

# **Efficient tuning of automated machine learning pipelines** Nguyen, D.A.

#### Citation

Nguyen, D. A. (2024, October 9). *Efficient tuning of automated machine learning pipelines*. Retrieved from https://hdl.handle.net/1887/4094132

Version:	Publisher's Version
License:	<u>Licence agreement concerning inclusion of doctoral</u> <u>thesis in the Institutional Repository of the University</u> <u>of Leiden</u>
Downloaded from:	https://hdl.handle.net/1887/4094132

**Note:** To cite this publication please use the final published version (if applicable).



# Efficient AutoML via Combinational Sampling

In the previous chapter, the CASH approach converted the ML pipeline optimization problem into a hyperparameter optimization (HPO) problem, where the choice of algorithms was modelled as an additional categorical hyperparameter. In this manner, algorithms and their local hyperparameters are referred to at the same level. Consequently, this approach renders the resulting initial sampling less robust. Unlike the CASH approach, in this study, we used a new hyperparameter class to model the choice of the algorithm under the operator. Additionally, we propose a novel initial sampling approach to maximize the coverage of the AutoML search space to help BO construct a robust surrogate model. We experimented with both experimental scenarios of AutoML with two operators and six operators over 117 benchmark datasets, as introduced in Section 4. The results of our experiments demonstrate that the performance of BO is significantly improved using our sampling approach.

The remainder of this chapter is organized as follows. First, the motivation and introduction are provided in Section 7.1. Next, our contributions are highlighted in Section 7.2, Section 7.3 lays out the experimental setup. The experimental results are discussed in Section 7.4. Finally, the chapter is concluded, and further work is outlined in Section 7.5.

### 7.1 Introduction

Recall that existing AutoML approaches (e.g., [39], [40]) can be considered as optimization processes for which the best ML pipeline is searched. Each pipeline includes an architecture and a set of hyperparameter settings.

Bayesian Optimization (BO) is a commonly used approach in AutoML as it has been successfully used in hyperparameter optimization (HPO) problems and plays a role of an optimizer in many AutoML frameworks, e.g., AUTO-SKLEARN [39], Auto-Weka [40], and Hyperopt-sklearn (HPsklearn) [42]. BO is an efficient global optimization approach (in terms of the number of function evaluations) in which the trade-off between local exploitation and global exploration is well handled. Therefore, in this work, we focus on improving the BO by using it to solve the AutoML optimization problem. Traditionally, the AutoML optimization problem is treated as a HPO process, where the optimizer is inherited from the HPO domain. As HPO was originally developed to find the best hyperparameter setting from a single algorithm, it naturally does not consider the choice of algorithm. The choice of algorithm is then modelled as an extra categorical hyperparameter. Consequently, this HPO-based approach in handling the choice of algorithm mismatches the nature of the AutoML optimization problem.

The search space in the AutoML approach is largely owing to the many possible algorithm choices for pipeline operators. However, including many algorithms in the search space naturally leads BO to slow convergence or to get stuck in a local optimum [30], [31], [35]. One reason is that the initial sampling step in AutoML is typically restricted to a small budget, which is much smaller than the number of possible pipelines that can be constructed in the search space. The reason for this setting is that the effectiveness of BO becomes evident mainly in the later stages of optimization when it learns to produce better configurations. Many well-known sampling approaches, for example, the discrepancy-based quasi-random (quasi-random) sampling [30], the Latin Hypercube (LHD) sampling [239], have been employed for the initial sampling in this optimization context. However, they have shown themselves to be insufficiently robust [31], [240], [241] because they have been used in conjunction with the traditional approach of solving an AutoML optimization problem, which, as explained above, consists of converting it to a HPO, thus rendering obscure the differences between the choice of an algorithm and the choice of the algorithm's parameters.

Additionally, to construct a robust surrogate model, BO requires good coverage of the search space [31], but as the number of algorithms increases, the number of samples required to cover the search space increases exponentially. Previous studies [46], [47] pointed out that some algorithms can be grouped based on their technical behaviors.

To assess this theory, in this chapter, we propose a new two-fold approach to improve BO used in AutoML optimization:

• Group the similar operator algorithms when allocating initial sampling budget, e.g., the grouping of linear classifiers vs. the grouping of rule-based

Hyperparameter A	nnotation	Description
Continuous F Ordinal	[oatParam(min, max)]	Choose a float value in range of $[min, max] \cap \mathbb{R}$ Choose a integer value in range of $[min \ max] \cap \mathbb{Z}$
Nominal	$\operatorname{ateroricalParam}(C_1, \ldots, C_n)$	Choose a value in set $\{C_1, \ldots, C_n\}$
*Algorithm A	$\operatorname{lgorithmChoice}(A_1,\ldots,A_n)$	Choose a value in set $\{A_1, \ldots, A_n\}$
Locidor and H		urban a UrmonDonom has shilden
Infeasible F	orderivation at an $(a, content, t^{value}), (baran_1, \dots, baran_f)$ orbiddenParam(( $Param_1, \{P_{value(s)}\}$ ), ( $Param_2, \{P_{valve(s)}^2\}$ ))	when a hyper adam has cultured when the combination of $Param_1$ and $Param_2$ is forbidden
* Grouping H	$\operatorname{sperParam}(\operatorname{value}_{n}^{1},\ldots,\operatorname{value}_{v_{1}}^{1}),\ldots,\operatorname{value}_{n}^{n})$	each group <sub>i</sub> can be of any type: { <i>Continuous</i> , <i>Ordinal</i> , <i>Nominal</i> , <i>Algorithm</i> }
	udnorg Idnorg	

classifiers [46]. Table 7.1 summarises different hyperparameter classes with their semantics in our work.

• Building on top of other sampling approaches, we propose a novel sampling method that aims to allocate reasonable budgets for each set of algorithms to maximize the coverage of sampling areas in terms of the grouping of algorithms to provide a robust surrogate model. In other words, our proposed approach is complementary to other sampling approaches, rather than competitive, with the aim of optimizing the performance of the search space of AutoML.

## 7.2 The Proposed Approaches for Automated Machine Learning

In this section, we first introduce our proposed combination-based sampling approach for increasing the efficiency and robustness of AutoML. Next, we introduce a new BO Python library for AutoML optimization and an AutoML framework that implements this paradigm.

# 7.2.1 Novel combination-based initial sampling for Bayesian optimization for AutoML optimization

The central idea of our approach is to provide optimized coverage of the algorithmhyperparameter search space already during the initial sampling of BO in order to characterize the response surface more accurately.

To properly analyze this discussion, we need to utilize the notations that were introduced in Chapter 1 (Section 1.1.1). These notations are crucial for our ongoing analysis and were discussed in detail in their original context in Chapter 1 to ensure a better understanding. Given a search space denoted by  $\mathcal{M}$  includes the sequence of z operators  $\mathbb{O} = \mathcal{O}_1 \times \ldots \times \mathcal{O}_z$  and its corresponding hyperparameter spaces  $\Lambda = \Lambda_{\mathcal{O}_1} \cup \ldots \cup \Lambda_{\mathcal{O}_z}$ , as defined in Section 1.1.1 of Chapter 1.

A grouping of algorithms of operator  $\mathcal{O}_i$  assumes that the set of all algorithms  $\{\emptyset, \mathcal{A}_i^1, \ldots, \mathcal{A}_i^{n_i}\}$  available to be employed for operator<sup>1</sup>  $\mathcal{O}_i$  can be partitioned into  $g_i$  non-empty and non-overlapping subsets, according to their inner workings<sup>2</sup>:  $\{G_i^1, \ldots, G_i^{g_i}\}, g_i \leq n_i + 1^3$ . Such partitioning is called a grouping of algorithms.

<sup>&</sup>lt;sup>1</sup> if i < z or  $\{\mathcal{A}_z^1, \dots, \mathcal{A}_z^{n_z}\}$  if i = z.

<sup>&</sup>lt;sup>2</sup>or any other user-defined logic.

<sup>&</sup>lt;sup>3</sup> if i < z and  $g_z < n_z$  otherwise.

The operator can then be represented as  $\mathcal{O}_i = \{G_i^1, \ldots, G_i^{g_i}\}$ . According to our proposed combination-based initial sampling method (see Algorithm 9), the sequence of pipeline operators  $\mathbb{O} = \mathcal{O}_1 \times \ldots \times \mathcal{O}_z$  should be sampled in BO from the domain space of sets  $\{G_1^1, \ldots, G_1^{m_1}\} \times \ldots \times \{G_z^1, \ldots, G_z^{n_z}\}$  and the total initial sampling budget should be split equally per group. The main idea behind such sampling budget reallocation is the potential exploitation of similarities between algorithms within the group: sampling fewer of the same algorithms frees up the budget to be distributed to other (different) algorithms, thus improving the coverage of algorithm-hyperparameter search space at an earlier stage of BO.

As an input parameter for our method, we require a number of data points  $B_{init}$ for the initial sampling and a maximum number of combinations  $K, K \leq B_{init}$ . If K exceeds the maximum number of possible combinations computed from the input operation steps  $k = \prod_{i=1}^{z} |\mathcal{O}_i|$ , then we use Algorithm 10 to randomly regroup algorithms in operators to ensure  $k \leq K$ . The proposed sampling algorithm, presented in Algorithm 9, consists of the three following steps:

- 1. Generate the list of combinations: List all k possible combinations of groups for all z operators; apply RANDOMREGROUPING until k is small enough  $(k \le K)$  (lines 2 - 4).
- 2. Allocate budget to combinations: first allocate budget to all combinations based on the number of algorithms and hyperparameters behind (lines 5 10). Then, if there is any remaining budget  $B_{remain}$ , randomly allocate  $B_{remain}$  to the top  $\frac{k}{\eta}$  combinations ordered by their size (i.e. the number of algorithms and hyperparameters in the combination). We take the size of the combinations into account to give larger combinations a higher chance of getting a larger budget.
- 3. Sampling configurations: for each combination s, an existing sampling approach (e.g., LHD, quasi-random, here we use quasi-random) is used to generate a trial sequence  $s_j = (G_1, \ldots, G_z)$  (lines 12 16); Lastly, the generated configurations must be verified by CHECKFORBIDDEN<sup>4</sup>.

Lastly, the generated configurations are shuffled to remove a potential impact of grouped configurations based on combinations. This is highly recommended

<sup>4</sup> An external function that verifies a combination of algorithms/a configuration with the forbidden rules defined by the user.

#### Algorithm 9: Combination-based sampling

**Input:**  $\mathbb{O}$ : sequence of operators,  $\Lambda$ : hyperparameter spaces,  $B_{init}$ : number of initial samples, K: maximum number of combinations of grouping of algorithms over operators,  $\eta = 2$ : proportion of combinations to be chosen to assign more budget if any remaining budgets are available. **Output:**  $\Theta$ : set of configurations // 1-Generating combinations 1  $k = \prod_{i=1}^{z} |\mathcal{O}_i|$  // maximum number of possible combinations 2 if k > K then 3  $(\mathbb{O}, k) = \text{RANDOMREGROUPING}(\mathbb{O}, K) // \text{Algorithm 10}$ // Create a list s of all possible combinations from  $\mathbb O$ 4  $s = \{s_1, \ldots, s_k\} = \{G_1^1, \ldots, G_1^{g_1}\} \times \ldots \times \{G_z^1, \ldots, G_z^{g_z}\}$ // 2-Allocate budgets to k combinations 5  $l_c = rac{B_{init}}{k}$  // number of inital samples per combination 6  $m=rac{1}{k}\sum_{i=1}^k \left(|\Lambda_{s_i}|\!+\!|s_i|
ight)$  //  $|s_i|$  is the number of all unique algorithms and  $|\Lambda_{s_i}|$  is the number of hyperparameters 7  $\Theta = \emptyset$  // set of initial configurations s foreach  $j \in \{1, \ldots, k\}$  do  $l_j = \lfloor l_c imes rac{|\Lambda_{s(j)}| + |s_j|}{m} 
floor // l_j$  is the number of samples for the 9 combination  $s_j$  $l_j = \begin{cases} 1, & \text{if } l_j = 0. \\ l_j, & \text{otherwise.} \end{cases}$ 10 11 if  $B_{remain} = B_{init} - \sum_{j=1}^{k} l_j > 0$  then // Randomly allocate  $B_{remain}$  to the top  $rac{k}{\eta}$  combinations based on the number of algorithms and hyperparameters // 3-Sampling Configurations 12 foreach  $j \in \{1, ..., k\}$  do  $\Theta_j = \emptyset$  // feasible configurations in the  $j^{th}$  combination 13 while  $|\Theta_j| \leq l_j$  do 14  $\Theta_i = \Theta_i \cup \text{SAMPLING}(s_i, \Lambda_i, l_i - |\Theta_i|)$ 15 // SAMPLING is done via an existing approach, here we choose quasi-random sampling with minor adjustments foreach  $\lambda \in \Theta_i$  do 16 if CHECKFORBIDDEN( $\lambda$ ) then 17  $| \quad \Theta_j = \Theta_j \setminus \lambda$ 18  $\Theta = \Theta \cup \Theta_i$ 19

#### Algorithm 10: Random Regrouping

**Input:**  $\mathbb{O} = \left( \{G_1^1, \dots, G_1^{g_1}\} \times \dots \times \{G_z^1 \dots G_z^{g_z}\} \right)$ : sequence of operators, K: number of combinations **Output:**  $\mathbb{O}_{new}$ : new sequence of operators, k: new number of combinations 1 k = K // number of all possible combinations 2  $S = \emptyset$  // split solutions 3  $C_1 = \{1, \ldots, g_1\}, \ldots, C_z = \{1, \ldots, g_z\} //$  set of possible groupings of  $\mathcal{O}_{i \in \{1,\ldots,z\}}$ // List out all split solutions 4 Create a list of all possible splits  $H = \{h_i\}$  where  $h_i = (c_1, \ldots, c_z) : c_i \in C_i \forall j$ // Select split solutions which can produce k combinations,  $k \le K$ 5 while  $S = \emptyset$  do  $S = \{h = (c_1, \dots, c_z) \in H : (\prod_j c_j) = k\}$ 6 if  $S = \emptyset$  then 7 | k = k - 18 9  $s_{chosen} \sim \mathcal{U}(S) / /$  randomly choose one solution 10  $\mathbb{O}_{new} = (\emptyset_1, \dots, \emptyset_z), i = 1$ 11 foreach  $c_i \in s_{chosen}$  do //  $c_i$  is the number of groups to be created  $n_i = |\mathcal{O}_i|$  // number of groups in the  $i^{th}$  operator 12 if  $c_i = n_i$  then 13  $| \mathcal{O}_i = \{\{G_i^1\}, \dots, \{G_i^{n_i}\}\} / / \text{ when } c_i = n_i$  $\mathbf{14}$ else if  $c_i = 1$  then 15  $\mathcal{O}_i = \{G_i^1, \dots, G_i^{n_i}\} //$  merging all predefined groups 16 17 else  $G = \emptyset, \mathcal{O}_0 = \mathcal{O}_i$ 18 while  $n_i > 0$  do 19  $G_0 = \emptyset, n_{size} = \left\lceil \frac{n_i}{c} \right\rceil$ 20 if  $n_i > n_{size}$  then 21  $G_0 = \{\text{Random pick } n_{size} \text{ items in } \mathcal{O}_0\}$ 22 23 else  $G_0 = \{\mathcal{O}_0\}$ 24  $G = G \cup G_0, \, \mathcal{O}_0 = \mathcal{O}_0 \setminus G_0, \, n_i = |\mathcal{O}_0|$  $\mathbf{25}$  $\mathcal{O}_i = \{G\}$ 26  $\mathbb{O}_{new}^{(i)} = \mathcal{O}_i, \ i = i+1$ 27 28 return  $\mathbb{O}_{new}, k$ 

since, in some cases, the computational optimization budget, i.e., the run time limit, can run out before finishing this initialization step.

The RANDOMREGROUPING method used in Algorithm 9 is presented in Algorithm 10. For a sequence of operators that consists of multiple groupings of algorithms, it produces, via a regrouping, k combinations of operators ( $k \leq K$ ) using the following steps:

- 1. Step 1 (lines 3-9): Based on the number of the grouping in operators, we list out all possible solutions of regrouping to have k combinations.
- 2. Step 2 (line 10): Randomly choose one solution  $s_{chosen} = (c_1, \ldots, c_z)$  where  $c_i$  is the number of groupings to be created for the operator  $\mathcal{O}_i$ .
- 3. Step 3 (lines 11 27): For each operator  $\mathcal{O}_i$ , we randomly group algorithms into  $c_i$  groups.

#### 7.2.2 A New Optimization Library for AutoML Optimization

To take advantage of the new sampling approach introduced in Section 7.2.1, we introduce a BO library for AutoML optimization, named BO4ML<sup>5</sup>, where the new sampling approach is implemented. In this work, we use the Tree-structured Parzen Estimator (TPE) implemented in Hyperopt [153] for the surrogate model and Expected improvement (EI) [156] for the acquisition function.

#### 7.3 Experimental Setup

This study examines the two experiments introduced in Chapter 4 (Section 4.2 and Section 4.3). In both scenarios, we compare the performance of Bayesian optimization (see Chapter 3. Section 3.1.3) with and without our proposed initial sampling approach.

The first experiment uses similar parameter settings as in Chapter 5; we select two different values of the initial sample size 20 and 50. We use a budget of 500 function evaluations. The 5-fold cross-validation approach and the averaged geometric mean values over 10 repetitions are reported. The selected classification algorithms are not grouped together. The resampling techniques are grouped into

<sup>&</sup>lt;sup>5</sup> The library is published at https://github.com/ECOLE-ITN/NguyenSSCI2021.

four groups "Over-resampling", "Under-resampling", "Combine-resampling", and "No-sampling", as suggested in [47], [48].

In the second experiment, we used a budget of 50 samples for the initial sampling. All the experiments performed 10 runs with different random seeds, with a time limit of 1 hour. The performance of a single configuration is limited to 10 minutes with 4-folds cross-validation on training data, i.e., the evaluation of a fold is allowed to take 150 seconds. The evaluation of a configuration is aborted and returns zero if any folds have an error, for example, due to an infeasible configuration or timeout. Then, the average accuracy values for the test data over 10 runs are reported. Finally, the selected algorithms are grouped, according to the suggestions in [46], [151].

The implementation of the proposed methods is published in a git-repository<sup>5</sup> and PyPi-repository<sup>6</sup>. The experiment scripts for the reproducibility of the reported results are provided in a git-repository<sup>7</sup>.

## 7.4 Results and Discussion

In this section, we report and discuss the results obtained from using the above experimental setups. Our experiments has two objectives. First, we compare the performance of Bayesian optimization with the help of our proposed sampling approach with that without our contributions in terms of AutoML optimization for class-imbalance problems, with a search space of two operators. Second, we compared them against state-of-the-art AutoML frameworks with a search space of six operators.

#### 7.4.1 First experimental results

The results of the first experiment are presented in Table 7.2 to illustrate the performance of BO with and without the help of our proposed approach for two different initial sample sizes, that is, 20 (left, not shaded) and 50 (right, gray shaded). In both scenarios, the best performance for the corresponding dataset is highlighted in bold. A method that performs significantly worse than the best method according to the Wilcoxon signed-rank test with  $\alpha = 0.05$  is underlined. A value labeled \* indicates the best result obtained for the corresponding dataset. The two extra rows at the end display the additional summaries. The first extra

<sup>&</sup>lt;sup>6</sup>https://pypi.org/project/BO4ML.

<sup>&</sup>lt;sup>7</sup>https://github.com/ECOLE-ITN/NguyenSSCI2021/tree/assets/SSCI-Experiments.

Table 7.2: Average geometric mean (rounded to 4 decimals) based on two different initial sampling settings, i.e., the Hyperopt approach and our approach (BO4ML), over 10 repetitions for the 44 examined datasets, ordered by increasing imbalance ratio (#IR) value. The two extra rows display summaries for each scenario, i.e., 20 and 50 initial samples: (1) *Highest performance* shows the number of times the optimizer achieved the highest value. (2) *Significantly better performance* shows the number of times the optimizer was significantly better than the competitor.

Dataset	#IR	20 initial	samples	50 initia	l samples
		Hyperopt	BO4ML	Hyperopt	BO4ML
glass1	1.82	0.7935	0.7944	*0.7970	0.7944
ecoli-0_vs_1	1.86	0.9864	0.9864	0.9864	*0.9868
wisconsin	1.86	0.9814	0.9817	0.9818	*0.9819
pima	1.87	*0.7712	0.7696	0.7703	0.7707
iris0	2	*1	*1	*1	*1
glass0	2.06	0.8777	0.8748	0.8740	*0.8853
yeast1	2.46	0.7319	0.7332	0.7321	$^{*}0.7345$
haberman	2.78	*0.7049	0.7012	0.6991	0.7040
vehicle2	2.88	0.9908	*0.9927	0.9912	0.9918
vehicle1	2.9	0.8690	0.8684	0.8713	*0.8735
vehicle3	2.99	0.8463	0.8486	0.8416	*0.8506
glass-0-1-2-3 vs 4-5-6	3.2	*0.9567	0.9539	0.9534	0.9553
vehicle0	3.25	*0.9876	0.9867	0.9867	0.9867
ecoli1	3.36	0.9038	*0.9053	0.9050	0.9043
new-thyroid1	5.14	0.9980	0.9972	*0.9983	0.9966
new-thyroid2	5.14	*0.9972	0.9964	0.9952	0.9966
ecoli2	5.46	0.9363	0.9353	*0.9365	0.9360
segment0	6.02	*0.9993	0.9992	0.9992	0.9992
glass6	6.38	0.9488	0.9514	*0.9518	0.9511
veast3	8.1	0.9423	0.9421	0.9427	*0.9441
ecoli3	8.6	0.9038	0.9059	0.9064	*0.9072
page-blocks0	8.79	*0.9475	0.9472	0.9464	0.9457
veast-2 vs 4	9.08	0.9549	0.9542	*0.9554	0.9531
$v_{east-0.5-6-7-9}$ vs 4	9.35	0.8245	0.8177	*0.8261	0.8193
vowel0	9.98	0.9567	*0.9593	0.9525	0.9561
glass-0-1-6 vs 2	10.29	0.8404	0.8421	0.8334	*0.8460
glass2	11.59	*0.8504	0.8461	0.8462	0.8471
shuttle-c0-vs-c4	13.87	*1	*1	*1	*1
veast-1 vs 7	14.3	0.7991	0.8013	*0.8033	0.8010
glass4	15.46	*0.9390	0.9230	0.9299	0.9324
ecoli4	15.8	0.9712	0.9694	0.9632	*0.9737
page-blocks-1-3 vs 4	15.86	0.9931	0.9874	$\frac{1}{0.9917}$	*0.9944
abalone9-18	16.4	*0.8899	0.8829	0.8856	0.8859
glass-0-1-6 vs 5	19.44	0.9494	*0.9571	0.9564	0.9565
shuttle-c2-vs-c4	20.5	*1	*1	*1	*1
veast-1-4-5-8 vs 7	22.1	0.6989	0.7024	*0.7052	0.7045
glass5	22.78	0.9589	0.9558	0.9591	*0.9595
veast-2 vs 8	23.1	0.8136	*0.8348	0.8136	0.8150
veast4	28.1	$\overline{0.8764}$	0.8788	0.8782	*0.8788
veast-1-2-8-9 vs 7	30.57	0.7500	0.7489	0.7397	*0.7538
veast5	32.73	*0.9802	0.9798	*0.9802	0.9800
ecoli-0-1-3-7 vs 2-6	39.14	*0.9265	0.9076	0.9113	0.8982
veast6	41.4	0.8953	0.8918	0.8939	*0.8955
abalone19	129.44	0.7958	0.7974	0.7992	*0.7998
Highest p	erformance	15	8	12	19
Significantly better p	erformance	0	$\tilde{2}$	0	4



Figure 7.1: Illustration of the number of samples allocated to different combinations of methods in random sampling implemented in Hyperopt (top rows) vs. our proposed approach (bottom rows). The cases with 20 (left) and 50 (right) samples are shown.

row shows the number of times each scenario achieved the highest value over 44 datasets. The last extra row indicates the number of times each approach was significantly better than the others in the group. From the table, we can observe the following.

- In 20 initial samples scenario, Hyperopt achieved the best results on 28/44 cases, and our approach on 20/44 cases. However, our approach significantly wins on 2 tested cases, i.e., "ecoli3" and "yeast-2\_vs\_8" and is not significantly worse than Hyperopt in any tested cases.
- In the second scenario, our approach achieves the highest value on 31/44 cases and Hyperopt- on 16/44 cases. Similarly, our approach is not significantly worse than Hyperopt in any tested cases but significantly better on 4 examined datasets, i.e., "glass0", "yeast1", "ecoli4", "yeast-1-2-8-9\_vs\_7".

To investigate the sampling behavior of both approaches for the initial sample sizes, we provided two plots: Figure 7.1 shows the distributions of the samples and Figure 7.2 shows the distributions at the level of the individual algorithms. In both plots, the case with 20 samples is shown on the left and 50 on the right of the plot,



Figure 7.2: Illustration on the distribution of samples obtained via initial sampling methods on the level of individual methods, i.e., under resampling has 11 algorithms, combine resampling has 2 algorithms, over resampling has 7 algorithms and no resampling. The left part shows the case with 20 samples, while the case with 50 samples is shown on the right.

respectively. Looking at these figures, we observe that our approach samples all combinations of groupings over two operators for both sample sizes. By contrast, the sampling strategy used in the Hyperopt samples has less coverage in terms of these combinations. This is because we consider the choice of algorithms in operators to be different from categorical parameters, whereas Hyperopt does not. The plots clearly explain why BO performs better with the help of our approach.

#### 7.4.2 Results of second experiment

The results of the second experiment are presented in Table 7.3. The third and fourth columns show our experimental results, i.e., TPE with and without our sampling approach. The remaining columns contain the results obtained using other AutoML frameworks according to [22]. This table reports the average accuracy over

	<b>OpenML IDs</b>	Our Expe	eriments		Existin	g AutoML f	rameworl	<b>ks</b> [22]	
#TaskID	Dataset	TPE with	<b>TPE</b> without	Auto-s	klearn	HPsklearn	TPOT	ATM	H2O
	Name (#ID)	our sampling (RobustAutoML)	our sampling (Auto-Hyperopt)	SMAC	Random				
3	kr-vs-kp (3)	0.99656	0.99510	0.98986	0.99062	0.99051	0.99431	0.99326	0.99426
12	mfeat-factors (12)	0.98417	0.98117	0.97767	0.97633	0.94758	0.97333	0.98178	0.97433
15	breast-w $(15)$	0.97952	0.97048	0.96875	0.95873	0.96000	0.96571	0.98474	0.96286
23	$\operatorname{cmc}(23)$	0.57285	0.55158	0.54638	0.53262	0.53047	0.55882	0.58100	0.53733
$^{24}$	mushroom (24)	н	1	1	0.99993	1	1	1	0.99848
29	credit-approval (29)	0.88744	0.86522	0.87289	0.85507	0.85956	0.86377	0.89133	0.86184
31	credit-g (31)	0.76600	0.72733	0.73433	0.72400	0.70121	0.74400	0.76578	0.74867
41	sick $(42)$	0.95171	0.92878	0.91954	0.91911	0.92585	0.92732	0.94504	0.93122
53	soybean (54)	0.86457	0.83858	0.82008	0.81969	0.75787	0.81811	0.81522	0.82717
2079	vehicle (188)	0.69502	0.66018	0.63886	0.62670	0.64072	0.65566	0.64190	0.65570
3021	eucalyptus (38)	0.99152	0.98737	0.98288	0.98550	0.97438	0.98746	'	0.98419
3543	irish $(451)$	1	1	0.99019	0.99081	0.99404	0.99091	1	0.97967
3560	analcatdata_dmf (469)	0.23042	0.21125	0.20365	0.20382	0.19139	0.20833	0.27028	0.19542
3561	profb $(470)$	0.63119	0.64752	0.65687	0.64563	0.63762	0.66832	0.71221	0.71089
3904	jm1~(1053)	0.82404	0.81393	0.81344	0.81126	0.80998	0.81810	0.82100	0.74819
3917	kc1 (1067)	0.87393	0.85972	0.85118	0.85340	0.84044	0.86019	0.86856	0.80869
3945	KDDCup09_appete (1111)	0.98323	0.98197	0.98244	0.98228	0.98189	0.98182	ı	0.96555
3946	KDDCup09_churn (1112)	0.92901	0.92624	0.92725	0.92586	0.92599	0.92624	ı	0.78802
3948	KDDCup09_upsell (1114)	0.94345	0.94116	0.95094	0.95030	0.95068	0.95085	1	0.93415
7592	airlines (1590)	0.86251	0.85769	0.86938	0.87013	0.86727	0.87089	0.85448	0.86656
7593	bank-marketing (1596)	0.70278	0.80902	0.96395	0.89143	0.95227	0.94542	0.66390	0.92908
9910	blood-transfusi (4134)	0.80107	0.78073	0.78890	0.77762	0.77798	0.80249	0.77087	0.80044
9952	cnae-9 (1489)	0.91319	0.90826	0.89716	0.89205	0.89273	0.90450	0.89963	0.89205
9955	first-order-the (1492)	0.67167	0.65146	0.65172	0.62795	0.54667	0.61146	0.61097	0.56435
2000000000000000000000000000000000000	nomao $(1486)$	0.96525	0.95924	0.96903	0.96656	0.96891	0.97026	0.96055	0.97146
9981	phoneme $(1468)$	0.95093	0.94228	0.94167	0.93117	0.94012	0.94784	0.96049	0.95216
9985	one-hundred-pla (1475)	0.61029	0.59853	0.59695	0.58601	0.58293	0.61291	0.60272	0.61656
10101	adult (1464)	0.80667	0.76578	0.76667	0.77778	0.78044	0.78711	0.81956	0.73378
14952	covertype (4534)	0.97094	0.96623	0.96590	0.96244	0.96964	0.96913	0.96464	0.97160
14954	Bioresponse $(6332)$	0.83642	0.81111	0.79012	0.76173	0.76667	0.81009	0.81701	0.78333
14965	Amazon_employee (1461)	0.90307	0.90007	0.90447	0.90398	0.90451	0.90705	0.89957	0.90060
14967	PhishingWebsite (23380)	1	1	0.98265	0.99841	0.97131	1	ı	1
14968	GesturePhaseSeg (6332)	0.83580	0.80432	0.77353	0.77058	0.75823	0.81173	0.79155	0.80000
14969	MiceProtein (4538)	0.64001	0.61864	0.67733	0.65004	0.67272	0.67586	0.66217	0.70165
34538	cylinder-bands (4550)	1	70666.0	1	0.99907	0.99983	1	1	1
34539	cylinder-bands (4135)	0.94825	0.94557	0.94761	0.94444	0.94750	0.94891	0.94606	0.95114
125920	cjs (23381)	0.63133	0.56200	0.56667	0.55556	0.56844	0.56867	0.66978	0.58400
							continue	d on the n	ext page

ordered by #Task id. renetitions for the 73 OnenMIL datasets  $\overline{\mathbf{0}}$ ( ol o ..... ÷ Ŀ < Table 7.3:

### 7.4 Results and Discussion

Significar	Number	189356	189355	189354	168912	168911	168910	168909	168908	168868	168338	168337	168335	168332	168331	168330	168329	167141	167140	167125	167124	167121	167120	167119	146825	146824	146822	146821	146820	146819	146818	146817	146800	146607	146606	146212	146195			#TaskID		
nt wins over other approaches	of cases achieved	helena (41147)	jannis (41167)	dionis (1169)	volkert (41146)	robert (41143)	fabert $(41164)$	dilbert (41163)	riccardo (41142)	guillermo (41138)	MiniBooNE (41161)	albert $(41159)$	sylvine (41150)	jasmine (41165)	christine (41166)	APSFailure (41168)	jungle_chess_2p (41169)	Fashion-MNIST (40701)	climate-model-s (40670)	segment $(40978)$	wilt (40927)	steel-plates-fa $(40923)$	Australian (23517)	mfeat-pixel (41027)	Internet-Advert $(40996)$	car (40979)	MiceProtein (40984)	$CIFAR_{10} (40975)$	Devnagari-Scrip (40983)	churn (40994)	shuttle $(40981)$	dna (40982)	connect-4 (40966)	SpeedDating $(40536)$	numerai28.6 (23512)	higgs (40685)	dresses-sales (40668)		Name (#ID)	Dataset	OpenML IDs	+ + (
23	28	0.64810	0.73916	0.64626	0.95709	0.83259	0.67565	0.95040	0.73659	0.98985	0.98303	0.78205	0.92248	0.42507	0.60445	0.68479	0.31690	0.96273	0.96485	0.97713	0.39675	0.86652	0.52257	0.84647	0.84300	0.98433	0.94473	0.99441	0.97355	0.95802	0.88647	0.80497	0.99969	0.86611	0.70605	0.99965	0.77358	(RobustAutoML) (	our compline	TPE with	Our Exper	
1	5	0.64737	0.68112	0.63957	0.94655	0.80748	0.66177	$\frac{0.94437}{0.94437}$	0.72565	0.98900	0.98035	0.72707	0.91035	0.38497	0.59520	0.66670	0.29294	0.95367	0.95397	0.97033	0.37813	0.74910	0.52134	0.83956	0.83891	0.98117	0.93189	0.98748	0.97906	0.93951	0.85845	0.78216	0.99321	0.85871	0.69761	0.99945	0.77321	Auto-Hyperopt)	our campling	<b>TPE</b> without	iments	
9	11	0.68314	1	0.66665	0.93921	0.82009	0.70255	0.98357	0.74754	0.99287	0.74757	0.64227	0.94334	0.44843	0.66933	0.71814	0.30692	0.95620	0.95962	0.97774	1	0.74009	0.51926	0.86775	0.87844	0.97783	0.93088	0.97264	0.98612	0.94074	0.87053	0.78268	0.99043	0.86291	0.72296	0.99978	0.82109		SMAC	Auto-s		
I	T	0.66709	I	0.59845	0.94753	0.80603	0.67395	0.94793	0.73081	0.99137	0.75042	I	0.92891	0.39922	0.63762	0.69273	0.29566	0.95313	0.95889	0.97114	ı	0.02169	0.51939	0.85378	0.84450	0.97367	0.93333	0.97958	0.98581	0.92407	0.85556	0.76364	0.99506	0.86225	0.71930	0.99968	0.79628	Landom	Random	klearn	Existin	
1	3	0.66694	0.77971	0.65080	0.94675	0.80078	0.69104	0.97243	0.71630	0.99360	0.82518	0.74347	0.87477	0.34203	0.65451	0.68494	0.28741	0.94533	0.96109	0.97358	0.32093	0.86438	0.52033	0.88691	0.85060	0.98121	0.90664	0.98786	0.95289	0.92593	0.86913	0.75955	0.96380	0.86661	0.70743	0.99253	0.82886			HPsklearn	g AutoML f	
6	9	0.66110	ı	0.66895	0.95533	0.82366	0.68336	0.96254	0.72645	0.99339	0.98495	0.72548	0.93850	ı	0.65075	0.69642	0.33576	0.96000	0.95931	0.97398	0.29429	1	0.52082	0.88735	0.78089	0.96883	0.94055	0.99422	0.98540	0.94547	0.86184	0.79091	0.99506	0.86392	0.72031	0.99974	0.84123			TPOT	ramework	
13	18	0.80064	0.38666	0.63671	0.93476	0.79911	0.67357	0.95391	0.72169	0.97097	0.90729	0.66063	0.90234	0.35252	0.67940	0.63788	0.32108	0.95007	0.95282	0.96900	0.32001	0.89470	0.51941	0.87540	0.82114	0.97750	0.92564	0.96763	0.98657	0.96975	0.89050	0.76415	1	0.86128	0.67135	0.99955	0.77698			ATM	ks [22]	
13	15	0.64798	ı	0.61266	0.92510	0.80906	0.71752	0.96988	0.72811	0.99369	0.95625	0.81928	0.94604		0.67841	0.71786	ı	0.95370	0.96904	,	0.36389	0.58220	0.50635	0.90047	0.87341	0.97600	0.94185	0.99191	0.98574	0.93642	0.87633	0.78062	0.99551	0.84968	0.71281	0.99987	0.86500			H2O		

7. Efficient AutoML via Combinational Sampling



Figure 7.3: Comparison of all approaches against each other with the Nemenyi test with 5% significance level.

10 repetitions to illustrate the performance differences between the two implemented approaches in our AutoML framework<sup>8</sup>, i.e., TPE with (ROBUSTAUTOML) and without (Auto-Hyperopt) our sampling approach, to compare them against other well-known AutoML frameworks, i.e., AUTO-SKLEARN-SMAC (Auto-sklearn) and AUTO-SKLEARN-Random search (Random), HPSKLEARN, TPOT, ATM, and H2O. Values in bold indicate the highest values in the corresponding dataset. Underline values indicate significantly different results from the best method according to a Wilcoxon signed-rank test with p < 0.05. The two extra rows at the end show the additional summaries. The first extra row shows the number of times each approach achieved the highest performance over 73 examined datasets. The last row presents the number of cases in which these methods significantly outperformed the other compared methods.

The results allow the following insights:

• Comparing the results of approaches using the search space of AUTO-SKLEARN includes our two approaches, AUTO-SKLEARN and Random Search. First, it is not surprising that all Bayesian optimization approaches perform better than random search in most tested cases. This has been demonstrated in other studies [40], [47]. Second, AUTO-SKLEARN won more tested cases than AUTO-HYPEROPT with the same search space. A possible explanation for this might be that HYPEROPT lacks support for k-fold cross-validation yet, while SMAC, the BO variant used in AUTO-SKLEARN, uses racing algorithms to skip performing on unnecessary folds. Consequently, within the

 $<sup>^8{\</sup>rm For}$  readability, ROBUSTAUTOML stands for TPE with our sampling approach, and AUTO-HYPEROPT stands for the original version of TPE implemented by Hyperopt without our improvement.

#### 7. Efficient AutoML via Combinational Sampling

same budget of time, AUTO-HYPEROPT evaluated a much smaller number of configurations than AUTO-SKLEARN. Lastly, the experimental results clearly indicate that the performance of TPE with the help of our sampling approach significantly improves.

- From the results of three approaches using TPE, we can observe that: Firstly, comparing the two approaches that do not use our sampling, i.e., HPsklearn vs. AUTO-HYPEROPT, we can conclude that the search space of AUTO-SKLEARN does not improve the final performance of TPE. Secondly, the results clearly demonstrate that significant improvement was achieved with the help of our sampling approach. Our approach outperforms others 23 times, significantly winning AUTO-HYPEROPT in 16 cases and HPSKLEARN in 20 cases. Furthermore, in all 3 cases where AUTO-HYPEROPT achieves the highest results, e.g., tasks 24, 3543, and 14967, both our approach and AUTO-HYPEROPT get maximum accuracy in those cases. On the other hand, HPSKLEARN got the highest results in 3 cases, e.g., tasks 24, 146607, 189355, but never performed significantly better than our approach in any of those.
- Overall, our proposed approach shows the best results in more cases than all other approaches compared, namely 28/73. Moreover, according to the results of the Wilcoxon signed-rank test, our approach also significantly outperforms other compared approaches in 23/73 test cases. However, AUTO-HYPEROPT, without our improvement, does not win for any of the datasets.

When all approaches are compared, Friedman's test reveals a significant difference in average accuracy with  $p = 6.35 \cdot 10^{-11}$ . Thus, we performed a post-hoc multiple comparison test with the Nemenyi test ( $\alpha = 0.05$ ), shown in Figure 7.3. Approaches that have a distance higher than CD<sup>9</sup> are considered significantly different. According to this figure, we conclude that ROBUSTAUTOML is better than both TPE-based approaches and better than five other AutoML frameworks, such as, H20, ATM, AUTO-HYPEROPT, HPSKLEARN, and Random Search.

## 7.5 Conclusions and Future Work

In this chapter, we formulated AutoML as an optimization process for the machine learning pipeline. Then, we built on this paradigm, we proposed a new class for modeling the choice of algorithms and the concept of grouping algorithms. Second,

 $<sup>^{9}</sup>$ Critical Difference, here CD=1.2288.

a robust sampling approach for Bayesian optimization for AutoML optimization problems was introduced; Third, a BO approach for AutoML optimization was presented, where our proposed sampling approach and new hyperparameter classes were implemented. Lastly, a robust AutoML framework was presented which takes advantage of the proposed BO approach mentioned above.

The experimental results demonstrate the effectiveness of our approaches in two independent experiments over 117 datasets. The results clearly show significant improvement achieved by using our approach.

There are several interesting research directions for extending this study. First, we intend to apply the proposed sampling approach to other AutoML frameworks. Additionally, we plan to apply some pruning approaches such as Hyperband [35] and racing algorithm to reduce the time for evaluating configurations that are not promising by evaluating fewer folds.