	6.76
(e) LOIO-CV ( $\mathcal{A}_{dlvat}$ )					(f) LOPO-CV ( $\mathcal{A}_{dlvat}$ )					(g) LOPOAD-CV ( $\mathcal{A}_{dlvat}$ )					(h) RI-CV ( $\mathcal{A}_{dlvat}$ )				
System	$n = 2$	$n = 3$	$n = 5$	$n = 10$	System	$n = 2$	$n = 3$	$n = 5$	$n = 10$	System	$n = 2$	$n = 3$	$n = 5$	$n = 10$	System	$n = 2$	$n = 3$	$n = 5$	$n = 10$
AS-50n	5.78	6.63	4.88	4.60	AS-50n	12.81	61.92	124.22	366.02	AS-50n	9.24	132.13	63.63	182.67	AS-50n	5.66	6.75	5.70	4.71
AS-100n	8.52-	9.44-	7.66-	7.24-	AS-100n	15.46-	1676.44	127.15	542.89-	AS-100n	11.06-	127.44	7.94	185.24-	AS-100n	8.75-	9.07-	8.30-	7.30-
SLSQP-AS	3.42+	4.48+	2.46+	2.33+	SLSQP-AS	10.04+	59.50	121.63+	363.43+	SLSQP-AS	6.74+	129.94	61.16+	180.40+	SLSQP-AS	3.37+	4.57+	3.37+	2.45+
SMAC-AS	5.94	7.70	6.81-	7.25-	SMAC-AS	12.75	61.81	126.11	368.73-	SMAC-AS	8.98	115.02	65.33-	185.32-	SMAC-AS	5.87	7.68-	7.79-	7.38-
SBS	5.58	9.79	4.44	6.67	SBS	5.58	9.79	4.44	6.67	SBS	5.58	9.79	4.44	6.67	SBS	5.58	9.79	4.44	6.67
(i) LOIO-CV ( $\mathcal{A}_{jped}$ )					(j) LOPO-CV ( $\mathcal{A}_{jped}$ )					(k) LOPOAD-CV ( $\mathcal{A}_{jped}$ )					(l) RI-CV ( $\mathcal{A}_{jped}$ )				
System	$n = 2$	$n = 3$	$n = 5$	$n = 10$	System	$n = 2$	$n = 3$	$n = 5$	$n = 10$	System	$n = 2$	$n = 3$	$n = 5$	$n = 10$	System	$n = 2$	$n = 3$	$n = 5$	$n = 10$
AS-50n	5.99	7.60	5.03	4.69	AS-50n	16.22	1674.33	146.95	362.62	AS-50n	8.19	111.15	5.34	4.59	AS-50n	5.60	6.48	4.68	4.40
AS-100n	8.92-	9.72-	7.84-	7.35-	AS-100n	16.88	1679.97-	149.85-	192.24	AS-100n	11.36-	111.98	8.00-	7.44-	AS-100n	8.50-	10.17-	7.66-	7.10-
SLSQP-AS	3.57+	5.07+	2.59+	2.43+	SLSQP-AS	11.85+	1671.79	144.19+	359.85+	SLSQP-AS	5.82+	107.28	2.85+	2.28+	SLSQP-AS	3.11+	4.17+	2.22+	2.13+
SMAC-AS	6.24	7.71	7.03-	7.42-	SMAC-AS	16.23	1675.28	148.80-	365.36-	SMAC-AS	8.45	109.74	7.46-	7.32-	SMAC-AS	5.71	7.60-	6.71-	7.15-
SBS	6.05	9.81	4.50	6.72	SBS	6.05	9.81	4.50	6.72	SBS	6.05	9.81	4.50	6.72	SBS	6.05	9.81	4.50	6.72
(m) LOIO-CV ( $\mathcal{A}_{bmtpt}$ )					(n) LOPO-CV ( $\mathcal{A}_{bmtpt}$ )					(o) LOPOAD-CV ( $\mathcal{A}_{bmtpt}$ )					(p) RI-CV ( $\mathcal{A}_{bmtpt}$ )				
System	$n = 2$	$n = 3$	$n = 5$	$n = 10$	System	$n = 2$	$n = 3$	$n = 5$	$n = 10$	System	$n = 2$	$n = 3$	$n = 5$	$n = 10$	System	$n = 2$	$n = 3$	$n = 5$	$n = 10$
AS-50n	7.41	36.15	96.86	158.07	AS-50n	12.30	947.32	109.69	323.98	AS-50n	12.93	13.28	7.19	158.48	AS-50n	8.73	95.37	95.92	160.57
AS-100n	7.82	18.40+	7.07	160.06-	AS-100n	14.23-	944.64	106.29	324.18	AS-100n	14.82-	12.54	7.41	160.06-	AS-100n	8.44	239.83	97.68	160.78
SLSQP-AS	5.55	34.76	95.28+	156.52+	SLSQP-AS	10.20+	743.19	103.99	314.43+	SLSQP-AS	11.41+	11.73	5.29+	156.84+	SLSQP-AS	7.19+	93.95	93.76+	157.20+
SMAC-AS	6.77	36.83	97.93-	159.64-	SMAC-AS	11.51	716.73	111.16	325.62	SMAC-AS	10.76+	13.79	8.27	160.10-	SMAC-AS	8.65	45.85	97.18	162.17
SBS	11.53	35.92	102.10	315.32	SBS	11.53	35.92	102.10	315.32	SBS	11.53	35.92	102.10	315.32	SBS	11.53	35.92	102.10	315.32
(q) LOIO-CV ( $\mathcal{A}_{mk}$ )					(r) LOPO-CV ( $\mathcal{A}_{mk}$ )					(s) LOPOAD-CV ( $\mathcal{A}_{mk}$ )					(t) RI-CV ( $\mathcal{A}_{mk}$ )				
System	$n = 2$	$n = 3$	$n = 5$	$n = 10$	System	$n = 2$	$n = 3$	$n = 5$	$n = 10$	System	$n = 2$	$n = 3$	$n = 5$	$n = 10$	System	$n = 2$	$n = 3$	$n = 5$	$n = 10$
AS-50n	3.20	8.46	66.46	34.00	AS-50n	3.11	240.73	131.75	164.48	AS-50n	3.36	6.52	65.99	35.64	AS-50n	3.01	5.72	68.25	66.84
AS-100n	4.59-	9.44-	66.15	34.76-	AS-100n	4.46-	241.04	132.36	165.14	AS-100n	4.76-	7.17-	66.86	67.56-	AS-100n	4.33-	7.22-	68.42	67.77
SLSQP-AS	2.05+	7.63+	65.53	33.35+	SLSQP-AS	1.92+	239.85	130.87	163.76	SLSQP-AS	2.21+	5.68+	65.09	34.86+	SLSQP-AS	1.87+	4.89+	67.40	66.20
SMAC-AS	3.29	8.70	67.00	34.52-	SMAC-AS	3.15	8.53+	132.13	165.03	SMAC-AS	3.42	6.11	66.38	36.20-	SMAC-AS	3.07-	5.59	68.48	67.38
SBS	1.64	31.99	68.74	67.44	SBS	1.64	31.99	68.74	67.44	SBS	1.64	31.99	68.74	67.44	SBS	1.64	31.99	68.74	67.44

TABLE S.13: Median of 31 “mean relSP1” values of the pairwise regression-based algorithm selection system with and without the pre-solvers. “AS-50n” is the algorithm selection system with the sample size  $50 \times n$ . “AS-100n” is the algorithm selection system with the sample size  $100 \times n$ . “SLSQP-AS” is the algorithm selection system using SLSQP as a pre-solver. “SMAC-AS” is the algorithm selection system using SMAC as a pre-solver.

(a) LOIO-CV ( $\mathcal{A}_{kt}$ )					(b) LOPO-CV ( $\mathcal{A}_{kt}$ )					(c) LOPOAD-CV ( $\mathcal{A}_{kt}$ )					(d) RI-CV ( $\mathcal{A}_{kt}$ )				
System	$n = 2$	$n = 3$	$n = 5$	$n = 10$	System	$n = 2$	$n = 3$	$n = 5$	$n = 10$	System	$n = 2$	$n = 3$	$n = 5$	$n = 10$	System	$n = 2$	$n = 3$	$n = 5$	$n = 10$
AS-50n	18.65	72.91	9.48	11.10	AS-50n	42.90	246.10	150.56	188.74	AS-50n	19.30	144.11	8.74	11.35	AS-50n	16.14	73.11	8.98	11.20
AS-100n	25.46	77.57	12.19-	14.60-	AS-100n	44.35	1685.43-	293.02-	192.32-	AS-100n	19.12	148.93	11.56-	15.32-	AS-100n	26.69-	71.44	11.99-	14.65-
SLSQP-AS	14.53+	69.91	5.45+	7.35+	SLSQP-AS	33.00	242.60	146.98+	185.27+	SLSQP-AS	14.96+	141.01	5.07+	7.68+	SLSQP-AS	11.25+	68.85	5.18+	7.38+
SMAC-AS	18.57	71.47	11.39-	14.42-	SMAC-AS	42.60	30.71+	152.14-	192.05-	SMAC-AS	19.33	75.74+	10.42-	14.69-	SMAC-AS	16.21	71.91	10.88-	14.52-
SBS	5.81	9.82	4.49	6.76	SBS	5.81	9.82	4.49	6.76	SBS	5.81	9.82	4.49	6.76	SBS	5.81	9.82	4.49	6.76

(e) LOIO-CV ( $\mathcal{A}_{dlvat}$ )					(f) LOPO-CV ( $\mathcal{A}_{dlvat}$ )					(g) LOPOAD-CV ( $\mathcal{A}_{dlvat}$ )					(h) RI-CV ( $\mathcal{A}_{dlvat}$ )				
System	$n = 2$	$n = 3$	$n = 5$	$n = 10$	System	$n = 2$	$n = 3$	$n = 5$	$n = 10$	System	$n = 2$	$n = 3$	$n = 5$	$n = 10$	System	$n = 2$	$n = 3$	$n = 5$	$n = 10$
AS-50n	22.17	19.97	7.86	9.96	AS-50n	25.53	30.26	8.35	188.05	AS-50n	15.54	135.69	9.46	10.22	AS-50n	18.39	16.82	7.97	10.07
AS-100n	20.70	62.68-	10.33-	9.31	AS-100n	29.43	71.90	11.65-	190.84-	AS-100n	17.20	135.65	10.72-	13.83-	AS-100n	21.33-	19.44	10.47-	10.16-
SLSQP-AS	19.00	16.83	4.70+	6.98+	SLSQP-AS	20.79	27.63	5.20+	185.14+	SLSQP-AS	12.62+	132.79	6.07+	7.44+	SLSQP-AS	14.87+	14.06+	4.78+	6.94+
SMAC-AS	22.09	17.62	9.60-	12.61-	SMAC-AS	24.40	18.38+	10.02-	190.70-	SMAC-AS	14.97	68.77+	11.20-	12.89-	SMAC-AS	16.80	16.79	9.70-	12.73-
SBS	5.58	9.79	4.44	6.67	SBS	5.58	9.79	4.44	6.67	SBS	5.58	9.79	4.44	6.67	SBS	5.58	9.79	4.44	6.67

(i) LOIO-CV ( $\mathcal{A}_{jped}$ )					(j) LOPO-CV ( $\mathcal{A}_{jped}$ )					(k) LOPOAD-CV ( $\mathcal{A}_{jped}$ )					(l) RI-CV ( $\mathcal{A}_{jped}$ )				
System	$n = 2$	$n = 3$	$n = 5$	$n = 10$	System	$n = 2$	$n = 3$	$n = 5$	$n = 10$	System	$n = 2$	$n = 3$	$n = 5$	$n = 10$	System	$n = 2$	$n = 3$	$n = 5$	$n = 10$
AS-50n	17.46	31.77	11.31	10.49	AS-50n	24.20	60.95	150.88	12.89	AS-50n	19.57	66.80	10.31	10.31	AS-50n	17.42	21.05	11.76	8.75
AS-100n	23.17-	24.33	14.17-	12.19-	AS-100n	31.61	61.91	155.07-	23.22-	AS-100n	29.39-	66.42	13.84-	13.47-	AS-100n	21.52-	21.39	16.29-	11.49-
SLSQP-AS	13.27+	27.08	7.57+	7.12+	SLSQP-AS	18.58+	57.88	147.28+	9.60+	SLSQP-AS	15.29+	64.18	6.67+	6.76+	SLSQP-AS	13.84+	16.80	7.85+	5.43+
SMAC-AS	17.42	29.74	12.88-	13.22-	SMAC-AS	24.14	36.94+	153.09-	15.63-	SMAC-AS	19.05	34.54+	12.20-	13.08-	SMAC-AS	16.26	17.99	13.46-	11.48-
SBS	6.05	9.81	4.50	6.72	SBS	6.05	9.81	4.50	6.72	SBS	6.05	9.81	4.50	6.72	SBS	6.05	9.81	4.50	6.72

(m) LOIO-CV ( $\mathcal{A}_{bmtip}$ )					(n) LOPO-CV ( $\mathcal{A}_{bmtip}$ )					(o) LOPOAD-CV ( $\mathcal{A}_{bmtip}$ )					(p) RI-CV ( $\mathcal{A}_{bmtip}$ )				
System	$n = 2$	$n = 3$	$n = 5$	$n = 10$	System	$n = 2$	$n = 3$	$n = 5$	$n = 10$	System	$n = 2$	$n = 3$	$n = 5$	$n = 10$	System	$n = 2$	$n = 3$	$n = 5$	$n = 10$
AS-50n	7.18	36.21	10.02	7.20	AS-50n	18.49	268.47	102.23	164.71	AS-50n	12.38	32.27	5.10	163.98	AS-50n	7.20	23.79	8.03	7.49
AS-100n	8.26-	28.16	10.06	8.93-	AS-100n	21.16	252.01	104.30-	318.32-	AS-100n	12.65	33.35	6.47-	166.37-	AS-100n	8.52-	24.77	9.57	9.54-
SLSQP-AS	4.22+	33.98	6.58+	3.38+	SLSQP-AS	14.90	265.82	99.61+	158.71+	SLSQP-AS	9.45+	29.53	2.13+	158.51+	SLSQP-AS	4.17+	21.31	5.31+	3.46+
SMAC-AS	6.17+	36.47	11.16-	8.78-	SMAC-AS	15.70	266.91	103.33-	166.33	SMAC-AS	11.08	32.37	5.89-	165.56	SMAC-AS	6.02+	23.89	8.84	9.09-
SBS	11.53	35.92	102.10	315.32	SBS	11.53	35.92	102.10	315.32	SBS	11.53	35.92	102.10	315.32	SBS	11.53	35.92	102.10	315.32

(q) LOIO-CV ( $\mathcal{A}_{mk}$ )					(r) LOPO-CV ( $\mathcal{A}_{mk}$ )					(s) LOPOAD-CV ( $\mathcal{A}_{mk}$ )					(t) RI-CV ( $\mathcal{A}_{mk}$ )				
System	$n = 2$	$n = 3$	$n = 5$	$n = 10$	System	$n = 2$	$n = 3$	$n = 5$	$n = 10$	System	$n = 2$	$n = 3$	$n = 5$	$n = 10$	System	$n = 2$	$n = 3$	$n = 5$	$n = 10$
AS-50n	4.74	10.14	5.07	3.14	AS-50n	5.02	240.59	70.34	36.89	AS-50n	4.71	240.34	67.54	35.51	AS-50n	4.39	8.21	8.39	3.07
AS-100n	6.15-	11.26-	5.96	3.73-	AS-100n	6.18-	242.67-	69.91	37.24	AS-100n	5.94-	241.88-	68.58	4.55	AS-100n	6.21-	9.34-	6.20	4.05-
SLSQP-AS	3.29+	9.06+	4.08+	1.85+	SLSQP-AS	3.65+	239.51+	69.23+	35.58+	SLSQP-AS	3.19+	239.33+	66.45+	34.73+	SLSQP-AS	3.04+	7.16+	7.48+	1.81+
SMAC-AS	4.62	10.18	5.48	3.67-	SMAC-AS	4.91	8.31+	70.75	37.44	SMAC-AS	4.35	8.04+	67.87	36.04-	SMAC-AS	4.33	8.49	8.70	3.60-
SBS	1.64	31.99	68.74	67.44	SBS	1.64	31.99	68.74	67.44	SBS	1.64	31.99	68.74	67.44	SBS	1.64	31.99	68.74	67.44

TABLE S.14: Median of 31 “mean relSP1” values of the clustering-based algorithm selection system with and without the pre-solvers. “AS-50n” is the algorithm selection system with the sample size  $50 \times n$ . “AS-100n” is the algorithm selection system with the sample size  $100 \times n$ . “SLSQP-AS” is the algorithm selection system using SLSQP as a pre-solver. “SMAC-AS” is the algorithm selection system using SMAC as a pre-solver.

(a) LOIO-CV ( $\mathcal{A}_{kt}$ )					(b) LOPO-CV ( $\mathcal{A}_{kt}$ )					(c) LOPOAD-CV ( $\mathcal{A}_{kt}$ )					(d) RI-CV ( $\mathcal{A}_{kt}$ )				
System	$n = 2$	$n = 3$	$n = 5$	$n = 10$	System	$n = 2$	$n = 3$	$n = 5$	$n = 10$	System	$n = 2$	$n = 3$	$n = 5$	$n = 10$	System	$n = 2$	$n = 3$	$n = 5$	$n = 10$
AS-50n	761.68	15.39	289.31	373.30	AS-50n	11.57	1745.19	431.31	903.53	AS-50n	12.79	65.33	292.15	728.66	AS-50n	9.94	21.44	150.46	544.56
AS-100n	764.26–	22.39–	291.94	552.62	AS-100n	14.41–	3349.86–	433.75–	743.96	AS-100n	14.86–	522.32	296.93–	909.25–	AS-100n	12.52–	1673.66–	151.76–	374.35
SLSQP-AS	759.25+	12.25+	286.81	370.24	SLSQP-AS	8.65+	1741.12	428.40+	899.90	SLSQP-AS	9.42+	62.57	289.29+	725.37	SLSQP-AS	7.01+	16.93	147.84+	541.73
SMAC-AS	761.59	13.98	291.01	376.76	SMAC-AS	11.45	1686.97	432.90	906.86	SMAC-AS	12.32	20.33+	294.08	732.13	SMAC-AS	9.79	16.41	152.37	548.03
SBS	5.81	9.82	4.49	6.76	SBS	5.81	9.82	4.49	6.76	SBS	5.81	9.82	4.49	6.76	SBS	5.81	9.82	4.49	6.76

(e) LOIO-CV ( $\mathcal{A}_{dlvat}$ )					(f) LOPO-CV ( $\mathcal{A}_{dlvat}$ )					(g) LOPOAD-CV ( $\mathcal{A}_{dlvat}$ )					(h) RI-CV ( $\mathcal{A}_{dlvat}$ )				
System	$n = 2$	$n = 3$	$n = 5$	$n = 10$	System	$n = 2$	$n = 3$	$n = 5$	$n = 10$	System	$n = 2$	$n = 3$	$n = 5$	$n = 10$	System	$n = 2$	$n = 3$	$n = 5$	$n = 10$
AS-50n	761.19	14.40	125.47	372.41	AS-50n	11.11	1688.18	185.85	902.62	AS-50n	12.11	65.26	128.37	727.80	AS-50n	9.87	21.31	68.03	543.79
AS-100n	764.49–	21.74–	127.96	548.54	AS-100n	13.74–	3348.24–	188.02	742.80	AS-100n	14.34–	521.69	133.03–	907.68–	AS-100n	12.34–	1673.55–	69.72–	371.10
SLSQP-AS	758.65+	12.06+	122.98	370.04	SLSQP-AS	7.59+	1685.66	183.06+	899.70	SLSQP-AS	9.26+	62.53	125.58+	725.18	SLSQP-AS	7.20+	16.87	65.61+	541.54
SMAC-AS	760.88	13.24	127.12	375.20	SMAC-AS	10.38	1685.50	187.32	905.28	SMAC-AS	11.89	20.66+	130.36	730.59	SMAC-AS	9.76	16.06	70.37	546.53
SBS	5.58	9.79	4.44	6.67	SBS	5.58	9.79	4.44	6.67	SBS	5.58	9.79	4.44	6.67	SBS	5.58	9.79	4.44	6.67

(i) LOIO-CV ( $\mathcal{A}_{jped}$ )					(j) LOPO-CV ( $\mathcal{A}_{jped}$ )					(k) LOPOAD-CV ( $\mathcal{A}_{jped}$ )					(l) RI-CV ( $\mathcal{A}_{jped}$ )				
System	$n = 2$	$n = 3$	$n = 5$	$n = 10$	System	$n = 2$	$n = 3$	$n = 5$	$n = 10$	System	$n = 2$	$n = 3$	$n = 5$	$n = 10$	System	$n = 2$	$n = 3$	$n = 5$	$n = 10$
AS-50n	762.17	13.82	288.74	363.34	AS-50n	11.79	1732.04	431.24	723.13	AS-50n	15.06	142.39	292.87	728.71	AS-50n	10.15	15.68	149.28	188.46
AS-100n	764.68–	21.81–	291.68	192.17	AS-100n	14.41–	3348.73–	433.68–	561.51	AS-100n	17.38	147.02	296.18	730.28	AS-100n	12.54–	25.19–	151.08	191.96
SLSQP-AS	759.20+	11.24+	286.17	360.85	SLSQP-AS	7.99+	1729.63	428.29+	720.14	SLSQP-AS	11.78+	139.85	289.53+	726.02	SLSQP-AS	7.21+	12.14+	146.73+	185.87
SMAC-AS	761.92	12.75	290.39	366.09	SMAC-AS	11.32	1686.70	432.81	725.85	SMAC-AS	14.73	21.18+	294.38	731.56	SMAC-AS	9.99	13.95	150.77	191.28
SBS	6.05	9.81	4.50	6.72	SBS	6.05	9.81	4.50	6.72	SBS	6.05	9.81	4.50	6.72	SBS	6.05	9.81	4.50	6.72

(m) LOIO-CV ( $\mathcal{A}_{bmtpt}$ )					(n) LOPO-CV ( $\mathcal{A}_{bmtpt}$ )					(o) LOPOAD-CV ( $\mathcal{A}_{bmtpt}$ )					(p) RI-CV ( $\mathcal{A}_{bmtpt}$ )				
System	$n = 2$	$n = 3$	$n = 5$	$n = 10$	System	$n = 2$	$n = 3$	$n = 5$	$n = 10$	System	$n = 2$	$n = 3$	$n = 5$	$n = 10$	System	$n = 2$	$n = 3$	$n = 5$	$n = 10$
AS-50n	18.67	245.56	280.60	472.73	AS-50n	71.05	710.99	561.92	781.26	AS-50n	135.26	1171.26	557.19	1089.57	AS-50n	17.96	249.09	288.69	470.16
AS-100n	25.99	473.96–	284.69	471.36	AS-100n	69.49	944.60–	559.23	781.36	AS-100n	143.88	1403.79	568.98	1092.12	AS-100n	19.49	478.87–	286.96	316.55
SLSQP-AS	16.31	243.84	279.09	468.12	SLSQP-AS	68.56	708.94	559.82	778.11	SLSQP-AS	133.05	1169.20	555.16	1086.80	SLSQP-AS	15.67	247.25	286.60	467.83
SMAC-AS	17.38	246.04	281.62	474.33	SMAC-AS	70.72	710.81	558.96	782.83	SMAC-AS	133.23	963.42	558.31	1091.14	SMAC-AS	17.09	249.35	289.76	471.74
SBS	11.53	35.92	102.10	315.32	SBS	11.53	35.92	102.10	315.32	SBS	11.53	35.92	102.10	315.32	SBS	11.53	35.92	102.10	315.32

(q) LOIO-CV ( $\mathcal{A}_{mk}$ )					(r) LOPO-CV ( $\mathcal{A}_{mk}$ )					(s) LOPOAD-CV ( $\mathcal{A}_{mk}$ )					(t) RI-CV ( $\mathcal{A}_{mk}$ )				
System	$n = 2$	$n = 3$	$n = 5$	$n = 10$	System	$n = 2$	$n = 3$	$n = 5$	$n = 10$	System	$n = 2$	$n = 3$	$n = 5$	$n = 10$	System	$n = 2$	$n = 3$	$n = 5$	$n = 10$
AS-50n	3.05	239.94	130.58	66.68	AS-50n	3.28	240.12	194.70	99.41	AS-50n	11.50	31.00	67.80	100.27	AS-50n	3.09	8.91	131.70	66.71
AS-100n	4.44–	8.69	132.60	67.53–	AS-100n	4.53–	240.40	256.47–	100.56	AS-100n	13.09–	32.37–	68.37–	100.03	AS-100n	4.41–	8.45	134.69	67.20–
SLSQP-AS	1.93+	239.07	129.71	65.94+	SLSQP-AS	2.02+	239.08+	193.85	98.54	SLSQP-AS	8.99+	29.49+	66.54+	99.47	SLSQP-AS	1.95+	8.01	130.69	65.95+
SMAC-AS	3.13–	8.57+	131.00	67.19–	SMAC-AS	3.23	7.95+	195.10	99.92	SMAC-AS	11.05	30.99	68.27	100.78	SMAC-AS	3.22	8.45	132.17	67.23–
SBS	1.64	31.99	68.74	67.44	SBS	1.64	31.99	68.74	67.44	SBS	1.64	31.99	68.74	67.44	SBS	1.64	31.99	68.74	67.44

TABLE S.15: Median of 31 “mean relSP1” values of the five algorithm selection system with different sample size. “AS-50n”, “AS-25n”, “AS-100n”, and “AS-200n” are the algorithm selection system with the sample size  $50 \times n$ ,  $25 \times n$ ,  $100 \times n$ , and  $200 \times n$ , respectively.  $\mathcal{A}_{kt}$  was used in this comparison.

(a) LOIO-CV (Classification)					(b) LOPO-CV (Classification)					(c) LOPOAD-CV (Classification)					(d) RI-CV (Classification)				
System	$n = 2$	$n = 3$	$n = 5$	$n = 10$	System	$n = 2$	$n = 3$	$n = 5$	$n = 10$	System	$n = 2$	$n = 3$	$n = 5$	$n = 10$	System	$n = 2$	$n = 3$	$n = 5$	$n = 10$
AS-50n	7.54	10.10	4.75	4.99	AS-50n	33.77	6686.64	429.01	367.30	AS-50n	12.07	1695.37	288.15	361.93	AS-50n	7.63	6.57	4.38	4.73
AS-25n	7.34	9.46	3.39+	3.45+	AS-25n	763.59	6683.24	426.90+	364.83+	AS-25n	11.43	3338.98	426.24	360.37	AS-25n	7.47	7.41	3.16+	3.15+
AS-100n	10.84-	12.86-	7.79-	8.39-	AS-100n	33.53	6683.59	572.14-	370.59	AS-100n	13.92-	3351.36-	430.87-	365.25-	AS-100n	10.15-	10.11-	7.34-	8.04-
AS-200n	16.23-	18.57-	13.72-	14.98-	AS-200n	784.39-	6688.03	578.21-	550.86-	AS-200n	20.82-	3360.46-	436.98-	371.99-	AS-200n	15.91-	15.64-	13.20-	14.67-
SBS	5.81	9.82	4.49	6.76	SBS	5.81	9.82	4.49	6.76	SBS	5.81	9.82	4.49	6.76	SBS	5.81	9.82	4.49	6.76

(e) LOIO-CV (Regression)					(f) LOPO-CV (Regression)					(g) LOPOAD-CV (Regression)					(h) RI-CV (Regression)				
System	$n = 2$	$n = 3$	$n = 5$	$n = 10$	System	$n = 2$	$n = 3$	$n = 5$	$n = 10$	System	$n = 2$	$n = 3$	$n = 5$	$n = 10$	System	$n = 2$	$n = 3$	$n = 5$	$n = 10$
AS-50n	9.48	67.95	7.18	9.98	AS-50n	16.73	134.31	7.74	10.08	AS-50n	10.11	191.74	8.43	10.20	AS-50n	6.31	11.96	6.82	9.98
AS-25n	8.62	67.11	5.90+	8.36+	AS-25n	22.88	120.51	6.11+	8.47+	AS-25n	13.32	190.08	6.17+	8.55+	AS-25n	5.72	10.76	5.65+	8.35+
AS-100n	10.70-	70.79	10.23-	13.21-	AS-100n	24.55-	136.81	10.81-	13.44-	AS-100n	15.66-	169.12	10.17-	13.52-	AS-100n	8.99-	13.53-	9.78-	13.10-
AS-200n	16.01-	78.58-	15.92-	19.69-	AS-200n	40.89-	93.39	16.83-	20.19-	AS-200n	23.75-	198.55	16.47-	20.18-	AS-200n	14.70-	19.11-	15.75-	19.70-
SBS	5.81	9.82	4.49	6.76	SBS	5.81	9.82	4.49	6.76	SBS	5.81	9.82	4.49	6.76	SBS	5.81	9.82	4.49	6.76

(i) LOIO-CV (P-classification)					(j) LOPO-CV (P-classification)					(k) LOPOAD-CV (P-classification)					(l) RI-CV (P-classification)				
System	$n = 2$	$n = 3$	$n = 5$	$n = 10$	System	$n = 2$	$n = 3$	$n = 5$	$n = 10$	System	$n = 2$	$n = 3$	$n = 5$	$n = 10$	System	$n = 2$	$n = 3$	$n = 5$	$n = 10$
AS-50n	5.94	8.13	5.01	5.28	AS-50n	15.75	3337.17	146.83	363.62	AS-50n	9.76	131.42	5.42	5.29	AS-50n	5.43	6.86	5.10	5.44
AS-25n	7.96	23.61	3.77+	3.52+	AS-25n	15.16	3336.94	146.30	361.43	AS-25n	8.56	116.85	4.20+	181.68-	AS-25n	4.12+	6.48	3.46+	3.73+
AS-100n	8.75-	11.33-	7.88-	8.59-	AS-100n	16.66	3340.90	149.85-	366.78	AS-100n	11.22-	124.18	8.01-	8.60-	AS-100n	8.37-	10.98-	7.85-	8.64-
AS-200n	14.39-	51.43-	13.82-	15.08-	AS-200n	23.83-	1796.86	155.72-	373.55-	AS-200n	18.12-	86.35	13.91-	193.29-	AS-200n	14.51-	15.14-	13.88-	15.25-
SBS	5.81	9.82	4.49	6.76	SBS	5.81	9.82	4.49	6.76	SBS	5.81	9.82	4.49	6.76	SBS	5.81	9.82	4.49	6.76

(m) LOIO-CV (P-regression)					(n) LOPO-CV (P-regression)					(o) LOPOAD-CV (P-regression)					(p) RI-CV (P-regression)				
System	$n = 2$	$n = 3$	$n = 5$	$n = 10$	System	$n = 2$	$n = 3$	$n = 5$	$n = 10$	System	$n = 2$	$n = 3$	$n = 5$	$n = 10$	System	$n = 2$	$n = 3$	$n = 5$	$n = 10$
AS-50n	18.65	72.91	9.48	11.10	AS-50n	42.90	246.10	150.56	188.74	AS-50n	19.30	144.11	8.74	11.35	AS-50n	16.14	73.11	8.98	11.20
AS-25n	16.57	72.58	7.98+	9.57+	AS-25n	24.51+	459.41	149.45+	187.35+	AS-25n	15.63	138.43+	7.57+	10.07+	AS-25n	16.90	67.00	7.44+	9.58+
AS-100n	25.46	77.57	12.19-	14.60-	AS-100n	44.35	1685.43-	293.02-	192.32-	AS-100n	19.12	148.93	11.56-	15.32-	AS-100n	26.69-	71.44	11.99-	14.65-
AS-200n	34.93-	80.70	18.57-	20.96-	AS-200n	778.91-	1694.53-	299.44-	199.00-	AS-200n	25.03-	160.18-	17.91-	23.11-	AS-200n	27.46-	61.64	18.20-	21.00-
SBS	5.81	9.82	4.49	6.76	SBS	5.81	9.82	4.49	6.76	SBS	5.81	9.82	4.49	6.76	SBS	5.81	9.82	4.49	6.76

(q) LOIO-CV (Clustering)					(r) LOPO-CV (Clustering)					(s) LOPOAD-CV (Clustering)					(t) RI-CV (Clustering)				
System	$n = 2$	$n = 3$	$n = 5$	$n = 10$	System	$n = 2$	$n = 3$	$n = 5$	$n = 10$	System	$n = 2$	$n = 3$	$n = 5$	$n = 10$	System	$n = 2$	$n = 3$	$n = 5$	$n = 10$
AS-50n	761.68	15.39	289.31	373.30	AS-50n	11.57	1745.19	431.31	903.53	AS-50n	12.79	65.33	292.15	728.66	AS-50n	9.94	21.44	150.46	544.56
AS-25n	14.59+	12.77+	290.42	907.90-	AS-25n	9.75+	1683.60+	430.48	903.52	AS-25n	10.84	57.76	292.26	903.10	AS-25n	8.53+	16.57+	148.58+	551.80
AS-100n	764.26-	22.39-	291.94	552.62	AS-100n	14.41-	3349.86-	433.75-	743.96	AS-100n	14.86-	522.32	296.93-	909.25-	AS-100n	12.52-	1673.66-	151.76-	374.35
AS-200n	771.55-	1681.08-	161.99	377.17	AS-200n	25.41-	3365.33-	438.28-	731.18	AS-200n	21.60-	1680.42-	302.06-	911.14	AS-200n	19.81-	1681.66-	156.30	203.05+
SBS	5.81	9.82	4.49	6.76	SBS	5.81	9.82	4.49	6.76	SBS	5.81	9.82	4.49	6.76	SBS	5.81	9.82	4.49	6.76

TABLE S.16: Median of 31 “mean relSP1” values of the five algorithm selection systems using the sample  $\mathcal{X}$  and the union of  $\mathcal{X}$  and  $\mathcal{Y}$ , which a set of all solutions found so far by SLSQP. “S-AS” is the algorithm selection system that uses SLSQP as a pre-solver and computes features based on  $\mathcal{X}$ . “S-AS-U” is the algorithm selection system that uses SLSQP as a pre-solver and computes features based on  $\mathcal{X} \cup \mathcal{Y}$ .  $\mathcal{A}_{kt}$  was used in this comparison.

(a) LOIO-CV (Classification)					(b) LOPO-CV (Classification)					(c) LOPOAD-CV (Classification)					(d) RI-CV (Classification)				
System	$n = 2$	$n = 3$	$n = 5$	$n = 10$	System	$n = 2$	$n = 3$	$n = 5$	$n = 10$	System	$n = 2$	$n = 3$	$n = 5$	$n = 10$	System	$n = 2$	$n = 3$	$n = 5$	$n = 10$
S-AS	759.25	12.25	286.81	370.24	S-AS	8.65	1741.12	428.40	899.90	S-AS	9.42	62.57	289.29	725.37	S-AS	7.01	16.93	147.84	541.73
S-AS-U	89.33	9.70+	148.15+	182.65+	S-AS-U	13.25-	18.22+	427.64	539.17+	S-AS-U	19.14-	1704.11-	567.81-	717.44+	S-AS-U	8.70-	10.78+	146.29	7.99+
SBS	5.81	9.82	4.49	6.76	SBS	5.81	9.82	4.49	6.76	SBS	5.81	9.82	4.49	6.76	SBS	5.81	9.82	4.49	6.76

(e) LOIO-CV (Regression)					(f) LOPO-CV (Regression)					(g) LOPOAD-CV (Regression)					(h) RI-CV (Regression)				
System	$n = 2$	$n = 3$	$n = 5$	$n = 10$	System	$n = 2$	$n = 3$	$n = 5$	$n = 10$	System	$n = 2$	$n = 3$	$n = 5$	$n = 10$	System	$n = 2$	$n = 3$	$n = 5$	$n = 10$
S-AS	7.02	65.50	4.23	6.88	S-AS	10.83	131.85	4.58	6.90	S-AS	7.08	189.23	5.44	6.99	S-AS	3.54	9.58	3.97	6.86
S-AS-U	6.54	10.60+	4.72-	2.58+	S-AS-U	30.97-	121.95	4.19+	6.73+	S-AS-U	10.26-	119.31+	5.01	6.96	S-AS-U	4.64-	8.35	3.98	2.61+
SBS	5.81	9.82	4.49	6.76	SBS	5.81	9.82	4.49	6.76	SBS	5.81	9.82	4.49	6.76	SBS	5.81	9.82	4.49	6.76

(i) LOIO-CV (P-classification)					(j) LOPO-CV (P-classification)					(k) LOPOAD-CV (P-classification)					(l) RI-CV (P-classification)				
System	$n = 2$	$n = 3$	$n = 5$	$n = 10$	System	$n = 2$	$n = 3$	$n = 5$	$n = 10$	System	$n = 2$	$n = 3$	$n = 5$	$n = 10$	System	$n = 2$	$n = 3$	$n = 5$	$n = 10$
S-AS	3.53	4.85	2.54	2.49	S-AS	11.58	3334.35	144.06	360.35	S-AS	7.36	128.10	2.92	2.50	S-AS	2.98	4.25	2.65	2.66
S-AS-U	3.16	107.32-	2.80-	2.28+	S-AS-U	9.83	1730.36+	3.96+	183.07+	S-AS-U	7.80	216.73	2.87	2.52	S-AS-U	2.88	4.55	2.72	3.69
SBS	5.81	9.82	4.49	6.76	SBS	5.81	9.82	4.49	6.76	SBS	5.81	9.82	4.49	6.76	SBS	5.81	9.82	4.49	6.76

(m) LOIO-CV (P-regression)					(n) LOPO-CV (P-regression)					(o) LOPOAD-CV (P-regression)					(p) RI-CV (P-regression)				
System	$n = 2$	$n = 3$	$n = 5$	$n = 10$	System	$n = 2$	$n = 3$	$n = 5$	$n = 10$	System	$n = 2$	$n = 3$	$n = 5$	$n = 10$	System	$n = 2$	$n = 3$	$n = 5$	$n = 10$
S-AS	14.53	69.91	5.45	7.35	S-AS	33.00	242.60	146.98	185.27	S-AS	14.96	141.01	5.07	7.68	S-AS	11.25	68.85	5.18	7.38
S-AS-U	19.40	178.01-	5.37	7.95-	S-AS-U	763.00	365.94	146.72	7.67+	S-AS-U	19.91-	187.13	5.08	8.34-	S-AS-U	16.99-	76.16	5.46	7.86-
SBS	5.81	9.82	4.49	6.76	SBS	5.81	9.82	4.49	6.76	SBS	5.81	9.82	4.49	6.76	SBS	5.81	9.82	4.49	6.76

(q) LOIO-CV (Clustering)					(r) LOPO-CV (Clustering)					(s) LOPOAD-CV (Clustering)					(t) RI-CV (Clustering)				
System	$n = 2$	$n = 3$	$n = 5$	$n = 10$	System	$n = 2$	$n = 3$	$n = 5$	$n = 10$	System	$n = 2$	$n = 3$	$n = 5$	$n = 10$	System	$n = 2$	$n = 3$	$n = 5$	$n = 10$
S-AS	759.25	12.25	286.81	370.24	S-AS	8.65	1741.12	428.40	899.90	S-AS	9.42	62.57	289.29	725.37	S-AS	7.01	16.93	147.84	541.73
S-AS-U	89.33	9.70+	148.15+	182.65+	S-AS-U	13.25-	18.22+	427.64	539.17+	S-AS-U	19.14-	1704.11-	567.81-	717.44+	S-AS-U	8.70-	10.78+	146.29	7.99+
SBS	5.81	9.82	4.49	6.76	SBS	5.81	9.82	4.49	6.76	SBS	5.81	9.82	4.49	6.76	SBS	5.81	9.82	4.49	6.76

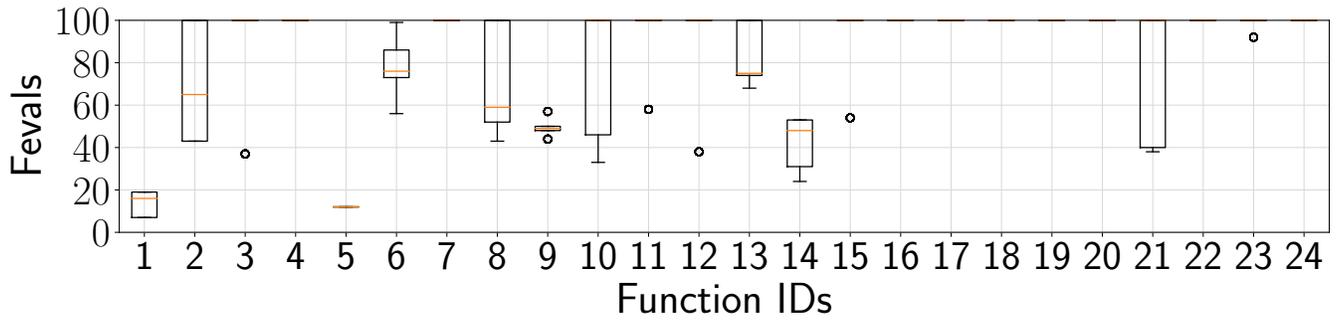
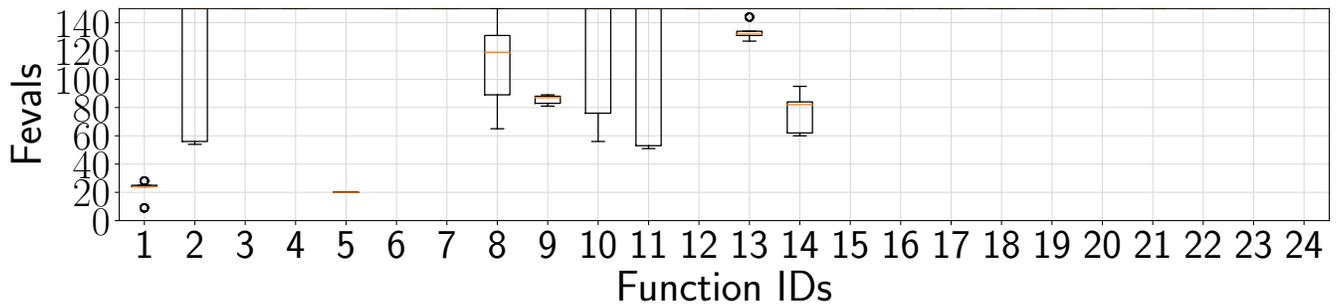
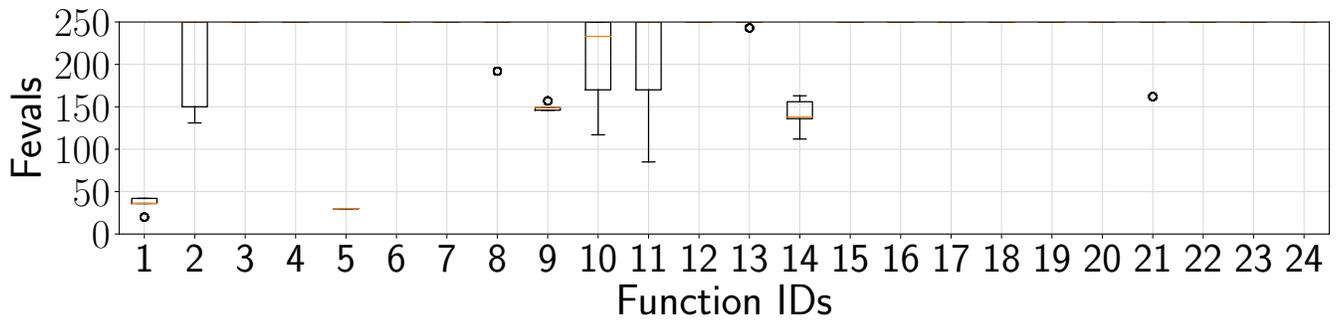
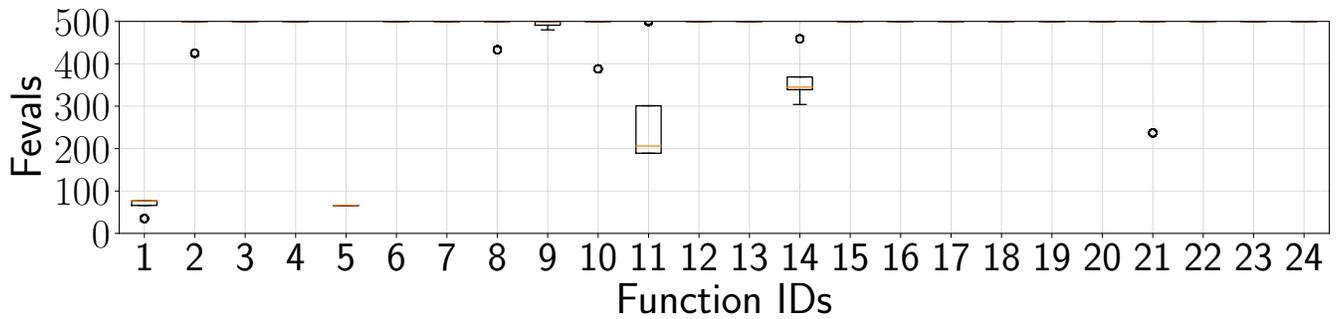
(a)  $n = 2$ (b)  $n = 3$ (c)  $n = 5$ (d)  $n = 10$ 

Fig. S.18: Distribution of the number of function evaluations used in the pre-solving phase (SLSQP). For each function, the results on 5 instances over 31 runs are shown.

TABLE S.17: Results of the five algorithm selection systems for  $n \in \{2, 3, 5, 10\}$ . Tables (a)–(o) show the performance score values of the five systems using the five algorithm portfolios ( $\mathcal{A}_{kt}$ , ...,  $\mathcal{A}_{mk}$ ) for the four cross-validation methods, respectively.

(a) LOIO-CV ( $\mathcal{A}_{kt}$ )					(b) LOPO-CV ( $\mathcal{A}_{kt}$ )					(c) LOPOAD-CV ( $\mathcal{A}_{kt}$ )					(d) RI-CV ( $\mathcal{A}_{kt}$ )				
$n$					$n$					$n$					$n$				
2 3 5 10					2 3 5 10					2 3 5 10					2 3 5 10				
Classification	1	1	0	0	Classification	3	4	3	3	Classification	2	4	3	3	Classification.	2	0	0	0
Regression	2	3	2	2	Regression	1	0	0	0	Regression	0	0	1	0	Regression	1	2	1	2
P-classification	0	0	1	1	P-classification	0	2	1	2	P-classification	0	0	0	0	P-classification	0	1	1	1
P-regression	3	4	3	3	P-regression	3	0	2	1	P-regression	4	1	1	1	P-regression	4	3	2	3
Clustering	3	2	4	4	Clustering	0	2	3	4	Clustering	0	0	3	4	Clustering	3	3	4	4
(e) LOIO-CV ( $\mathcal{A}_{dlvat}$ )					(f) LOPO-CV ( $\mathcal{A}_{dlvat}$ )					(g) LOPOAD-CV ( $\mathcal{A}_{dlvat}$ )					(h) RI-CV ( $\mathcal{A}_{dlvat}$ )				
$n$					$n$					$n$					$n$				
2 3 5 10					2 3 5 10					2 3 5 10					2 3 5 10				
Classification.	1	0	0	0	Classification.	1	4	3	2	Classification.	0	0	3	3	Classification.	2	0	0	0
Regression	2	3	2	2	Regression	1	0	0	0	Regression	0	0	0	0	Regression	1	2	1	2
P-classification	0	0	1	1	P-classification	1	0	2	2	P-classification	0	0	2	0	P-classification	0	1	1	1
P-regression	3	3	3	3	P-regression	3	0	1	1	P-regression	4	0	0	1	P-regression	4	3	2	3
Clustering	3	2	4	4	Clustering	0	3	3	4	Clustering	1	0	4	4	Clustering	3	3	4	4
(i) LOIO-CV ( $\mathcal{A}_{jped}$ )					(j) LOPO-CV ( $\mathcal{A}_{jped}$ )					(k) LOPOAD-CV ( $\mathcal{A}_{jped}$ )					(l) RI-CV ( $\mathcal{A}_{jped}$ )				
$n$					$n$					$n$					$n$				
2 3 5 10					2 3 5 10					2 3 5 10					2 3 5 10				
Classification	1	1	0	0	Classification	4	4	3	3	Classification	1	4	4	3	Classification.	2	0	0	0
Regression	1	4	2	2	Regression	1	1	0	0	Regression	1	1	1	0	Regression	1	2	2	2
P-classification	0	0	0	1	P-classification	1	2	1	2	P-classification	0	0	0	0	P-classification	0	1	1	1
P-regression	3	3	3	2	P-regression	2	0	2	1	P-regression	2	0	1	0	P-regression	4	3	3	2
Clustering	3	2	4	3	Clustering	0	3	3	4	Clustering	1	0	3	4	Clustering	2	3	4	4
(m) LOIO-CV ( $\mathcal{A}_{bmtpt}$ )					(n) LOPO-CV ( $\mathcal{A}_{bmtpt}$ )					(o) LOPOAD-CV ( $\mathcal{A}_{bmtpt}$ )					(p) RI-CV ( $\mathcal{A}_{bmtpt}$ )				
$n$					$n$					$n$					$n$				
2 3 5 10					2 3 5 10					2 3 5 10					2 3 5 10				
Classification	0	3	0	2	Classification	0	2	3	3	Classification	3	0	0	0	Classification.	0	3	0	0
Regression	3	1	2	0	Regression	0	0	0	0	Regression	0	3	3	3	Regression	2	1	3	1
P-classification	0	0	1	0	P-classification	0	2	0	2	P-classification	0	0	2	1	P-classification	1	1	1	3
P-regression	0	0	1	0	P-regression	1	1	0	0	P-regression	0	2	1	2	P-regression	0	0	1	1
Clustering	4	4	4	4	Clustering	4	2	4	4	Clustering	4	4	4	4	Clustering	4	3	4	4
(q) LOIO-CV ( $\mathcal{A}_{mk}$ )					(r) LOPO-CV ( $\mathcal{A}_{mk}$ )					(s) LOPOAD-CV ( $\mathcal{A}_{mk}$ )					(t) RI-CV ( $\mathcal{A}_{mk}$ )				
$n$					$n$					$n$					$n$				
2 3 5 10					2 3 5 10					2 3 5 10					2 3 5 10				
Classification	1	0	0	2	Classification	2	3	2	0	Classification	0	0	0	0	Classification.	0	0	0	0
Regression	1	3	2	1	Regression	0	0	0	0	Regression	3	2	4	1	Regression	1	2	2	2
P-classification	3	1	2	0	P-classification	0	2	3	4	P-classification	0	1	0	1	P-classification	1	0	2	3
P-regression	4	2	1	0	P-regression	4	0	0	0	P-regression	2	4	2	1	P-regression	4	2	1	1
Clustering	0	2	4	4	Clustering	2	0	4	3	Clustering	4	1	2	4	Clustering	3	2	4	3

TABLE S.18: Results of the five algorithm selection systems for  $n \in \{2, 3, 5, 10\}$ . Tables (a)–(o) show the performance score values of the five systems using the six algorithm portfolios ( $\mathcal{A}_{1s2}, \dots, \mathcal{A}_{1s12}$ ) for the four cross-validation methods, respectively.

(a) LOIO-CV ( $\mathcal{A}_{1s2}$ )					(b) LOPO-CV ( $\mathcal{A}_{1s2}$ )					(c) LOPOAD-CV ( $\mathcal{A}_{1s2}$ )					(d) RI-CV ( $\mathcal{A}_{1s2}$ )				
$n$					$n$					$n$					$n$				
2 3 5 10					2 3 5 10					2 3 5 10					2 3 5 10				
Classification.	0	1	0	0	Classification.	4	2	0	2	Classification.	2	0	0	0	Classification.	0	0	0	0
Regression	3	0	2	2	Regression	0	1	1	0	Regression	0	2	2	2	Regression	1	0	1	2
P-classification	0	1	0	0	P-classification	0	0	3	2	P-classification	2	0	0	0	P-classification	1	0	3	2
P-regression	3	0	2	2	P-regression	0	2	1	0	P-regression	0	2	3	2	P-regression	1	0	1	2
Clustering	0	0	4	2	Clustering	0	0	4	2	Clustering	2	2	3	4	Clustering	4	2	3	1
(e) LOIO-CV ( $\mathcal{A}_{1s4}$ )					(f) LOPO-CV ( $\mathcal{A}_{1s4}$ )					(g) LOPOAD-CV ( $\mathcal{A}_{1s4}$ )					(h) RI-CV ( $\mathcal{A}_{1s4}$ )				
$n$					$n$					$n$					$n$				
2 3 5 10					2 3 5 10					2 3 5 10					2 3 5 10				
Classification.	2	0	0	0	Classification.	0	0	3	3	Classification.	0	0	2	2	Classification.	0	0	0	0
Regression	3	2	2	3	Regression	0	0	0	0	Regression	1	1	1	0	Regression	2	1	1	2
P-classification	1	3	0	0	P-classification	0	2	2	2	P-classification	1	1	0	2	P-classification	0	1	1	1
P-regression	1	4	3	3	P-regression	4	4	1	0	P-regression	3	3	1	1	P-regression	4	4	2	2
Clustering	0	1	2	2	Clustering	0	0	4	3	Clustering	3	1	3	4	Clustering	3	1	4	2
(i) LOIO-CV ( $\mathcal{A}_{1s6}$ )					(j) LOPO-CV ( $\mathcal{A}_{1s6}$ )					(k) LOPOAD-CV ( $\mathcal{A}_{1s6}$ )					(l) RI-CV ( $\mathcal{A}_{1s6}$ )				
$n$					$n$					$n$					$n$				
2 3 5 10					2 3 5 10					2 3 5 10					2 3 5 10				
Classification.	2	4	0	1	Classification.	0	4	3	4	Classification.	0	0	4	3	Classification.	0	0	0	0
Regression	1	1	0	2	Regression	2	0	0	0	Regression	2	1	1	0	Regression	1	2	1	2
P-classification	1	1	0	0	P-classification	1	2	2	2	P-classification	1	0	0	2	P-classification	1	1	1	1
P-regression	0	2	3	3	P-regression	1	0	1	1	P-regression	3	1	2	1	P-regression	3	2	3	3
Clustering	2	0	4	3	Clustering	3	2	4	3	Clustering	3	2	3	4	Clustering	3	4	4	2
(m) LOIO-CV ( $\mathcal{A}_{1s8}$ )					(n) LOPO-CV ( $\mathcal{A}_{1s8}$ )					(o) LOPOAD-CV ( $\mathcal{A}_{1s8}$ )					(p) RI-CV ( $\mathcal{A}_{1s8}$ )				
$n$					$n$					$n$					$n$				
2 3 5 10					2 3 5 10					2 3 5 10					2 3 5 10				
Classification.	0	4	0	0	Classification.	0	4	3	2	Classification.	0	3	3	2	Classification.	0	0	0	0
Regression	4	1	2	2	Regression	2	0	1	1	Regression	2	1	2	0	Regression	1	2	2	2
P-classification	1	0	1	1	P-classification	1	3	2	3	P-classification	0	0	0	3	P-classification	2	0	1	3
P-regression	2	2	3	1	P-regression	1	1	0	0	P-regression	3	1	0	0	P-regression	3	4	2	1
Clustering	3	0	3	4	Clustering	1	2	3	3	Clustering	4	3	4	4	Clustering	3	2	4	4
(q) LOIO-CV ( $\mathcal{A}_{1s10}$ )					(r) LOPO-CV ( $\mathcal{A}_{1s10}$ )					(s) LOPOAD-CV ( $\mathcal{A}_{1s10}$ )					(t) RI-CV ( $\mathcal{A}_{1s10}$ )				
$n$					$n$					$n$					$n$				
2 3 5 10					2 3 5 10					2 3 5 10					2 3 5 10				
Classification.	0	4	0	0	Classification.	1	3	2	2	Classification.	3	3	3	2	Classification.	0	0	0	0
Regression	2	0	2	1	Regression	2	0	0	0	Regression	1	1	1	0	Regression	1	1	1	1
P-classification	1	0	0	3	P-classification	4	3	4	2	P-classification	0	0	0	2	P-classification	2	3	1	3
P-regression	3	2	3	1	P-regression	0	0	0	0	P-regression	2	1	0	0	P-regression	2	3	2	2
Clustering	3	0	4	4	Clustering	0	2	3	2	Clustering	3	3	4	4	Clustering	2	1	4	4
(u) LOIO-CV ( $\mathcal{A}_{1s12}$ )					(v) LOPO-CV ( $\mathcal{A}_{1s12}$ )					(w) LOPOAD-CV ( $\mathcal{A}_{1s12}$ )					(x) RI-CV ( $\mathcal{A}_{1s12}$ )				
$n$					$n$					$n$					$n$				
2 3 5 10					2 3 5 10					2 3 5 10					2 3 5 10				
Classification.	2	4	0	0	Classification.	4	4	4	2	Classification.	3	4	3	2	Classification.	2	0	0	0
Regression	1	0	2	1	Regression	2	0	0	0	Regression	1	1	0	0	Regression	0	1	1	1
P-classification	0	0	0	3	P-classification	2	3	3	2	P-classification	0	0	1	2	P-classification	0	1	3	3
P-regression	3	3	3	1	P-regression	0	0	0	0	P-regression	2	1	0	0	P-regression	4	3	2	1
Clustering	2	1	4	4	Clustering	0	2	2	4	Clustering	3	3	4	4	Clustering	2	4	4	4

TABLE S.19: Results of the five algorithm selection systems for  $n \in \{2, 3, 5, 10\}$ . Tables (a)–(o) show the performance score values of the five systems using the three algorithm portfolios ( $\mathcal{A}_{1s14}$ , ...,  $\mathcal{A}_{1s18}$ ) for the four cross-validation methods, respectively.

(a) LOIO-CV ( $\mathcal{A}_{1s14}$ )					(b) LOPO-CV ( $\mathcal{A}_{1s14}$ )					(c) LOPOAD-CV ( $\mathcal{A}_{1s14}$ )					(d) RI-CV ( $\mathcal{A}_{1s14}$ )				
$n$					$n$					$n$					$n$				
2 3 5 10					2 3 5 10					2 3 5 10					2 3 5 10				
Classification.	2	4	0	0	Classification.	4	4	4	2	Classification.	3	3	3	3	Classification.	2	0	0	0
Regression	1	1	2	1	Regression	3	0	0	0	Regression	1	0	2	0	Regression	0	2	1	1
P-classification	0	0	1	3	P-classification	0	3	1	3	P-classification	0	1	0	2	P-classification	0	1	1	3
P-regression	3	3	3	1	P-regression	0	0	1	1	P-regression	2	0	0	1	P-regression	3	4	3	1
Clustering	4	2	4	4	Clustering	0	2	3	3	Clustering	3	4	4	4	Clustering	2	3	4	4

(e) LOIO-CV ( $\mathcal{A}_{1s16}$ )					(f) LOPO-CV ( $\mathcal{A}_{1s16}$ )					(g) LOPOAD-CV ( $\mathcal{A}_{1s16}$ )					(h) RI-CV ( $\mathcal{A}_{1s16}$ )				
$n$					$n$					$n$					$n$				
2 3 5 10					2 3 5 10					2 3 5 10					2 3 5 10				
Classification.	3	1	0	0	Classification.	1	3	4	2	Classification.	3	1	3	3	Classification.	2	0	0	0
Regression	0	0	2	1	Regression	1	0	0	0	Regression	0	1	0	0	Regression	0	1	1	1
P-classification	2	2	1	2	P-classification	1	2	2	3	P-classification	0	1	1	2	P-classification	1	3	1	3
P-regression	0	1	3	2	P-regression	0	0	1	1	P-regression	0	0	2	1	P-regression	3	2	3	1
Clustering	3	3	4	4	Clustering	4	2	2	4	Clustering	4	4	4	4	Clustering	4	4	4	4

(i) LOIO-CV ( $\mathcal{A}_{1s18}$ )					(j) LOPO-CV ( $\mathcal{A}_{1s18}$ )					(k) LOPOAD-CV ( $\mathcal{A}_{1s18}$ )					(l) RI-CV ( $\mathcal{A}_{1s18}$ )				
$n$					$n$					$n$					$n$				
2 3 5 10					2 3 5 10					2 3 5 10					2 3 5 10				
Classification.	3	1	0	0	Classification.	1	4	4	2	Classification.	3	3	3	3	Classification.	2	0	0	0
Regression	0	0	2	1	Regression	1	0	0	0	Regression	0	1	0	0	Regression	0	1	1	1
P-classification	2	2	1	3	P-classification	1	2	1	2	P-classification	0	1	1	2	P-classification	1	3	1	3
P-regression	0	1	3	1	P-regression	0	0	2	1	P-regression	0	0	2	0	P-regression	3	2	2	1
Clustering	3	2	4	4	Clustering	4	2	3	4	Clustering	4	4	4	4	Clustering	4	4	4	4

TABLE S.20: Results of the five algorithm selection systems for  $n \in \{2, 3, 5, 10\}$ . Tables (a)–(o) show the Friedman test-based average rankings of the five systems using the five algorithm portfolios ( $\mathcal{A}_{kt}$ , ...,  $\mathcal{A}_{mk}$ ) for the four cross-validation methods, respectively.

(a) LOIO-CV ( $\mathcal{A}_{kt}$ )					(b) LOPO-CV ( $\mathcal{A}_{kt}$ )					(c) LOPOAD-CV ( $\mathcal{A}_{kt}$ )					(d) RI-CV ( $\mathcal{A}_{kt}$ )				
$n$					$n$					$n$					$n$				
	2	3	5	10		2	3	5	10		2	3	5	10		2	3	5	10
Classification.	1.19	1.52	2.53	2.53	Classification.	4.03	4.82	2.77	2.77	Classification.	2.84	4.71	3.71	2.74	Classification.	2.02	1.55	2.24	2.21
Regression	3.47	3.79	2.53	2.53	Regression	2.94	3.16	1.74	1.76	Regression	3.00	3.03	3.06	2.48	Regression	3.27	3.77	2.18	2.21
P-classification	2.11	1.65	2.53	2.53	P-classification	2.29	1.82	2.66	2.85	P-classification	2.97	1.71	1.79	2.48	P-classification	1.98	1.65	4.06	3.85
P-regression	3.92	4.26	2.53	2.53	P-regression	3.47	3.21	3.90	2.74	P-regression	3.65	2.98	2.52	2.48	P-regression	4.11	4.02	2.18	2.21
Clustering	4.31	3.79	4.87	4.87	Clustering	2.27	1.98	3.92	4.87	Clustering	2.55	2.56	3.92	4.81	Clustering	3.61	4.02	4.34	4.52
(e) LOIO-CV ( $\mathcal{A}_{dlvat}$ )					(f) LOPO-CV ( $\mathcal{A}_{dlvat}$ )					(g) LOPOAD-CV ( $\mathcal{A}_{dlvat}$ )					(h) RI-CV ( $\mathcal{A}_{dlvat}$ )				
$n$					$n$					$n$					$n$				
	2	3	5	10		2	3	5	10		2	3	5	10		2	3	5	10
Classification.	1.52	1.50	2.53	2.53	Classification.	3.60	4.89	2.95	2.35	Classification.	2.68	1.71	3.50	2.53	Classification.	1.50	1.45	2.19	2.21
Regression	3.27	3.85	2.53	2.53	Regression	2.84	3.23	1.77	1.71	Regression	2.81	3.87	2.97	2.47	Regression	3.34	3.87	2.19	2.21
P-classification	1.92	1.69	2.53	2.53	P-classification	2.87	1.77	3.42	2.81	P-classification	2.89	2.45	1.63	2.47	P-classification	1.74	1.61	4.06	3.85
P-regression	4.19	4.10	2.53	2.53	P-regression	3.94	3.00	2.66	3.29	P-regression	4.18	3.73	3.02	2.76	P-regression	4.76	3.71	2.19	2.21
Clustering	4.10	3.85	4.87	4.87	Clustering	1.76	2.11	4.19	4.84	Clustering	2.45	3.24	3.89	4.77	Clustering	3.66	4.35	4.35	4.52
(i) LOIO-CV ( $\mathcal{A}_{jped}$ )					(j) LOPO-CV ( $\mathcal{A}_{jped}$ )					(k) LOPOAD-CV ( $\mathcal{A}_{jped}$ )					(l) RI-CV ( $\mathcal{A}_{jped}$ )				
$n$					$n$					$n$					$n$				
	2	3	5	10		2	3	5	10		2	3	5	10		2	3	5	10
Classification.	1.44	1.53	1.50	2.47	Classification.	3.52	4.90	3.05	2.82	Classification.	2.58	4.18	3.69	3.74	Classification.	1.87	1.50	1.50	2.29
Regression	3.29	4.06	3.77	2.53	Regression	3.53	3.34	2.02	1.81	Regression	3.32	3.16	3.15	2.13	Regression	3.47	3.76	3.84	2.40
P-classification	2.21	1.66	1.65	2.47	P-classification	2.52	1.71	3.00	2.89	P-classification	2.06	1.71	2.31	1.74	P-classification	1.74	1.56	1.74	2.34
P-regression	3.68	4.02	3.16	3.45	P-regression	3.52	2.81	2.84	2.98	P-regression	4.11	2.68	1.97	3.11	P-regression	4.23	4.15	3.15	3.27
Clustering	4.39	3.73	4.92	4.08	Clustering	1.92	2.24	4.10	4.50	Clustering	2.92	3.27	3.89	4.27	Clustering	3.69	4.03	4.77	4.69
(m) LOIO-CV ( $\mathcal{A}_{bmtpt}$ )					(n) LOPO-CV ( $\mathcal{A}_{bmtpt}$ )					(o) LOPOAD-CV ( $\mathcal{A}_{bmtpt}$ )					(p) RI-CV ( $\mathcal{A}_{bmtpt}$ )				
$n$					$n$					$n$					$n$				
	2	3	5	10		2	3	5	10		2	3	5	10		2	3	5	10
Classification.	2.68	2.85	2.63	4.37	Classification.	2.77	2.45	2.47	4.26	Classification.	2.53	2.61	2.61	2.81	Classification.	2.56	2.52	2.31	2.77
Regression	2.68	2.85	2.85	2.52	Regression	2.29	2.45	2.31	2.24	Regression	2.61	2.61	2.61	2.81	Regression	2.56	2.52	2.31	2.61
P-classification	2.76	2.85	2.63	2.76	P-classification	3.84	4.94	3.98	3.85	P-classification	2.85	2.61	2.61	2.73	P-classification	3.56	4.37	4.27	3.66
P-regression	2.68	2.85	2.63	1.95	P-regression	2.81	2.53	2.21	1.84	P-regression	2.63	2.61	2.61	2.73	P-regression	2.56	2.52	2.31	2.61
Clustering	4.21	3.58	4.26	3.40	Clustering	3.29	2.63	4.03	2.81	Clustering	4.37	4.55	4.55	3.94	Clustering	3.74	3.08	3.81	3.34
(q) LOIO-CV ( $\mathcal{A}_{mk}$ )					(r) LOPO-CV ( $\mathcal{A}_{mk}$ )					(s) LOPOAD-CV ( $\mathcal{A}_{mk}$ )					(t) RI-CV ( $\mathcal{A}_{mk}$ )				
$n$					$n$					$n$					$n$				
	2	3	5	10		2	3	5	10		2	3	5	10		2	3	5	10
Classification.	3.61	2.32	2.56	4.53	Classification.	3.84	2.03	2.15	2.98	Classification.	2.84	2.40	2.79	2.92	Classification.	1.48	2.06	2.26	2.74
Regression	3.21	2.52	2.69	2.68	Regression	2.58	2.45	2.24	2.66	Regression	2.63	3.31	2.79	2.92	Regression	3.37	2.18	2.26	2.74
P-classification	3.74	2.65	2.56	3.16	P-classification	2.68	3.63	4.13	4.27	P-classification	3.65	3.15	2.87	3.08	P-classification	3.35	3.98	4.18	3.79
P-regression	2.60	2.87	2.56	2.35	P-regression	2.66	3.63	2.29	2.42	P-regression	3.02	4.44	2.79	2.92	P-regression	2.82	2.61	2.26	2.74
Clustering	1.84	4.65	4.61	2.27	Clustering	3.24	3.26	4.19	2.66	Clustering	2.87	1.71	3.76	3.16	Clustering	3.97	4.16	4.05	2.98

TABLE S.21: Results of the five algorithm selection systems for  $n \in \{2, 3, 5, 10\}$ . Tables (a)–(o) show the Friedman test-based average rankings of the five systems using the six algorithm portfolios ( $\mathcal{A}_{1s2}$ , ...,  $\mathcal{A}_{1s12}$ ) for the four cross-validation methods, respectively.

(a) LOIO-CV ( $\mathcal{A}_{1s2}$ )					(b) LOPO-CV ( $\mathcal{A}_{1s2}$ )					(c) LOPOAD-CV ( $\mathcal{A}_{1s2}$ )					(d) RI-CV ( $\mathcal{A}_{1s2}$ )				
$n$					$n$					$n$					$n$				
2 3 5 10					2 3 5 10					2 3 5 10					2 3 5 10				
Classification.	3.19	2.92	2.82	2.97	Classification.	3.82	4.00	2.52	2.60	Classification.	2.82	3.32	2.16	2.92	Classification.	3.03	2.95	2.55	2.56
Regression	3.45	2.32	2.82	2.97	Regression	2.74	3.02	2.60	2.60	Regression	3.11	2.89	4.08	2.92	Regression	2.89	2.31	2.55	2.56
P-classification	3.19	2.92	2.82	2.97	P-classification	3.23	1.98	4.29	4.35	P-classification	2.82	3.32	2.16	2.92	P-classification	2.95	2.63	4.47	4.53
P-regression	3.40	2.40	2.82	2.97	P-regression	2.23	3.00	2.52	2.60	P-regression	2.44	2.65	4.42	2.92	P-regression	2.81	2.47	2.55	2.56
Clustering	1.76	4.44	3.71	3.13	Clustering	2.98	3.00	3.08	2.85	Clustering	3.81	2.82	2.18	3.32	Clustering	3.32	4.65	2.89	2.77
(e) LOIO-CV ( $\mathcal{A}_{1s4}$ )					(f) LOPO-CV ( $\mathcal{A}_{1s4}$ )					(g) LOPOAD-CV ( $\mathcal{A}_{1s4}$ )					(h) RI-CV ( $\mathcal{A}_{1s4}$ )				
$n$					$n$					$n$					$n$				
2 3 5 10					2 3 5 10					2 3 5 10					2 3 5 10				
Classification.	3.23	1.84	2.56	2.90	Classification.	1.85	2.19	4.15	4.21	Classification.	1.44	3.21	4.47	3.13	Classification.	1.42	1.85	2.47	2.45
Regression	3.79	2.05	2.56	2.90	Regression	3.39	2.58	1.53	1.90	Regression	3.81	2.58	3.61	2.60	Regression	3.21	1.87	2.32	2.35
P-classification	3.27	2.35	2.56	2.90	P-classification	2.76	3.26	3.84	3.05	P-classification	2.65	3.00	2.00	2.60	P-classification	2.60	2.34	4.19	4.65
P-regression	2.98	4.90	2.56	2.90	P-regression	4.27	4.05	2.02	1.98	P-regression	3.24	3.06	1.69	2.66	P-regression	4.60	4.71	2.32	2.35
Clustering	1.73	3.85	4.74	3.39	Clustering	2.73	2.92	3.47	3.85	Clustering	3.87	3.15	3.23	4.02	Clustering	3.18	4.23	3.69	3.19
(i) LOIO-CV ( $\mathcal{A}_{1s6}$ )					(j) LOPO-CV ( $\mathcal{A}_{1s6}$ )					(k) LOPOAD-CV ( $\mathcal{A}_{1s6}$ )					(l) RI-CV ( $\mathcal{A}_{1s6}$ )				
$n$					$n$					$n$					$n$				
2 3 5 10					2 3 5 10					2 3 5 10					2 3 5 10				
Classification.	3.73	2.03	2.50	2.68	Classification.	1.65	3.45	4.06	4.31	Classification.	1.61	3.60	4.63	3.00	Classification.	1.27	1.65	2.23	2.35
Regression	3.23	2.35	2.50	2.68	Regression	3.65	1.76	1.73	1.76	Regression	3.95	3.02	2.06	2.50	Regression	2.90	1.73	2.18	2.27
P-classification	3.32	2.35	2.50	2.68	P-classification	2.37	3.47	3.60	2.76	P-classification	1.85	3.15	2.92	2.50	P-classification	2.55	3.44	3.95	4.44
P-regression	1.60	5.00	2.50	2.68	P-regression	3.35	4.00	1.77	2.23	P-regression	3.68	2.58	2.08	2.74	P-regression	4.35	4.77	2.18	2.27
Clustering	3.13	3.26	5.00	4.29	Clustering	3.98	2.32	3.84	3.95	Clustering	3.90	2.66	3.31	4.26	Clustering	3.92	3.42	4.47	3.66
(m) LOIO-CV ( $\mathcal{A}_{1s8}$ )					(n) LOPO-CV ( $\mathcal{A}_{1s8}$ )					(o) LOPOAD-CV ( $\mathcal{A}_{1s8}$ )					(p) RI-CV ( $\mathcal{A}_{1s8}$ )				
$n$					$n$					$n$					$n$				
2 3 5 10					2 3 5 10					2 3 5 10					2 3 5 10				
Classification.	2.08	2.27	2.44	2.05	Classification.	2.29	4.42	3.81	3.89	Classification.	2.31	4.73	4.06	3.60	Classification.	1.16	2.23	2.02	2.58
Regression	4.15	2.29	2.56	3.58	Regression	3.35	1.84	2.71	2.60	Regression	4.18	2.44	2.15	2.63	Regression	2.89	2.26	2.06	2.34
P-classification	2.13	2.16	2.50	2.94	P-classification	2.85	3.39	3.65	3.40	P-classification	1.94	2.85	3.18	2.31	P-classification	3.00	2.06	4.05	3.71
P-regression	3.13	4.97	2.50	2.45	P-regression	3.23	3.39	1.03	1.47	P-regression	3.76	2.45	1.24	2.06	P-regression	4.08	4.97	1.97	2.34
Clustering	3.52	3.31	5.00	3.98	Clustering	3.27	1.97	3.81	3.65	Clustering	2.82	2.53	4.37	4.40	Clustering	3.87	3.48	4.90	4.03
(q) LOIO-CV ( $\mathcal{A}_{1s10}$ )					(r) LOPO-CV ( $\mathcal{A}_{1s10}$ )					(s) LOPOAD-CV ( $\mathcal{A}_{1s10}$ )					(t) RI-CV ( $\mathcal{A}_{1s10}$ )				
$n$					$n$					$n$					$n$				
2 3 5 10					2 3 5 10					2 3 5 10					2 3 5 10				
Classification.	3.02	2.19	2.42	3.18	Classification.	4.05	4.21	3.61	4.06	Classification.	4.39	4.61	4.24	3.73	Classification.	1.29	1.90	1.98	2.03
Regression	2.11	2.44	2.61	3.18	Regression	2.79	1.73	2.26	2.16	Regression	3.44	2.32	3.18	2.94	Regression	3.06	1.85	2.08	2.50
P-classification	3.27	2.13	2.48	3.10	P-classification	3.16	4.27	3.15	3.71	P-classification	1.60	2.82	3.02	2.85	P-classification	3.18	4.11	4.00	3.87
P-regression	2.97	4.97	2.48	3.18	P-regression	2.03	2.74	2.37	1.42	P-regression	3.02	2.29	1.06	1.15	P-regression	3.84	4.50	2.03	1.95
Clustering	3.63	3.27	5.00	2.37	Clustering	2.97	2.05	3.61	3.65	Clustering	2.56	2.95	3.50	4.34	Clustering	3.63	2.63	4.90	4.65
(u) LOIO-CV ( $\mathcal{A}_{1s12}$ )					(v) LOPO-CV ( $\mathcal{A}_{1s12}$ )					(w) LOPOAD-CV ( $\mathcal{A}_{1s12}$ )					(x) RI-CV ( $\mathcal{A}_{1s12}$ )				
$n$					$n$					$n$					$n$				
2 3 5 10					2 3 5 10					2 3 5 10					2 3 5 10				
Classification.	3.24	2.11	2.53	3.02	Classification.	4.08	3.45	3.42	3.92	Classification.	3.40	4.58	4.15	3.84	Classification.	1.35	1.77	1.98	2.31
Regression	2.18	2.82	2.56	3.02	Regression	2.92	3.06	2.34	1.97	Regression	3.65	2.27	2.94	2.71	Regression	3.53	2.13	2.03	2.87
P-classification	3.32	2.11	2.45	3.02	P-classification	3.10	3.03	3.35	3.77	P-classification	1.61	2.79	2.85	2.69	P-classification	2.52	4.19	4.53	2.68
P-regression	2.60	4.90	2.45	3.02	P-regression	2.08	3.44	2.47	1.47	P-regression	3.81	2.42	1.11	1.29	P-regression	3.77	4.31	1.98	2.31
Clustering	3.66	3.05	5.00	2.94	Clustering	2.82	2.02	3.42	3.87	Clustering	2.53	2.94	3.95	4.47	Clustering	3.82	2.60	4.47	4.84

TABLE S.22: Results of the five algorithm selection systems for  $n \in \{2, 3, 5, 10\}$ . Tables (a)–(o) show the Friedman test-based average rankings of the five systems using the three algorithm portfolios ( $\mathcal{A}_{\text{Is14}}$ , ...,  $\mathcal{A}_{\text{Is18}}$ ) for the four cross-validation methods, respectively.

(a) LOIO-CV ( $\mathcal{A}_{\text{Is14}}$ )					(b) LOPO-CV ( $\mathcal{A}_{\text{Is14}}$ )					(c) LOPOAD-CV ( $\mathcal{A}_{\text{Is14}}$ )					(d) RI-CV ( $\mathcal{A}_{\text{Is14}}$ )				
$n$					$n$					$n$					$n$				
	2	3	5	10		2	3	5	10		2	3	5	10		2	3	5	10
Classification.	1.29	1.37	1.39	3.00	Classification.	4.13	3.21	3.35	3.63	Classification.	4.89	3.44	3.89	3.84	Classification.	1.55	1.44	1.18	2.03
Regression	3.29	3.39	3.29	3.00	Regression	3.31	2.44	2.21	1.81	Regression	2.65	2.48	3.18	2.66	Regression	2.81	3.21	2.68	2.60
P-classification	2.11	1.69	1.73	3.00	P-classification	2.37	3.60	2.95	3.55	P-classification	1.40	3.05	2.84	2.56	P-classification	2.27	1.60	3.06	3.68
P-regression	4.00	4.90	4.27	3.00	P-regression	2.08	3.11	3.77	2.39	P-regression	3.52	1.97	1.21	1.77	P-regression	4.39	4.84	3.56	2.00
Clustering	4.31	3.65	4.32	3.00	Clustering	3.11	2.65	2.71	3.63	Clustering	2.55	4.06	3.89	4.16	Clustering	3.98	3.92	4.52	4.69
(e) LOIO-CV ( $\mathcal{A}_{\text{Is16}}$ )					(f) LOPO-CV ( $\mathcal{A}_{\text{Is16}}$ )					(g) LOPOAD-CV ( $\mathcal{A}_{\text{Is16}}$ )					(h) RI-CV ( $\mathcal{A}_{\text{Is16}}$ )				
$n$					$n$					$n$					$n$				
	2	3	5	10		2	3	5	10		2	3	5	10		2	3	5	10
Classification.	4.11	1.56	1.45	3.00	Classification.	3.87	2.42	3.02	3.50	Classification.	4.71	2.55	3.40	3.39	Classification.	1.45	1.61	1.21	2.13
Regression	2.50	3.81	3.06	3.00	Regression	2.98	3.18	2.53	1.94	Regression	3.27	3.19	2.26	2.32	Regression	1.95	3.66	2.34	2.50
P-classification	3.53	1.56	1.74	3.00	P-classification	3.44	3.19	2.71	3.42	P-classification	2.16	2.82	2.16	1.98	P-classification	3.23	1.53	3.03	3.60
P-regression	1.76	4.74	4.71	3.00	P-regression	1.37	2.98	4.11	2.65	P-regression	2.68	1.85	3.61	3.45	P-regression	3.71	4.81	4.02	2.03
Clustering	3.10	3.32	4.03	3.00	Clustering	3.34	3.23	2.63	3.50	Clustering	2.18	4.58	3.56	3.85	Clustering	4.66	3.39	4.40	4.74
(i) LOIO-CV ( $\mathcal{A}_{\text{Is18}}$ )					(j) LOPO-CV ( $\mathcal{A}_{\text{Is18}}$ )					(k) LOPOAD-CV ( $\mathcal{A}_{\text{Is18}}$ )					(l) RI-CV ( $\mathcal{A}_{\text{Is18}}$ )				
$n$					$n$					$n$					$n$				
	2	3	5	10		2	3	5	10		2	3	5	10		2	3	5	10
Classification.	4.06	1.53	1.48	3.00	Classification.	3.31	2.08	3.29	3.53	Classification.	4.47	2.35	3.45	3.32	Classification.	1.60	1.11	1.19	2.13
Regression	2.47	3.74	2.94	3.00	Regression	3.47	2.81	2.31	1.82	Regression	3.56	3.56	2.32	2.40	Regression	2.39	3.21	2.37	2.34
P-classification	3.39	1.53	1.77	3.00	P-classification	2.69	3.61	2.73	3.45	P-classification	2.29	2.81	2.18	2.06	P-classification	2.85	2.40	3.06	3.47
P-regression	1.94	4.74	4.71	3.00	P-regression	1.74	3.13	4.10	2.66	P-regression	2.34	1.73	3.45	3.26	P-regression	3.42	4.60	3.95	2.42
Clustering	3.15	3.45	4.10	3.00	Clustering	3.79	3.37	2.58	3.53	Clustering	2.34	4.55	3.60	3.95	Clustering	4.74	3.68	4.42	4.65

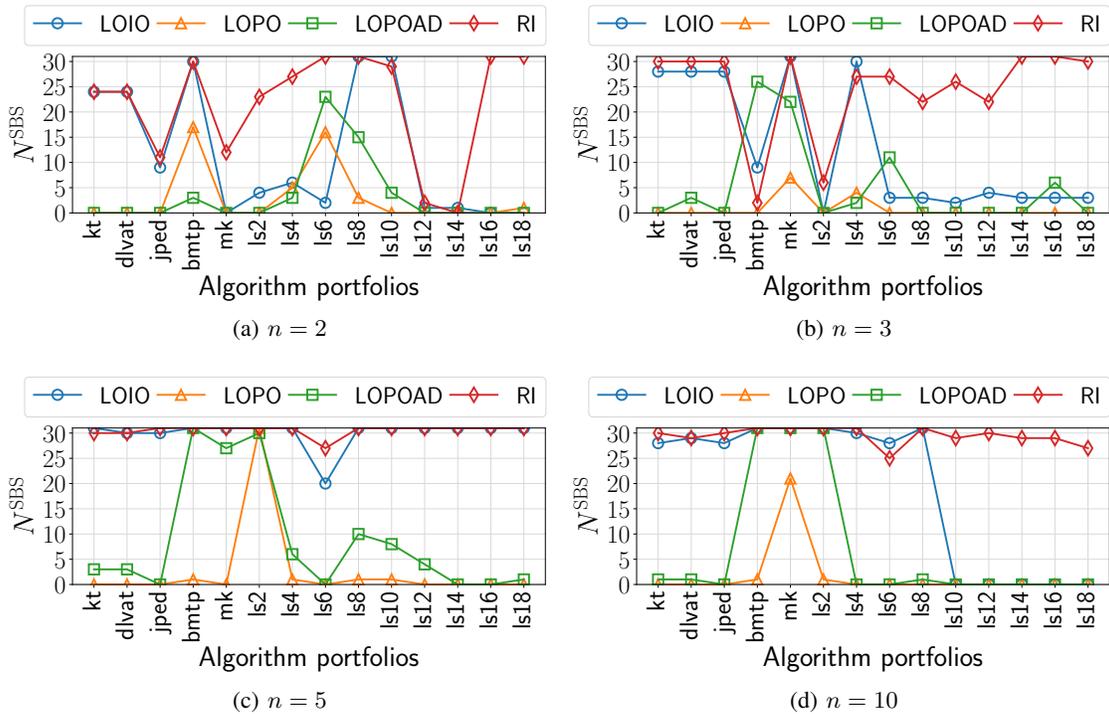


Fig. S.19: Number of times which the classification-based algorithm selection system outperforms the SBS in each algorithm portfolio.

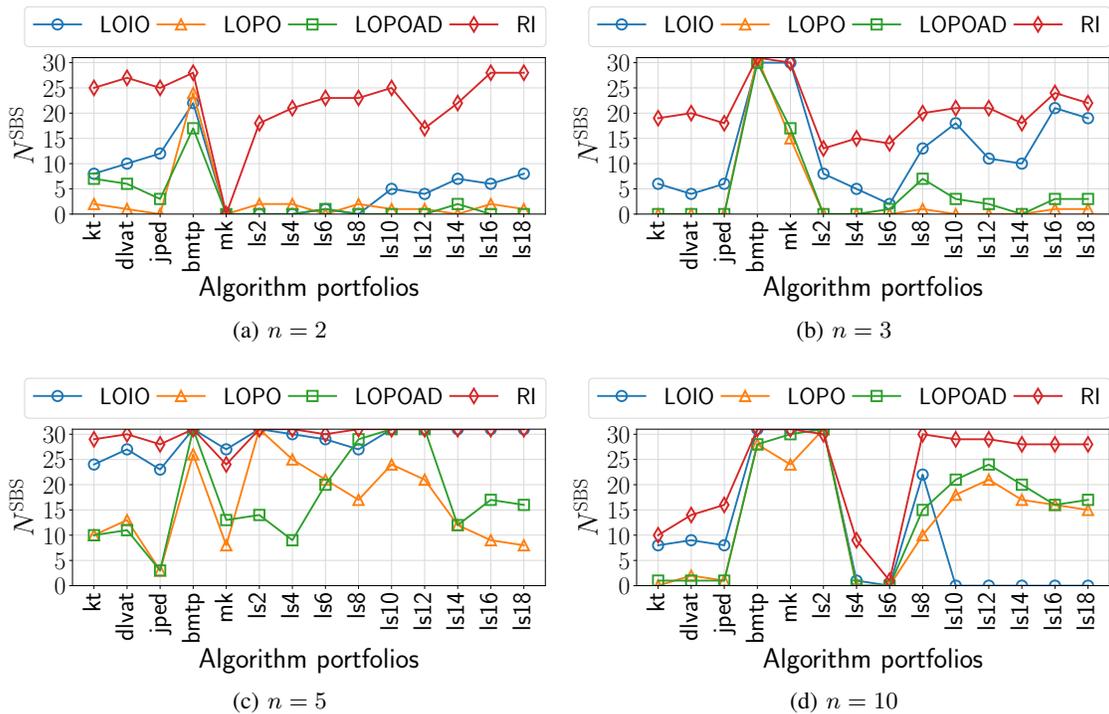


Fig. S.20: Number of times which the regression-based algorithm selection system outperforms the SBS in each algorithm portfolio.

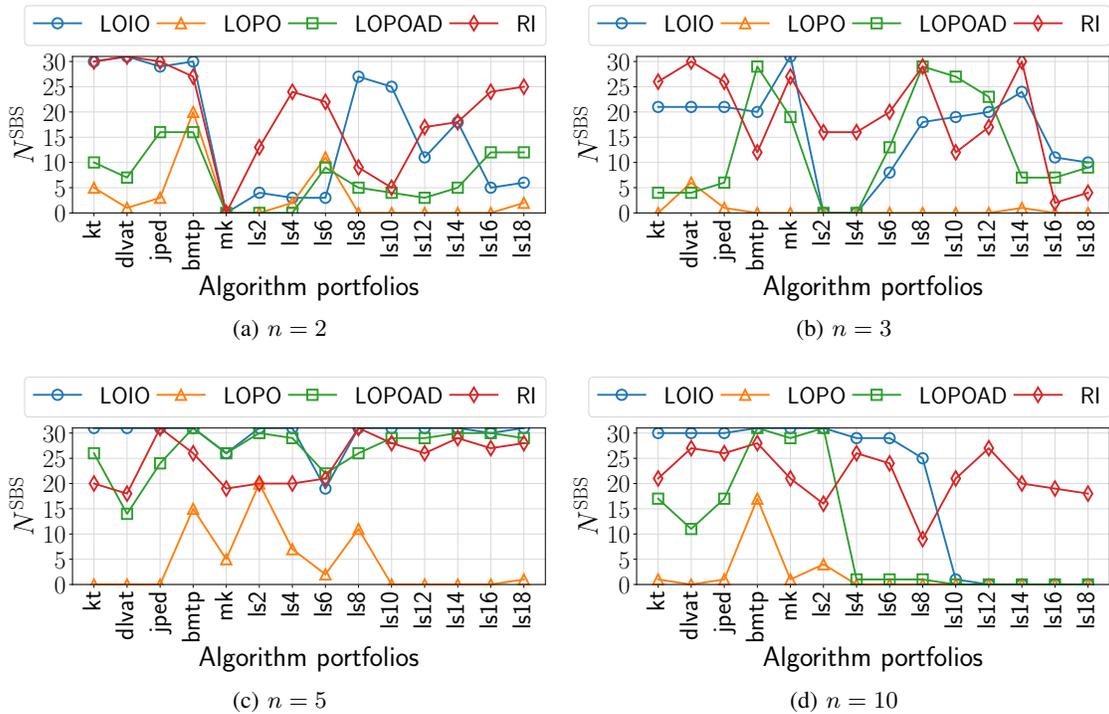


Fig. S.21: Number of times which the pairwise classification-based algorithm selection system outperforms the SBS in each algorithm portfolio.

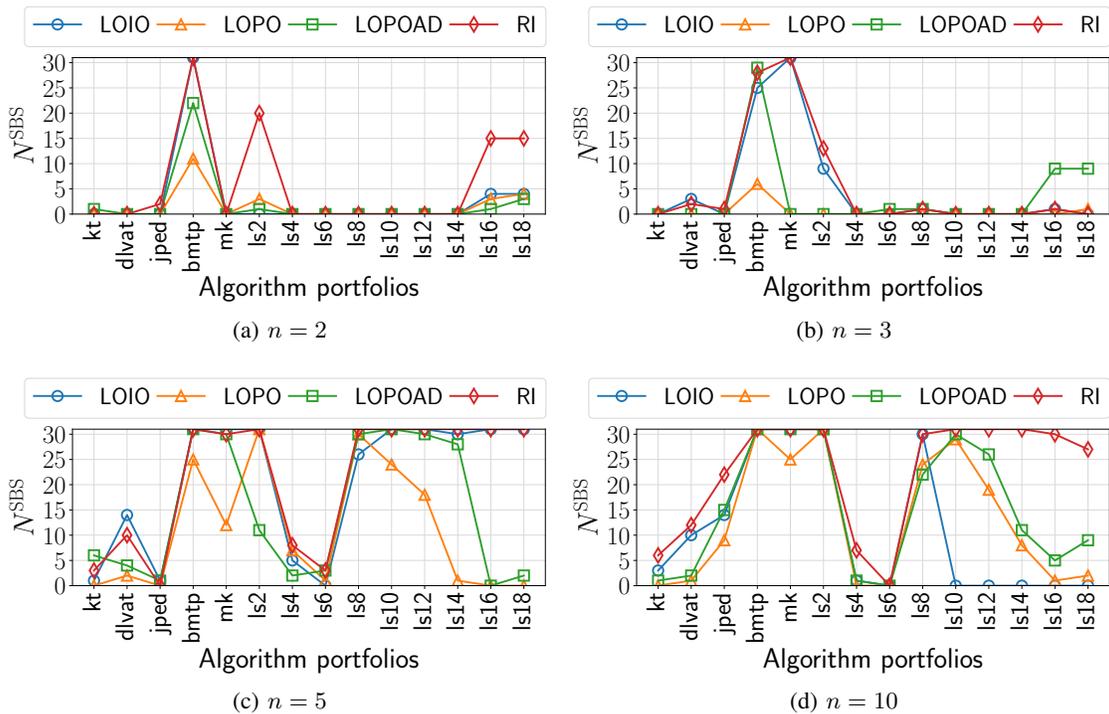


Fig. S.22: Number of times which the pairwise regression-based algorithm selection system outperforms the SBS in each algorithm portfolio.

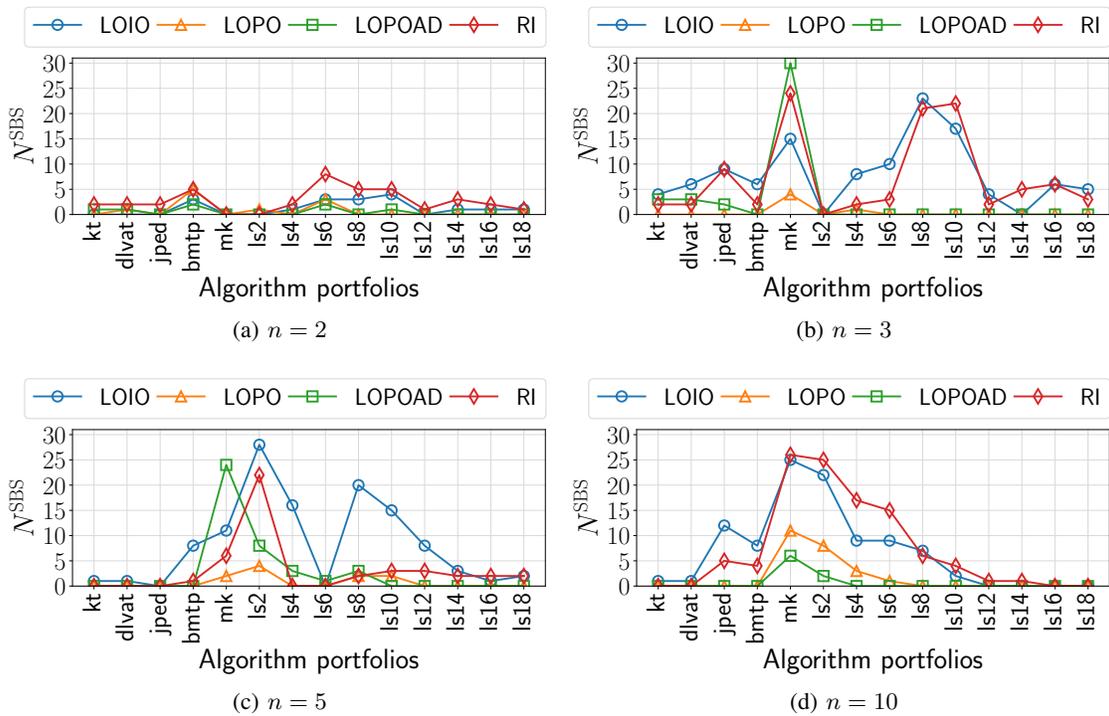


Fig. S.23: Number of times which the clustering-based algorithm selection system outperforms the SBS in each algorithm portfolio.

TABLE S.23: Results of the classification-based algorithm selection systems with the 14 algorithm portfolios for  $n \in \{2, 3, 5, 10\}$ . Tables (a)–(c) show the performance score values for the four cross-validation methods, respectively.

(a) LOIO-CV					(b) LOPO-CV					(c) LOPOAD-CV					(d) RI-CV				
	$n$																		
	2	3	5	10		2	3	5	10		2	3	5	10		2	3	5	10
$\mathcal{A}_{kt}$	0	2	0	0	$\mathcal{A}_{kt}$	5	12	5	1	$\mathcal{A}_{kt}$	4	6	3	6	$\mathcal{A}_{kt}$	6	2	4	5
$\mathcal{A}_{dlvat}$	0	1	1	0	$\mathcal{A}_{dlvat}$	4	4	5	1	$\mathcal{A}_{dlvat}$	3	0	3	3	$\mathcal{A}_{dlvat}$	6	3	4	7
$\mathcal{A}_{jped}$	3	2	0	0	$\mathcal{A}_{jped}$	9	4	8	1	$\mathcal{A}_{jped}$	4	6	10	7	$\mathcal{A}_{jped}$	12	2	3	5
$\mathcal{A}_{bmtmp}$	8	10	13	12	$\mathcal{A}_{bmtmp}$	5	6	7	7	$\mathcal{A}_{bmtmp}$	8	0	1	3	$\mathcal{A}_{bmtmp}$	13	13	12	11
$\mathcal{A}_{mk}$	6	4	11	13	$\mathcal{A}_{mk}$	3	2	3	2	$\mathcal{A}_{mk}$	6	0	3	1	$\mathcal{A}_{mk}$	7	11	13	12
$\mathcal{A}_{is2}$	9	5	3	4	$\mathcal{A}_{is2}$	2	0	0	0	$\mathcal{A}_{is2}$	2	3	0	0	$\mathcal{A}_{is2}$	4	11	10	8
$\mathcal{A}_{is4}$	10	0	4	3	$\mathcal{A}_{is4}$	1	0	1	1	$\mathcal{A}_{is4}$	1	1	2	2	$\mathcal{A}_{is4}$	3	7	9	8
$\mathcal{A}_{is6}$	11	9	7	5	$\mathcal{A}_{is6}$	0	2	2	6	$\mathcal{A}_{is6}$	0	1	4	3	$\mathcal{A}_{is6}$	2	7	8	8
$\mathcal{A}_{is8}$	4	9	6	5	$\mathcal{A}_{is8}$	1	4	3	8	$\mathcal{A}_{is8}$	1	8	3	7	$\mathcal{A}_{is8}$	0	6	3	8
$\mathcal{A}_{is10}$	4	9	10	10	$\mathcal{A}_{is10}$	7	4	3	7	$\mathcal{A}_{is10}$	9	11	6	8	$\mathcal{A}_{is10}$	0	6	3	1
$\mathcal{A}_{is12}$	7	9	11	7	$\mathcal{A}_{is12}$	12	13	10	12	$\mathcal{A}_{is12}$	6	13	7	9	$\mathcal{A}_{is12}$	8	7	2	0
$\mathcal{A}_{is14}$	0	6	3	7	$\mathcal{A}_{is14}$	11	4	10	7	$\mathcal{A}_{is14}$	10	8	10	9	$\mathcal{A}_{is14}$	5	0	2	0
$\mathcal{A}_{is16}$	12	6	6	7	$\mathcal{A}_{is16}$	8	4	10	8	$\mathcal{A}_{is16}$	11	6	10	12	$\mathcal{A}_{is16}$	5	0	0	0
$\mathcal{A}_{is18}$	13	6	6	10	$\mathcal{A}_{is18}$	5	6	10	12	$\mathcal{A}_{is18}$	10	10	10	12	$\mathcal{A}_{is18}$	6	2	0	0

TABLE S.24: Results of the regression-based algorithm selection systems with the 14 algorithm portfolios for  $n \in \{2, 3, 5, 10\}$ . Tables (a)–(c) show the performance score values for the four cross-validation methods, respectively.

(a) LOIO-CV					(b) LOPO-CV					(c) LOPOAD-CV					(d) RI-CV				
	$n$																		
	2	3	5	10		2	3	5	10		2	3	5	10		2	3	5	10
$\mathcal{A}_{kt}$	0	4	1	0	$\mathcal{A}_{kt}$	1	0	3	0	$\mathcal{A}_{kt}$	0	2	1	1	$\mathcal{A}_{kt}$	2	0	1	0
$\mathcal{A}_{dlvat}$	0	4	1	0	$\mathcal{A}_{dlvat}$	1	0	3	1	$\mathcal{A}_{dlvat}$	0	2	1	0	$\mathcal{A}_{dlvat}$	2	0	1	0
$\mathcal{A}_{jped}$	0	4	1	2	$\mathcal{A}_{jped}$	1	0	5	1	$\mathcal{A}_{jped}$	0	0	2	4	$\mathcal{A}_{jped}$	2	0	1	0
$\mathcal{A}_{bmtmp}$	7	13	11	6	$\mathcal{A}_{bmtmp}$	3	0	6	11	$\mathcal{A}_{bmtmp}$	5	1	12	12	$\mathcal{A}_{bmtmp}$	13	13	12	6
$\mathcal{A}_{mk}$	6	9	12	6	$\mathcal{A}_{mk}$	1	10	10	11	$\mathcal{A}_{mk}$	2	1	13	6	$\mathcal{A}_{mk}$	12	10	13	6
$\mathcal{A}_{is2}$	10	3	0	0	$\mathcal{A}_{is2}$	0	0	0	0	$\mathcal{A}_{is2}$	2	2	6	0	$\mathcal{A}_{is2}$	7	0	0	2
$\mathcal{A}_{is4}$	10	4	1	5	$\mathcal{A}_{is4}$	2	0	1	5	$\mathcal{A}_{is4}$	12	3	1	4	$\mathcal{A}_{is4}$	8	0	1	0
$\mathcal{A}_{is6}$	12	6	1	0	$\mathcal{A}_{is6}$	7	0	2	0	$\mathcal{A}_{is6}$	11	5	0	0	$\mathcal{A}_{is6}$	3	1	1	2
$\mathcal{A}_{is8}$	13	0	11	13	$\mathcal{A}_{is8}$	7	0	6	13	$\mathcal{A}_{is8}$	8	0	6	12	$\mathcal{A}_{is8}$	2	0	9	13
$\mathcal{A}_{is10}$	1	0	7	8	$\mathcal{A}_{is10}$	6	1	6	6	$\mathcal{A}_{is10}$	6	0	5	6	$\mathcal{A}_{is10}$	3	0	6	7
$\mathcal{A}_{is12}$	1	0	7	8	$\mathcal{A}_{is12}$	6	6	6	6	$\mathcal{A}_{is12}$	5	0	4	6	$\mathcal{A}_{is12}$	10	0	6	6
$\mathcal{A}_{is14}$	0	0	6	8	$\mathcal{A}_{is14}$	7	10	6	6	$\mathcal{A}_{is14}$	2	0	6	6	$\mathcal{A}_{is14}$	2	0	6	6
$\mathcal{A}_{is16}$	7	7	9	8	$\mathcal{A}_{is16}$	11	12	6	6	$\mathcal{A}_{is16}$	7	11	5	6	$\mathcal{A}_{is16}$	0	9	9	8
$\mathcal{A}_{is18}$	7	8	9	8	$\mathcal{A}_{is18}$	9	12	6	6	$\mathcal{A}_{is18}$	6	12	5	6	$\mathcal{A}_{is18}$	0	9	9	7

TABLE S.25: Results of the pairwise classification-based algorithm selection systems with the 14 algorithm portfolios for  $n \in \{2, 3, 5, 10\}$ . Tables (a)–(c) show the performance score values for the four cross-validation methods, respectively.

(a) LOIO-CV					(b) LOPO-CV					(c) LOPOAD-CV					(d) RI-CV				
	$n$																		
	2	3	5	10		2	3	5	10		2	3	5	10		2	3	5	10
$\mathcal{A}_{kt}$	0	0	0	0	$\mathcal{A}_{kt}$	0	5	1	3	$\mathcal{A}_{kt}$	0	1	0	1	$\mathcal{A}_{kt}$	0	0	1	1
$\mathcal{A}_{dlvat}$	0	0	0	0	$\mathcal{A}_{dlvat}$	1	0	5	3	$\mathcal{A}_{dlvat}$	1	1	6	1	$\mathcal{A}_{dlvat}$	2	1	2	1
$\mathcal{A}_{jped}$	0	0	0	0	$\mathcal{A}_{jped}$	0	3	1	3	$\mathcal{A}_{jped}$	0	1	0	1	$\mathcal{A}_{jped}$	0	0	0	0
$\mathcal{A}_{bmtmp}$	8	7	12	6	$\mathcal{A}_{bmtmp}$	8	11	5	6	$\mathcal{A}_{bmtmp}$	10	0	7	6	$\mathcal{A}_{bmtmp}$	13	11	11	7
$\mathcal{A}_{mk}$	6	6	12	6	$\mathcal{A}_{mk}$	1	4	9	7	$\mathcal{A}_{mk}$	8	5	11	7	$\mathcal{A}_{mk}$	11	9	13	12
$\mathcal{A}_{is2}$	9	5	3	4	$\mathcal{A}_{is2}$	0	0	0	0	$\mathcal{A}_{is2}$	2	5	0	0	$\mathcal{A}_{is2}$	6	5	5	5
$\mathcal{A}_{is4}$	9	7	4	0	$\mathcal{A}_{is4}$	0	1	0	1	$\mathcal{A}_{is4}$	3	8	0	1	$\mathcal{A}_{is4}$	3	5	2	1
$\mathcal{A}_{is6}$	11	9	5	4	$\mathcal{A}_{is6}$	0	5	1	1	$\mathcal{A}_{is6}$	0	7	0	5	$\mathcal{A}_{is6}$	3	6	2	1
$\mathcal{A}_{is8}$	4	4	11	6	$\mathcal{A}_{is8}$	5	5	5	8	$\mathcal{A}_{is8}$	2	1	6	13	$\mathcal{A}_{is8}$	9	4	5	12
$\mathcal{A}_{is10}$	4	4	6	8	$\mathcal{A}_{is10}$	10	10	12	7	$\mathcal{A}_{is10}$	2	1	5	8	$\mathcal{A}_{is10}$	11	9	6	7
$\mathcal{A}_{is12}$	6	4	6	8	$\mathcal{A}_{is12}$	10	8	12	7	$\mathcal{A}_{is12}$	2	1	5	8	$\mathcal{A}_{is12}$	4	6	11	5
$\mathcal{A}_{is14}$	0	0	6	8	$\mathcal{A}_{is14}$	5	9	6	7	$\mathcal{A}_{is14}$	1	7	5	8	$\mathcal{A}_{is14}$	2	0	5	6
$\mathcal{A}_{is16}$	12	11	6	8	$\mathcal{A}_{is16}$	10	10	10	7	$\mathcal{A}_{is16}$	10	12	7	8	$\mathcal{A}_{is16}$	3	10	5	6
$\mathcal{A}_{is18}$	12	11	6	8	$\mathcal{A}_{is18}$	9	10	6	7	$\mathcal{A}_{is18}$	10	11	7	8	$\mathcal{A}_{is18}$	3	11	5	6

TABLE S.26: Results of the pairwise regression-based algorithm selection systems with the 14 algorithm portfolios for  $n \in \{2, 3, 5, 10\}$ . Tables (a)–(c) show the performance score values for the four cross-validation methods, respectively.

(a) LOIO-CV					(b) LOPO-CV					(c) LOPOAD-CV					(d) RI-CV				
$n$																			
	2	3	5	10		2	3	5	10		2	3	5	10		2	3	5	10
$\mathcal{A}_{kt}$	0	3	2	1	$\mathcal{A}_{kt}$	2	1	4	3	$\mathcal{A}_{kt}$	4	2	0	2	$\mathcal{A}_{kt}$	1	3	1	2
$\mathcal{A}_{dlvat}$	9	0	1	1	$\mathcal{A}_{dlvat}$	2	0	1	4	$\mathcal{A}_{dlvat}$	9	1	3	2	$\mathcal{A}_{dlvat}$	9	0	1	0
$\mathcal{A}_{jped}$	0	0	4	0	$\mathcal{A}_{jped}$	2	0	4	1	$\mathcal{A}_{jped}$	2	0	3	0	$\mathcal{A}_{jped}$	1	0	4	0
$\mathcal{A}_{bmtmp}$	3	5	9	7	$\mathcal{A}_{bmtmp}$	7	11	6	8	$\mathcal{A}_{bmtmp}$	2	1	3	10	$\mathcal{A}_{bmtmp}$	8	5	8	7
$\mathcal{A}_{mk}$	0	3	9	1	$\mathcal{A}_{mk}$	1	10	8	8	$\mathcal{A}_{mk}$	1	13	11	8	$\mathcal{A}_{mk}$	1	4	11	2
$\mathcal{A}_{is2}$	12	1	0	0	$\mathcal{A}_{is2}$	0	2	0	0	$\mathcal{A}_{is2}$	0	2	3	0	$\mathcal{A}_{is2}$	0	0	0	1
$\mathcal{A}_{is4}$	10	10	1	3	$\mathcal{A}_{is4}$	10	6	1	2	$\mathcal{A}_{is4}$	12	12	0	5	$\mathcal{A}_{is4}$	13	10	1	1
$\mathcal{A}_{is6}$	10	10	5	5	$\mathcal{A}_{is6}$	8	4	3	1	$\mathcal{A}_{is6}$	10	1	1	1	$\mathcal{A}_{is6}$	11	5	4	3
$\mathcal{A}_{is8}$	7	5	13	8	$\mathcal{A}_{is8}$	9	5	4	10	$\mathcal{A}_{is8}$	5	1	0	13	$\mathcal{A}_{is8}$	9	5	5	13
$\mathcal{A}_{is10}$	1	5	7	9	$\mathcal{A}_{is10}$	2	3	4	2	$\mathcal{A}_{is10}$	2	1	4	6	$\mathcal{A}_{is10}$	4	5	8	10
$\mathcal{A}_{is12}$	0	7	5	9	$\mathcal{A}_{is12}$	2	3	6	3	$\mathcal{A}_{is12}$	4	1	3	6	$\mathcal{A}_{is12}$	3	5	6	8
$\mathcal{A}_{is14}$	1	5	5	9	$\mathcal{A}_{is14}$	2	5	9	8	$\mathcal{A}_{is14}$	4	1	5	9	$\mathcal{A}_{is14}$	4	7	6	7
$\mathcal{A}_{is16}$	0	12	10	13	$\mathcal{A}_{is16}$	2	10	12	12	$\mathcal{A}_{is16}$	5	5	13	10	$\mathcal{A}_{is16}$	1	12	11	10
$\mathcal{A}_{is18}$	1	12	10	9	$\mathcal{A}_{is18}$	2	12	12	12	$\mathcal{A}_{is18}$	3	4	12	8	$\mathcal{A}_{is18}$	1	12	11	7

TABLE S.27: Results of the clustering-based algorithm selection systems with the 14 algorithm portfolios for  $n \in \{2, 3, 5, 10\}$ . Tables (a)–(c) show the performance score values for the four cross-validation methods, respectively.

(a) LOIO-CV					(b) LOPO-CV					(c) LOPOAD-CV					(d) RI-CV				
$n$																			
	2	3	5	10		2	3	5	10		2	3	5	10		2	3	5	10
$\mathcal{A}_{kt}$	5	0	3	4	$\mathcal{A}_{kt}$	0	6	3	4	$\mathcal{A}_{kt}$	0	0	4	4	$\mathcal{A}_{kt}$	1	4	2	5
$\mathcal{A}_{dlvat}$	4	0	3	4	$\mathcal{A}_{dlvat}$	0	4	3	5	$\mathcal{A}_{dlvat}$	0	0	4	4	$\mathcal{A}_{dlvat}$	2	4	2	5
$\mathcal{A}_{jped}$	3	0	3	3	$\mathcal{A}_{jped}$	0	4	3	4	$\mathcal{A}_{jped}$	0	0	4	4	$\mathcal{A}_{jped}$	0	0	2	3
$\mathcal{A}_{bmtmp}$	3	13	13	7	$\mathcal{A}_{bmtmp}$	10	13	13	8	$\mathcal{A}_{bmtmp}$	13	12	10	7	$\mathcal{A}_{bmtmp}$	10	13	13	10
$\mathcal{A}_{mk}$	2	10	6	3	$\mathcal{A}_{mk}$	5	2	5	3	$\mathcal{A}_{mk}$	2	0	3	3	$\mathcal{A}_{mk}$	7	8	7	3
$\mathcal{A}_{is2}$	0	5	0	0	$\mathcal{A}_{is2}$	0	0	0	0	$\mathcal{A}_{is2}$	4	1	0	0	$\mathcal{A}_{is2}$	0	4	0	0
$\mathcal{A}_{is4}$	0	6	0	1	$\mathcal{A}_{is4}$	1	0	1	1	$\mathcal{A}_{is4}$	5	2	0	1	$\mathcal{A}_{is4}$	0	5	1	1
$\mathcal{A}_{is6}$	3	6	2	2	$\mathcal{A}_{is6}$	6	3	2	1	$\mathcal{A}_{is6}$	5	0	2	1	$\mathcal{A}_{is6}$	1	9	2	1
$\mathcal{A}_{is8}$	3	0	3	5	$\mathcal{A}_{is8}$	5	2	3	5	$\mathcal{A}_{is8}$	3	7	7	7	$\mathcal{A}_{is8}$	2	0	7	5
$\mathcal{A}_{is10}$	3	0	3	3	$\mathcal{A}_{is10}$	5	2	3	5	$\mathcal{A}_{is10}$	4	7	7	7	$\mathcal{A}_{is10}$	2	0	7	4
$\mathcal{A}_{is12}$	2	0	10	8	$\mathcal{A}_{is12}$	5	3	3	9	$\mathcal{A}_{is12}$	3	7	7	7	$\mathcal{A}_{is12}$	2	8	8	8
$\mathcal{A}_{is14}$	3	0	3	6	$\mathcal{A}_{is14}$	5	10	10	5	$\mathcal{A}_{is14}$	10	10	10	9	$\mathcal{A}_{is14}$	2	0	6	6
$\mathcal{A}_{is16}$	11	10	4	11	$\mathcal{A}_{is16}$	12	11	3	12	$\mathcal{A}_{is16}$	11	11	12	12	$\mathcal{A}_{is16}$	11	9	6	12
$\mathcal{A}_{is18}$	11	10	5	11	$\mathcal{A}_{is18}$	12	11	3	12	$\mathcal{A}_{is18}$	10	11	12	12	$\mathcal{A}_{is18}$	11	10	6	12

TABLE S.28: Results of the classification-based algorithm selection systems with the 14 algorithm portfolios for  $n \in \{2, 3, 5, 10\}$ . Tables (a)–(c) show the Friedman test-based average rankings for the four cross-validation methods, respectively.

(a) LOIO-CV					(b) LOPO-CV					(c) LOPOAD-CV					(d) RI-CV				
$n$					$n$					$n$					$n$				
	2	3	5	10		2	3	5	10		2	3	5	10		2	3	5	10
$\mathcal{A}_{kt}$	2.39	3.97	2.10	2.58	$\mathcal{A}_{kt}$	8.52	12.427.74	4.02		$\mathcal{A}_{kt}$	6.39	9.48	7.90	6.58	$\mathcal{A}_{kt}$	9.52	4.97	6.90	6.31
$\mathcal{A}_{divat}$	2.58	3.00	3.13	2.39	$\mathcal{A}_{divat}$	7.00	7.32	7.58	4.10	$\mathcal{A}_{divat}$	5.84	4.06	6.13	5.45	$\mathcal{A}_{divat}$	8.81	5.39	6.97	7.29
$\mathcal{A}_{jped}$	4.19	3.77	3.06	2.68	$\mathcal{A}_{jped}$	10.427.87	9.68	4.02		$\mathcal{A}_{jped}$	7.00	7.74	11.528.81		$\mathcal{A}_{jped}$	11.774.39	5.52	6.31	
$\mathcal{A}_{bmtmp}$	9.32	11.1913.2312.74			$\mathcal{A}_{bmtmp}$	7.39	9.06	9.06	9.61	$\mathcal{A}_{bmtmp}$	8.77	2.71	2.39	5.42	$\mathcal{A}_{bmtmp}$	13.4513.8412.9711.81			
$\mathcal{A}_{mk}$	7.29	5.90	12.2613.10		$\mathcal{A}_{mk}$	5.23	3.94	5.48	4.81	$\mathcal{A}_{mk}$	7.48	3.65	6.32	2.65	$\mathcal{A}_{mk}$	9.39	11.6113.4812.48		
$\mathcal{A}_{is2}$	9.06	6.87	4.74	4.97	$\mathcal{A}_{is2}$	5.45	1.71	1.00	1.23	$\mathcal{A}_{is2}$	5.71	4.90	1.06	1.00	$\mathcal{A}_{is2}$	5.58	11.5510.529.68		
$\mathcal{A}_{is4}$	9.03	2.23	6.16	4.10	$\mathcal{A}_{is4}$	3.16	1.58	2.23	4.03	$\mathcal{A}_{is4}$	3.55	4.06	3.29	2.97	$\mathcal{A}_{is4}$	4.87	9.68	9.39	9.10
$\mathcal{A}_{is6}$	11.4510.818.81	6.39			$\mathcal{A}_{is6}$	1.61	3.94	3.26	7.23	$\mathcal{A}_{is6}$	1.84	4.65	6.65	5.77	$\mathcal{A}_{is6}$	3.13	9.23	9.29	9.55
$\mathcal{A}_{is8}$	6.03	11.397.87	7.00		$\mathcal{A}_{is8}$	3.90	7.90	7.06	10.13	$\mathcal{A}_{is8}$	3.61	9.68	7.13	9.87	$\mathcal{A}_{is8}$	1.97	8.68	6.55	10.55
$\mathcal{A}_{is10}$	5.68	11.2911.1910.71			$\mathcal{A}_{is10}$	9.10	8.29	5.52	10.65	$\mathcal{A}_{is10}$	8.94	10.978.45	10.13		$\mathcal{A}_{is10}$	2.13	7.39	6.68	4.71
$\mathcal{A}_{is12}$	8.29	9.71	12.358.81		$\mathcal{A}_{is12}$	13.1013.2312.2911.90				$\mathcal{A}_{is12}$	9.77	13.879.68	10.94		$\mathcal{A}_{is12}$	10.459.35	6.35	4.00	
$\mathcal{A}_{is14}$	2.81	8.26	5.06	9.29	$\mathcal{A}_{is14}$	11.009.18	11.109.77			$\mathcal{A}_{is14}$	11.9010.3511.7710.84				$\mathcal{A}_{is14}$	8.23	2.61	4.77	3.74
$\mathcal{A}_{is16}$	13.168.29	8.00	9.13		$\mathcal{A}_{is16}$	10.299.18	11.6110.65			$\mathcal{A}_{is16}$	12.587.65	11.6111.68			$\mathcal{A}_{is16}$	6.97	2.58	3.00	3.87
$\mathcal{A}_{is18}$	13.718.32	7.03	11.13		$\mathcal{A}_{is18}$	8.84	9.39	11.3912.87		$\mathcal{A}_{is18}$	11.6111.2311.1012.90				$\mathcal{A}_{is18}$	8.74	3.74	2.61	5.61

TABLE S.29: Results of the regression-based algorithm selection systems with the 14 algorithm portfolios for  $n \in \{2, 3, 5, 10\}$ . Tables (a)–(c) show the Friedman test-based average rankings for the four cross-validation methods, respectively.

(a) LOIO-CV					(b) LOPO-CV					(c) LOPOAD-CV					(d) RI-CV				
$n$					$n$					$n$					$n$				
	2	3	5	10		2	3	5	10		2	3	5	10		2	3	5	10
$\mathcal{A}_{kt}$	5.29	6.66	3.84	3.11	$\mathcal{A}_{kt}$	5.19	5.81	4.52	2.60	$\mathcal{A}_{kt}$	4.55	7.55	4.81	3.18	$\mathcal{A}_{kt}$	6.29	6.03	4.06	3.11
$\mathcal{A}_{divat}$	4.87	7.13	3.97	2.98	$\mathcal{A}_{divat}$	5.65	5.56	4.03	3.31	$\mathcal{A}_{divat}$	4.48	7.55	4.19	2.85	$\mathcal{A}_{divat}$	6.48	6.23	3.52	3.52
$\mathcal{A}_{jped}$	4.26	6.69	4.81	3.90	$\mathcal{A}_{jped}$	6.77	5.95	5.55	3.73	$\mathcal{A}_{jped}$	5.71	4.90	6.48	5.16	$\mathcal{A}_{jped}$	6.03	5.87	3.74	3.71
$\mathcal{A}_{bmtmp}$	8.10	12.8712.137.13			$\mathcal{A}_{bmtmp}$	7.00	6.19	10.2611.26		$\mathcal{A}_{bmtmp}$	7.10	5.45	11.7113.00		$\mathcal{A}_{bmtmp}$	13.3913.1912.398.61			
$\mathcal{A}_{mk}$	6.68	10.5212.777.58			$\mathcal{A}_{mk}$	4.90	9.81	11.6111.68		$\mathcal{A}_{mk}$	4.97	7.81	13.559.71		$\mathcal{A}_{mk}$	11.8410.6513.719.55			
$\mathcal{A}_{is2}$	10.426.10	1.68	3.42		$\mathcal{A}_{is2}$	2.84	5.10	1.42	2.48	$\mathcal{A}_{is2}$	7.03	7.06	4.06	2.23	$\mathcal{A}_{is2}$	8.52	6.32	1.35	4.69
$\mathcal{A}_{is4}$	10.818.00	3.77	4.90		$\mathcal{A}_{is4}$	5.90	6.48	2.29	7.03	$\mathcal{A}_{is4}$	12.718.00	5.13	5.89		$\mathcal{A}_{is4}$	9.16	7.16	4.23	4.74
$\mathcal{A}_{is6}$	11.819.16	3.03	3.58		$\mathcal{A}_{is6}$	9.19	6.77	3.48	3.24	$\mathcal{A}_{is6}$	11.659.00	2.32	2.63		$\mathcal{A}_{is6}$	7.48	6.42	4.32	4.39
$\mathcal{A}_{is8}$	12.294.90	12.9412.45			$\mathcal{A}_{is8}$	8.71	4.97	10.2613.58		$\mathcal{A}_{is8}$	9.10	5.55	9.35	13.16	$\mathcal{A}_{is8}$	6.58	6.35	9.90	14.00
$\mathcal{A}_{is10}$	5.00	4.45	9.10	11.13	$\mathcal{A}_{is10}$	8.97	6.32	9.77	8.74	$\mathcal{A}_{is10}$	8.03	6.94	8.32	9.58	$\mathcal{A}_{is10}$	7.29	5.06	8.94	9.90
$\mathcal{A}_{is12}$	5.16	4.26	8.94	10.71	$\mathcal{A}_{is12}$	8.45	7.29	10.718.48		$\mathcal{A}_{is12}$	7.19	5.90	8.39	9.13	$\mathcal{A}_{is12}$	9.35	5.90	8.42	9.65
$\mathcal{A}_{is14}$	3.71	4.13	7.39	11.00	$\mathcal{A}_{is14}$	9.55	9.71	9.68	9.10	$\mathcal{A}_{is14}$	5.19	6.94	9.32	9.68	$\mathcal{A}_{is14}$	5.74	5.26	7.77	8.55
$\mathcal{A}_{is16}$	8.32	10.1310.1911.76			$\mathcal{A}_{is16}$	11.0012.5211.039.69				$\mathcal{A}_{is16}$	8.74	10.778.97	9.71		$\mathcal{A}_{is16}$	3.48	10.3211.4210.42		
$\mathcal{A}_{is18}$	8.29	10.0010.4511.34			$\mathcal{A}_{is18}$	10.8712.5210.3910.08				$\mathcal{A}_{is18}$	8.55	11.588.39	9.10		$\mathcal{A}_{is18}$	3.35	10.2311.2310.16		

TABLE S.30: Results of the pairwise classification-based algorithm selection systems with the 14 algorithm portfolios for  $n \in \{2, 3, 5, 10\}$ . Tables (a)–(c) show the Friedman test-based average rankings for the four cross-validation methods, respectively.

(a) LOIO-CV					(b) LOPO-CV					(c) LOPOAD-CV					(d) RI-CV				
$n$					$n$					$n$					$n$				
	2	3	5	10		2	3	5	10		2	3	5	10		2	3	5	10
$\mathcal{A}_{kt}$	3.10	3.95	2.50	2.19	$\mathcal{A}_{kt}$	4.06	7.40	3.23	4.82	$\mathcal{A}_{kt}$	6.58	6.65	4.10	3.35	$\mathcal{A}_{kt}$	2.68	3.44	3.39	4.58
$\mathcal{A}_{divat}$	2.23	4.26	2.63	3.61	$\mathcal{A}_{divat}$	6.52	2.68	7.74	4.68	$\mathcal{A}_{divat}$	6.74	6.94	7.13	4.39	$\mathcal{A}_{divat}$	4.48	3.39	4.90	3.74
$\mathcal{A}_{jped}$	3.42	3.89	1.63	2.42	$\mathcal{A}_{jped}$	5.13	4.08	3.58	4.18	$\mathcal{A}_{jped}$	4.23	5.65	5.00	3.16	$\mathcal{A}_{jped}$	2.84	3.15	1.32	2.71
$\mathcal{A}_{bmtmp}$	9.19	10.0312.908.77			$\mathcal{A}_{bmtmp}$	9.65	12.007.68	6.94		$\mathcal{A}_{bmtmp}$	11.3	3.94	10.297.84		$\mathcal{A}_{bmtmp}$	13.0612.0011.489.61			
$\mathcal{A}_{mk}$	7.74	8.13	13.269.10		$\mathcal{A}_{mk}$	6.06	5.81	10.2310.94		$\mathcal{A}_{mk}$	10.068.35	11.658.23			$\mathcal{A}_{mk}$	11.559.84	13.8112.77		
$\mathcal{A}_{is2}$	9.90	7.26	3.87	5.29	$\mathcal{A}_{is2}$	4.42	1.68	2.39	1.84	$\mathcal{A}_{is2}$	7.74	8.45	3.03	1.77	$\mathcal{A}_{is2}$	8.26	7.10	8.48	7.26
$\mathcal{A}_{is4}$	10.2310.164.58	2.68			$\mathcal{A}_{is4}$	5.32	2.77	3.26	2.90	$\mathcal{A}_{is4}$	8.39	9.48	3.00	3.84	$\mathcal{A}_{is4}$	6.26	7.77	4.77	4.16
$\mathcal{A}_{is6}$	11.6810.486.74	5.10			$\mathcal{A}_{is6}$	4.29	6.81	4.19	3.29	$\mathcal{A}_{is6}$	4.42	9.06	4.94	5.32	$\mathcal{A}_{is6}$	7.45	8.06	5.10	4.35
$\mathcal{A}_{is8}$	6.55	7.32	10.6810.29		$\mathcal{A}_{is8}$	7.58	7.81	7.94	11.58	$\mathcal{A}_{is8}$	6.21	5.16	9.55	12.74	$\mathcal{A}_{is8}$	9.52	5.48	8.03	12.90
$\mathcal{A}_{is10}$	6.39	6.94	9.35	11.45	$\mathcal{A}_{is10}$	11.8711.9712.529.97				$\mathcal{A}_{is10}$	6.66	5.97	8.71	10.77	$\mathcal{A}_{is10}$	10.3910.168.77	9.10		
$\mathcal{A}_{is12}$	7.58	7.90	10.1310.58		$\mathcal{A}_{is12}$	10.488.77	13.0311.10			$\mathcal{A}_{is12}$	6.61	5.84	9.03	10.77	$\mathcal{A}_{is12}$	8.13	8.45	12.267.81	
$\mathcal{A}_{is14}$	3.29	3.68	8.45	11.48	$\mathcal{A}_{is14}$	7.39	9.74	9.10	11.32	$\mathcal{A}_{is14}$	6.26	8.65	7.94	11.13	$\mathcal{A}_{is14}$	6.45	2.71	6.55	8.68
$\mathcal{A}_{is16}$	12.1910.169.29	11.13			$\mathcal{A}_{is16}$	12.0611.8111.2910.48				$\mathcal{A}_{is16}$	9.53	11.0210.3210.68			$\mathcal{A}_{is16}$	6.81	11.618.00	8.65	
$\mathcal{A}_{is18}$	11.5210.848.97	10.90			$\mathcal{A}_{is18}$	10.1611.688.84	10.97			$\mathcal{A}_{is18}$	10.219.85	10.3211.00			$\mathcal{A}_{is18}$	7.13	11.848.13	8.68	

TABLE S.31: Results of the pairwise regression-based algorithm selection systems with the 14 algorithm portfolios for  $n \in \{2, 3, 5, 10\}$ . Tables (a)–(c) show the Friedman test-based average rankings for the four cross-validation methods, respectively.

	(a) LOIO-CV				(b) LOPO-CV				(c) LOPOAD-CV				(d) RI-CV						
	$n$				$n$				$n$				$n$						
	2	3	5	10	2	3	5	10	2	3	5	10	2	3	5	10			
$\mathcal{A}_{kt}$	6.74	4.68	3.39	4.26	$\mathcal{A}_{kt}$	8.84	6.03	7.68	6.32	$\mathcal{A}_{kt}$	7.71	8.42	4.35	4.65	$\mathcal{A}_{kt}$	6.06	4.61	3.42	4.74
$\mathcal{A}_{divat}$	9.74	2.55	3.00	3.52	$\mathcal{A}_{divat}$	8.65	3.00	2.84	6.42	$\mathcal{A}_{divat}$	10.167	0.0	6.29	4.45	$\mathcal{A}_{divat}$	9.68	3.03	2.94	4.03
$\mathcal{A}_{jped}$	4.90	1.97	5.26	3.97	$\mathcal{A}_{jped}$	7.32	3.00	6.35	3.81	$\mathcal{A}_{jped}$	5.06	2.81	6.29	3.00	$\mathcal{A}_{jped}$	7.48	2.61	5.42	3.32
$\mathcal{A}_{bmtp}$	7.03	7.68	10.657	7.7	$\mathcal{A}_{bmtp}$	9.10	10.778	8.1	10.10	$\mathcal{A}_{bmtp}$	6.29	6.00	6.77	11.90	$\mathcal{A}_{bmtp}$	9.10	7.68	9.81	9.10
$\mathcal{A}_{mk}$	4.06	5.35	11.104	1.0	$\mathcal{A}_{mk}$	3.74	11.061	0.521	0.48	$\mathcal{A}_{mk}$	3.68	13.771	1.749	3.2	$\mathcal{A}_{mk}$	5.48	5.90	12.654	7.1
$\mathcal{A}_{Is2}$	11.942	7.4	1.06	2.48	$\mathcal{A}_{Is2}$	1.06	4.52	1.16	1.32	$\mathcal{A}_{Is2}$	3.74	7.06	7.29	1.87	$\mathcal{A}_{Is2}$	1.23	2.42	1.10	2.97
$\mathcal{A}_{Is4}$	11.611	1.033	3.2	4.90	$\mathcal{A}_{Is4}$	11.138	6.1	2.48	5.13	$\mathcal{A}_{Is4}$	12.031	0.2	3.97	6.35	$\mathcal{A}_{Is4}$	12.771	0.873	5.5	4.19
$\mathcal{A}_{Is6}$	10.321	0.587	4.5	5.48	$\mathcal{A}_{Is6}$	9.87	7.94	4.52	3.52	$\mathcal{A}_{Is6}$	10.587	3.5	5.45	3.52	$\mathcal{A}_{Is6}$	11.458	6.5	6.94	5.52
$\mathcal{A}_{Is8}$	8.39	8.19	13.529	9.4	$\mathcal{A}_{Is8}$	10.457	8.4	6.94	11.16	$\mathcal{A}_{Is8}$	8.00	6.74	3.94	13.32	$\mathcal{A}_{Is8}$	9.29	8.35	7.45	13.71
$\mathcal{A}_{Is10}$	5.32	7.58	8.74	11.94	$\mathcal{A}_{Is10}$	7.23	5.58	8.26	6.29	$\mathcal{A}_{Is10}$	5.74	6.68	7.39	8.29	$\mathcal{A}_{Is10}$	7.42	7.45	9.68	11.90
$\mathcal{A}_{Is12}$	5.26	9.26	7.32	11.13	$\mathcal{A}_{Is12}$	7.13	6.10	8.84	6.32	$\mathcal{A}_{Is12}$	7.55	6.55	7.52	7.39	$\mathcal{A}_{Is12}$	7.13	8.29	8.55	9.39
$\mathcal{A}_{Is14}$	6.74	7.55	7.45	11.03	$\mathcal{A}_{Is14}$	6.77	7.16	10.009	0.0	$\mathcal{A}_{Is14}$	7.23	6.00	7.55	10.71	$\mathcal{A}_{Is14}$	6.84	8.97	8.23	9.32
$\mathcal{A}_{Is16}$	5.58	12.901	1.481	2.97	$\mathcal{A}_{Is16}$	6.48	11.681	3.291	2.94	$\mathcal{A}_{Is16}$	8.97	8.42	13.741	1.65	$\mathcal{A}_{Is16}$	5.71	13.101	2.551	1.74
$\mathcal{A}_{Is18}$	7.35	12.941	1.261	1.52	$\mathcal{A}_{Is18}$	7.23	11.711	3.321	2.19	$\mathcal{A}_{Is18}$	8.26	7.90	12.718	5.8	$\mathcal{A}_{Is18}$	5.35	13.061	2.741	0.35

TABLE S.32: Results of the clustering-based algorithm selection systems with the 14 algorithm portfolios for  $n \in \{2, 3, 5, 10\}$ . Tables (a)–(c) show the Friedman test-based average rankings for the four cross-validation methods, respectively.

	(a) LOIO-CV				(b) LOPO-CV				(c) LOPOAD-CV				(d) RI-CV						
	$n$				$n$				$n$				$n$						
	2	3	5	10	2	3	5	10	2	3	5	10	2	3	5	10			
$\mathcal{A}_{kt}$	8.32	5.97	7.63	8.26	$\mathcal{A}_{kt}$	4.82	8.42	8.45	7.34	$\mathcal{A}_{kt}$	3.65	3.47	6.65	6.73	$\mathcal{A}_{kt}$	5.61	7.26	5.40	9.13
$\mathcal{A}_{divat}$	8.58	5.21	7.60	8.00	$\mathcal{A}_{divat}$	4.61	7.48	8.84	7.63	$\mathcal{A}_{divat}$	3.48	3.76	6.77	6.92	$\mathcal{A}_{divat}$	6.42	6.61	5.34	9.29
$\mathcal{A}_{jped}$	8.48	4.34	7.26	5.23	$\mathcal{A}_{jped}$	3.85	8.42	8.23	6.00	$\mathcal{A}_{jped}$	4.77	4.55	6.26	6.35	$\mathcal{A}_{jped}$	5.29	4.35	4.52	5.39
$\mathcal{A}_{bmtp}$	9.29	12.841	1.551	0.35	$\mathcal{A}_{bmtp}$	11.231	3.161	3.459	7.4	$\mathcal{A}_{bmtp}$	12.841	3.611	1.659	7.7	$\mathcal{A}_{bmtp}$	10.871	3.161	3.101	0.42
$\mathcal{A}_{mk}$	5.19	11.451	0.196	3.2	$\mathcal{A}_{mk}$	7.10	4.87	8.74	4.26	$\mathcal{A}_{mk}$	4.61	4.06	4.65	4.26	$\mathcal{A}_{mk}$	8.90	8.55	10.355	9.4
$\mathcal{A}_{Is2}$	3.66	6.77	1.60	1.69	$\mathcal{A}_{Is2}$	3.65	1.71	1.03	1.19	$\mathcal{A}_{Is2}$	6.00	4.55	1.65	1.26	$\mathcal{A}_{Is2}$	4.92	6.06	1.39	1.31
$\mathcal{A}_{Is4}$	3.66	6.55	1.82	2.32	$\mathcal{A}_{Is4}$	5.29	2.03	2.90	2.90	$\mathcal{A}_{Is4}$	7.45	5.16	1.87	2.13	$\mathcal{A}_{Is4}$	4.92	7.39	3.13	2.47
$\mathcal{A}_{Is6}$	7.90	8.16	3.68	3.27	$\mathcal{A}_{Is6}$	9.00	6.00	4.10	3.39	$\mathcal{A}_{Is6}$	6.84	5.35	3.19	3.16	$\mathcal{A}_{Is6}$	6.55	9.97	5.23	2.61
$\mathcal{A}_{Is8}$	6.65	4.73	7.53	8.71	$\mathcal{A}_{Is8}$	7.10	5.47	7.73	8.98	$\mathcal{A}_{Is8}$	7.47	9.02	7.98	8.87	$\mathcal{A}_{Is8}$	7.37	3.56	10.298	3.4
$\mathcal{A}_{Is10}$	7.00	4.44	8.47	6.84	$\mathcal{A}_{Is10}$	7.42	5.89	7.95	7.61	$\mathcal{A}_{Is10}$	7.85	9.24	8.15	8.58	$\mathcal{A}_{Is10}$	6.85	3.34	10.616	4.7
$\mathcal{A}_{Is12}$	7.56	4.94	11.321	0.03	$\mathcal{A}_{Is12}$	7.44	7.16	7.42	10.31	$\mathcal{A}_{Is12}$	6.97	7.32	9.03	9.84	$\mathcal{A}_{Is12}$	6.16	9.32	9.26	9.74
$\mathcal{A}_{Is14}$	7.08	6.39	7.61	8.97	$\mathcal{A}_{Is14}$	8.02	9.90	10.429	3.2	$\mathcal{A}_{Is14}$	9.32	9.90	10.481	0.58	$\mathcal{A}_{Is14}$	6.10	4.16	8.24	8.39
$\mathcal{A}_{Is16}$	10.551	1.489	1.6	12.74	$\mathcal{A}_{Is16}$	12.821	2.027	8.1	12.76	$\mathcal{A}_{Is16}$	12.051	2.611	3.161	3.10	$\mathcal{A}_{Is16}$	12.449	8.4	9.03	12.71
$\mathcal{A}_{Is18}$	11.061	1.749	5.8	12.26	$\mathcal{A}_{Is18}$	12.661	2.477	9.4	13.56	$\mathcal{A}_{Is18}$	11.691	2.391	3.521	3.45	$\mathcal{A}_{Is18}$	12.601	1.429	1.11	12.81