# Minimum variance threshold for $\epsilon$-lexicase selection

Guilherme Seidyo Imai Aldeia
Federal University of ABC
Santo Andre, São Paulo, Brazil
guilherme.aldeia@ufabc.edu.br

Fabrício Olivetti de França
Federal University of ABC
Santo Andre, São Paulo, Brazil
folivetti@ufabc.edu.br

William G. La Cava
Boston Children's Hospital
Harvard Medical School
Boston, Massachusetts, USA
william.lacava@childrens.harvard.edu

## ABSTRACT

Parent selection plays an important role in evolutionary algorithms, and many strategies exist to select the parent pool before breeding the next generation. Methods often rely on average error over the entire dataset as a criterion to select the parents, which can lead to an information loss due to aggregation of all test cases. Under $\epsilon$-lexicase selection, the population goes to a selection pool that is iteratively reduced by using each test individually, discarding individuals with an error higher than the elite error plus the median absolute deviation (MAD) of errors for that particular test case. In an attempt to better capture differences in performance of individuals on cases, we propose a new criteria that splits errors into two partitions that minimize the total variance within partitions. Our method was embedded into the FEAT symbolic regression algorithm, and evaluated with the SRBench framework, containing 122 black-box synthetic and real-world regression problems. The empirical results show a better performance of our approach compared to traditional $\epsilon$-lexicase selection in the real-world datasets while showing equivalent performance on the synthetic dataset.

## CCS CONCEPTS

• **Theory of computation → Evolutionary algorithms**; • **Computing methodologies → Genetic programming**; **Supervised learning by regression**.

## KEYWORDS

lexicase selection, minimum variance, genetic programming

## 1 INTRODUCTION

Genetic Programming (GP) [15] is a meta-heuristic created to evolve computer programs to solve a particular problem, inspired by a weak metaphor of the biological evolutionary process. One application that is particularly suited for GP is symbolic regression (SR) [14, 18], a supervised learning method that searches for a function in the parametric family of functions and the optimal parameter that best fits a set of observed examples. Given a set of $d$ inputs of dimension $\mathbf{x} \in \mathbb{R}^n$ and respective target values $\{\mathbf{x}_i, y_i\}_{i=1}^d = (\mathcal{X}, \mathbf{y})$, SR optimizes both the structure and parameters of a model, $\hat{f}(\mathbf{X}, \hat{\theta})$.

An essential aspect of any evolutionary algorithm is the selection of solutions that will be recombined to generate new *offspring* solutions. The main goal of the recombination operator is to combine different parts of highly fit parents with the expectation that the offspring will improve upon its parents. As such, this selection requires that the parents perform well at different dataset examples. Traditional approaches such as tournament selection [3, 22] only use the aggregated information the fitness function provides to decide which parents will reproduce. As such, it can choose parents

that perform well on average across the whole dataset while missing individuals that perform well on difficult subsets of the problem. Unsurprisingly, evidence shows that parent selection with the entire data ignores relevant cases, as less information is available [16].

An alternative approach is to make the selection by distinguishing between the individuals by rewarding those who excel in difficult subsets of test cases. This is the idea behind the lexicase selection [25], which has shown state-of-the-art performance in program synthesis and symbolic regression [6]. It is based on the assumption that problem modularity is identifiable to some degree by its individual fitness cases, each representing a circumstance in which the program must do well. To select one parent, the entire population goes into a selection pool. Then, a random test case is selected, and individuals who fail that test are eliminated from the pool. Randomly selecting the following test case is repeated until the pool contains only one individual. If the procedure runs out of test cases, a random individual from the pool is returned with uniform probability.

While the original lexicase is suitable for binary cases (correct or incorrect test cases), it required modifications to handle continuous values in the SR context by defining a threshold $\epsilon$ [20]. This has shown to be effective for regression compared to other selection techniques. For each test case, a threshold is established based on the median absolute deviation of errors, and any individual within the tolerance remains in the pool.

This paper proposes a new way of estimating $\epsilon$ by replacing the median absolute deviation with a threshold that splits the errors into two partitions, such that the total sum of variance for the partitions is minimized. This proposed threshold criterion corresponds to a type of information entropy, used commonly by decision tree methods [2]. By minimizing the number of different individuals on either side, it behaves as a 1-dimensional clustering. This represents a way of picking the individuals that remain in the pool that is more sensitive to similarities in performance among the pool.

We evaluate the proposed method by implementing it into FEAT [19], a symbolic regression framework that has been previously evaluated with $\epsilon$-lexicase selection [18] and successfully applied to real-world problem [21]. Our experiments are twofold: first, with 6 low-dimensional datasets, we provide an in-depth analysis of its performance characteristics during the evolution. Subsequently, we rigorously assessed the proposed selection process utilizing the SRBench [18], a unified framework designed to perform a relative performance evaluation of symbolic regression algorithms.

Our empirical findings show that our method presents a superior performance in terms of $R^2$ while keeping solutions within the same level of complexity when compared to the $\epsilon$-lexicase selection. In

SRBench, our method improved the FEAT algorithm in real-world problems and was still capable of increasing its ranking.

The remainder of the paper is organized as follows. §3 details our proposed modification by replacing the median absolute deviation with a more effective threshold optimization strategy. §4 presents the symbolic regression framework used in the experiments. §5 outlines experimental design, datasets, and evaluation criteria. Results and discussion are presented in §6 offering a comparative analysis of our approach against $\epsilon$-lexicase. Finally, §7 summarizes findings, discusses implications, and suggests future avenues for research in evolutionary algorithms and parent selection methodologies.

## 2 $\epsilon$-LEXICASE SELECTION

Let us define a test case as a pair $(\mathbf{x}_i, y_i)$, $i = 1, 2, \ldots, d$ from the data split used during fitness evaluation. We denote $t \in \mathcal{T}$ any possible test case. The error of an individual $n$, denoted as $e_t(n)$, is the absolute difference between its predicted value given $\mathbf{x}_t$ and the target value $y_t$. The vector of errors for a test case $t$, denoted as $\mathbf{e}_t$, is the concatenation of errors for all individuals from the set. The set $\mathcal{N}$ denotes the population, and the set $\mathcal{S}$ denotes the selection pool. Initially, La Cava et al. [20] proposed using a threshold in lexicase to handle continuous values. The method was further studied by them in later work [17], from which different implementations of $\epsilon$-lexicase were proposed: static, semi-dynamic, and dynamic.

In the static scenario, the criteria to stay in the pool uses $e_t^* + \epsilon$, where $e_t^*$ is the best error from the population to that case. In semi-dynamic and dynamic implementations, the threshold to stay in the pool is calculated as the error of the `elite` (smallest error from the pool) plus $\epsilon$. The difference between the semi-dynamic and the dynamic is that the former calculates $\epsilon$ over the entire population, while the latter calculates it for the pool. In any case, $\epsilon$ is calculated as the median absolute deviation (MAD) [23]:

$$\epsilon_t = \lambda(\mathbf{e}_t) = \text{median}(|\mathbf{e}_t - \text{median}(\mathbf{e}_t)|). \tag{1}$$

The MAD produces median-centered random variables with more robustness than standard deviation, as the median is more representative of the center than the mean in asymmetrical distributions [23]. The interpretation of the MAD as a criterion in the $\epsilon$-lexicase context is a dispersion measured from the center of the distribution.

## 3 MINIMUM VARIANCE THRESHOLD

In this paper, we propose the Minimum Variance Threshold (MVT) to replace MAD as a selection criteria for the static and dynamic $\epsilon$-lexicase algorithms. The MVT split the errors into two clusters, reducing the pool of individuals to those below the specific threshold. This criterion is inspired by regression trees [2]. Essentially, the parent selection method changes the selection distribution to specific cases, and we hypothesize that performing the selection based on clustering good and bad performing individuals can improve the overall performance of $\epsilon$-lexicase selection.

In each step of MVT lexicase, the threshold $\tau*$ for staying in the pool is such that the sum of the variance of the errors in the pool on either side of the threshold is minimized. To estimate the threshold $\tau*$, we calculate:

$$\tau^* = \min_{\min(\mathbf{e}_t) < \tau < \max(\mathbf{e}_t)} \left( \frac{\text{Var}(\mathbf{l})}{|\mathbf{l}|} + \frac{\text{Var}(\mathbf{r})}{|\mathbf{r}|} \right), \tag{2}$$

where $\mathbf{l} = [e_t : \mathbf{e}_t | e_t < \tau]$ and $\mathbf{r} = [e_t : \mathbf{e}_t | e_t \geq \tau]$ are the left-sided and right-sided partition, respectively.

The procedure is to search for a value that splits the errors by clustering them into one group with smaller error values (left) and one with greater error values (right). Figure 1 illustrates this process. Every time we perform a split, the remaining pool represents all individuals who consistently performed better in the error cases (that is, on the left side of the split). We see this process as a more intuitive way of splitting the pool.
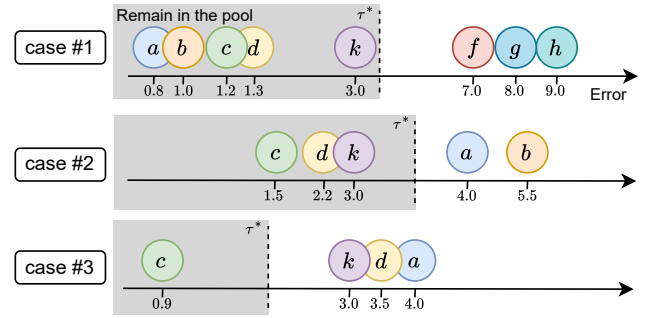


Figure 1: Process of consecutively splitting the pool of individuals into two clusters based on their error ($x$ axis). The process consists of randomly picking a test case, estimating $\tau^*$ by solving Eq. 2, and removing individuals with errors higher than the pool threshold. This is repeated until only one individual remains in the pool, or all training data was already used as singular test cases — returning one random individual from the remaining pool.

## 4 FEATURE ENGINEERING AUTOMATION TOOL

The Feature Engineering Automation Tool (FEAT) was proposed in [19] as an application of symbolic regression and was previously benchmarked in SRBench with good results [18]. The idea is to evolve a set of trees used as features in another machine learning (ML) model. We used FEAT as our SR framework to evaluate our proposed method. Our choice is motivated by semi-dynamic $\epsilon$-lexicase being the default selection method and, additionally, because the model was studied in previous works with different lexicase methods [19], also successfully applied to clinical decision modeling [21].

The algorithm follows the standard steps of an evolutionary algorithm, with some specific particularities besides the $\epsilon$-lexicase selection. First, it performs a multi-objective optimization, minimizing both the fitness and complexity of models. Second, it represents individuals as a collection of symbolic regression trees. Third, it uses a backpropagation algorithm to optimize the parameters of

Table 1: Comparison of $\epsilon$-lexicase variations.

| Strategy | Pool | Criteria | Characteristic |
|---|---|---|---|
| Static | pop ($\mathcal{N}$) | $e_t^* + \epsilon_t$ | The error vector is calculated previously with entire population. A binary mask matrix stores the global criteria to iterate through the random cases |
| Semi-dynamic | pop ($\mathcal{N}$) | elite $+ \epsilon_t$, where $\epsilon \leftarrow \lambda(\mathbf{e}_t)$ for $t \in \mathcal{T}$ | The error vector is calculated previously, but the decision criteria is based on the individuals that are in the pool, as it uses the elite error in the criteria |
| Dynamic | pool ($\mathcal{S}$) | elite $+ \epsilon_t$, where $\epsilon \leftarrow \lambda(\mathbf{e}_t(\mathcal{S}))$ | Criteria is dependent of the elite, and $\lambda$ is calculated with the selection pool |

the models. While the selection is made with $\epsilon$-lexicase, survival of the population is done using the non-dominated sorting genetic algorithm (NSGA2) [5].

Each individual is the combination of $\phi_0, \phi_1, \ldots, \phi_m$ expression trees representing one new feature. These features are used as inputs for any machine learning model (using a linear regression with $l_2$ regularization [11] by default). Figure 2 illustrates a possible individual.
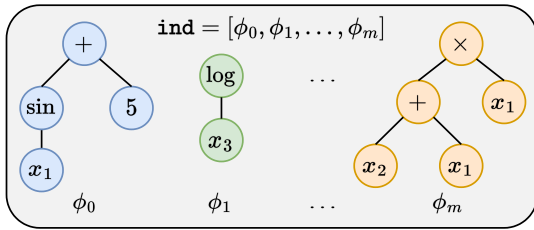


Figure 2: An individual in FEAT is a collection of symbolic regression trees as meta-features for any machine learning model.

The evolutionary loop of FEAT performs the evolution of a population of individuals as an application of evolutionary computing for feature engineering. The fitness of an individual is calculated by taking an ML model $f$ and training it with the collection of features corresponding to that individual, then evaluating the mean squared error between the predictions of the ML model $\widehat{\mathbf{y}} = f([\phi_0, \phi_1, \ldots, \phi_m])$ and observed values $\mathbf{y}$:

$$\text{MSE}(\widehat{\mathbf{y}}, \mathbf{y}) = \frac{1}{d} \sum_{i=1}^{d} (\widehat{y_i} - y_i)^2. \qquad (3)$$

The complexity of a model is defined recursively for a node $n$ with $k$ arguments as the combination of the complexity of its children with the complexity of the node itself:

$$C(n) = c_n * \left( \sum_{a=1}^{k} C(a) \right). \qquad (4)$$

The complexity of an individual is the sum of the complexity for each feature $\phi$ it contains, and the model size is the number of nodes for each feature $\phi$. The complexity of each feature increases exponentially as the depth increases but tends to reach higher values if deeper nodes are more complex. This way, the occurrence

of complicated operators is discouraged, and we focus on more interpretable models.

Crossover and mutation are implemented as follows. Crossover can either swap two subtrees between two individuals or swap two features (whole trees). For mutations, the algorithm uses *point*, *insert*, and *delete* mutations (randomly replace, insert, or delete one node, respectively, as well as *insert/delete* dimension, that creates or removes a new tree representing a dimension.

FEAT implements two mechanisms to split the training data further. The first, called *validation split*, internally separates the training data into two partitions: one for fitting the parameters and evaluating the model and validation data for logging and picking the final individual from the last population. The second, called *batch learning*, generates a random batch from the inner training partition to perform the fit and can be used to implement the downsampling strategy proposed in [8].

Mainly related to lexicase, previous work showed that downsampling the test cases can achieve similar or better performance [1, 9, 10], and $\epsilon$-lexicase was also shown to hold these benefits [8].

## 5 METHODS

Our experiments are twofold. First, we evaluate the FEAT algorithm with our proposed parent selection schema with 30 runs for 6 datasets. For the first set of experiments, we focus on performing longer runs and obtaining a larger sample size to calculate statistics and perform an in-depth analysis. We changed the NSGA2 survival step by replacing the original population with the offspring to isolate the effect of parenting selection, as NSGA2 is an elitist algorithm. The batch size was also turned off to assess the number of test cases each variation performed.

After the first batch of experiments, we evaluate it through a rigorous benchmark of symbolic regression algorithms known as SRBench [18], containing 122 black-box regression problems, taken from the Penn Machine Learning Benchmarks (PMLB) [24], which include the Friedman problems [7] (62 datasets). The Friedman contains synthetic problems designed to be tricky to solve by regression models, while the rest of the datasets contain real-world data. For the SRBench, we performed 10 runs for each dataset, as specified by the benchmark. We kept every hyper-parameter to correspond to the benchmark settings, allowing us to compare our results with the latest published result.

Regardless of whether it was from the first or second batch of experiments, every run was done by splitting the data into train

and test partitions in a ratio of .75/.25. Table 2 and Figure 3 show the dimensionality of each dataset used in the experiments.

**Table 2: Dimensionality of the six datasets used to perform an in-depth analysis.**

| Dataset | # samples | # features |
|---|---|---|
| airfoil | 1503 | 5 |
| concrete | 1030 | 8 |
| energy cooling | 768 | 8 |
| energy heating | 768 | 8 |
| Housing | 506 | 13 |
| Tower | 3135 | 25 |



**Figure 3: Dimensionality of the SRBench datasets.**

## 6 RESULTS AND DISCUSSION

Table 3 shows the hyper-parameter configuration used in the SRBench experiment, with values in parenthesis showing the settings used in the first 6 datasets. The batch size was selected after reading the results obtained from the first batch of experiments. The complexity for each operator specified in the hyper-parameter is depicted in Table 4, reminding the reader that these values are pre-defined by the FEAT package authors.

Original FEAT implementation is available at https://github.com/cavalab/feat. FEAT with the MVT selection is available at https://github.com/gAldeia/feat. All data and the source code for implementations, experiments, and post-processing analysis are available at https://github.com/gAldeia/srbench/tree/feat_split_benchmark.

## 6 RESULTS AND DISCUSSION

We divide this section into three parts. First, we analyze the convergence aspects of the different $\epsilon$-lexicase algorithms, focusing on how they affect the search convergence and the number of test cases they use. Then, we benchmark the proposed method with the SRBench framework, which consists of 23 machine learning methods (originally 21 plus our 2 proposed algorithms), of which 16 (14 plus ours) are symbolic regression algorithms. Finally, we also use the SRBench to asses how the proposed method scales with the number of features and number of samples of the dataset.

As for naming, FEAT($\epsilon$-lex) stands for the standard FEAT with default $\epsilon$-lexicase selection. D-Split and S-Split are our dynamic and static methods of estimating the threshold using FEAT.

**Table 3: FEAT hyper-parameters shared between all evaluated variations. Values in parenthesis indicates a configuration used just with the 6 datasets.**

| Parameter | Value |
|---|---|
| objectives | ["fitness","complexity"] |
| pop_size | 100 |
| gens | 100 (350) |
| cross_rate | 0.5 |
| ml | Linear ridge regression |
| max_depth | 6 (3) |
| backprop | True |
| iters | 10 |
| validation split | 0.25 |
| selection | NSGA2 (offspring) |
| batch_size | 200 (unlimited) |
| functions | $[+,-,*,\div,(\cdot)^2,(\cdot)^3,\sqrt{\cdot},\sin,\cos,e^{(\cdot)},\log]$ |

**Table 4: Complexity of each operator**

| Complexity | Operators |
|---|---|
| 1 | $+$, $-$ |
| 2 | $\div$, $*$, $(\cdot)^2$, $\sqrt{\cdot}$ |
| 3 | $\cos$, $\sin$, $(\cdot)^3$ |
| 4 | $e^{(\cdot)}$, $\log$ |

### 6.1 Behavior during the run

Figure 4 reports the minimum loss on the validation split for each dataset. Figure 5 reports the median number of test cases used to select the parent pool for each generation.

The D-Split shows a better convergence curve in two datasets (Airfoil and Energy Cooling) than other approaches, supporting the idea that the MVT can change the convergence curves. S-Split shows the worst convergence curve for all datasets. In some cases, such as in Concrete, Energy Cooling, and Tower datasets, it seems to reach a *plateau* with a higher error rate.

When looking at the number of test cases used in each generation, we notice that for all the datasets, the original $\epsilon$-lexicase requires a smaller number of cases, implying a faster execution time (as we will discuss later). Both D-Split and S-Split used more test cases, with S-Split using twice as many test cases as D-Split. This implies that many test cases have large clusters of good-performing individuals that would be ignored in the selection using the MAD, but with the MVT, they are kept in the pool for longer. We believe that this creates a more robust selection; thus, we have improvements over the original FEAT with the MAD threshold.

We believe that clustering the whole population can lead to high thresholds and a less effective selection pool, and the S-Split could be improved by implementing a down-sampling strategy that selects a subset of test cases where it is more likely to have significant differences between the individuals [1]. The D-split — which uses the pool to calculate the threshold— yielded a better

Figure 4: Convergence loss of the best individual on validation partition for the six problems.



Figure 5: Median number of test cases used to pick each parent for the six problems.

performance because the pool shrinks, removing individuals with a subpar overall goodness-of-fit and getting more stable results.

## 6.2 Performance on small datasets

To assess the final performance for the first 6 datasets, we report in Figure 6 the median $R^2$ ranking for each random seed used in the experiments, and in Figure 7, we show the size and complexity of the final models. We also included a critical differences diagram [12], showing each algorithm's median ranking, with horizontal lines connecting them when no statistical significance is observed from the samples.

Looking at the $R^2$ rank (Figure 6a) we see S-Split as the worse variant, with 3 asterisks indicating p-value $\leq 1 \times 10^{-3}$. The critical differences (6c) shows that D-Split is mainly classified as first among the other variants. When looking at the grouped $R^2$ for all datasets (Figure 6b), there is no statistical difference, but 75% of the distribution of the D-Split is above the median for the original $\epsilon$-lexicase. The worst performance is from the S-Split, with a large spread between $[0.3, 0.8]$.
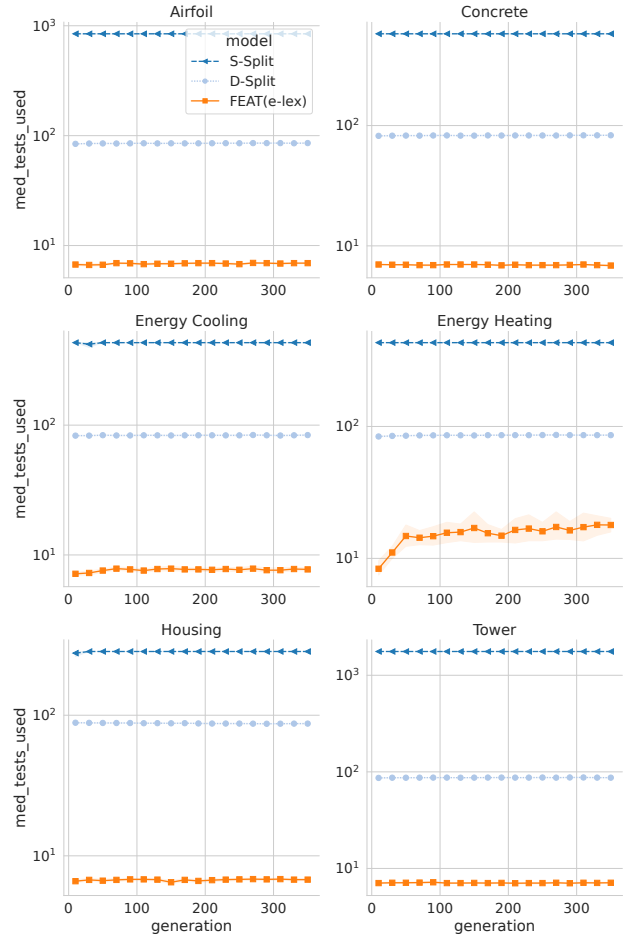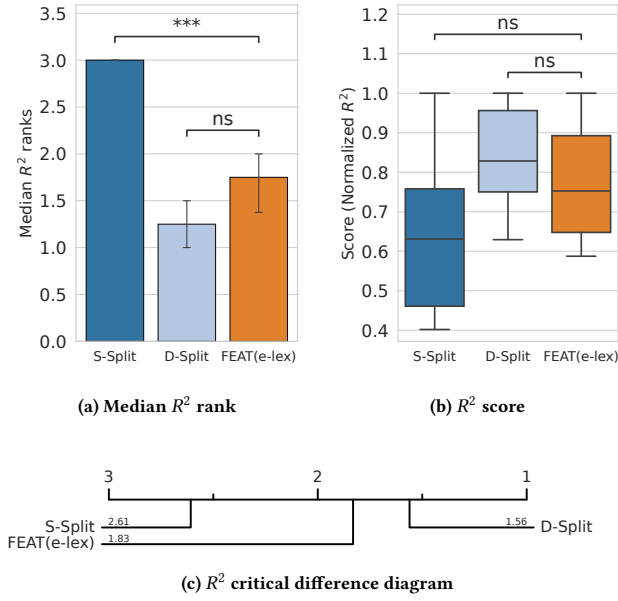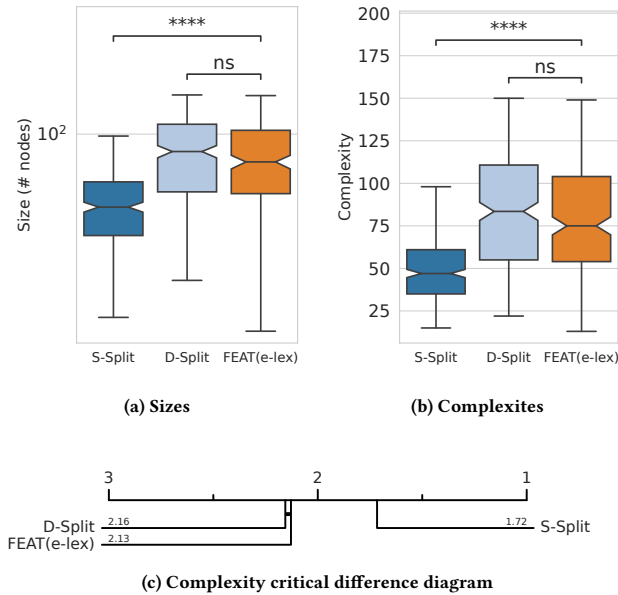
The worst performance of S-Split shows a better trade-off in terms of size and complexity (Figures 7a and 7b). The critical differences diagram shows that, while S-Split is the best algorithm in size and complexity, the other two variants, D-Split and $\epsilon$-lexicase, have no statistical difference.

## 6.3 Benchmarking with SRBench

The SRBench framework comprises several machine learning models, some of which are symbolic regression algorithms, meaning they could benefit from our proposed lexicase selection. We notice that, from the last published analysis of state-of-the-art symbolic regression methods, FEAT was one of the best-performing algorithms, originally ranked at 3 in terms of the root of $R^2$ on test partition, behind Operon [13] and SBP-GP [26]. Figure 8 reports the latest available results from SRBench, combined with our results of running our two variations, FEAT S-Split and FEAT D-Split, over the benchmark problems.

The FEAT D-Split outperforms the original FEAT by a modest difference, showing approximately similar RMSE, $R^2$, and model

(a) Median $R^2$ rank

(b) $R^2$ score

(c) $R^2$ critical difference diagram

**Figure 6: Grouped $R^2$ results for the 6 datasets.**



(a) Sizes

(b) Complexites

(c) Complexity critical difference diagram

**Figure 7: Grouped complexity results for the 6 datasets.**

size. FEAT S-Split performs worse than the other two variations of FEAT but still ranks 9th, outperforming more than half of the methods.

As pointed out in [4], the aggregated result of SRBench can mask the difference among the top algorithms. A division of the benchmarks into Friedman and non-Friedman datasets (roughly half and

half of the datasets) revealed that the best-ranked algorithms had a noticeable difference on the Friedman benchmarks. At the same time, they all behaved similarly in the non-Friedman sets. Following [4], we report in Figure 9 the percentage of times that each of the algorithms was ranked top 5 for the Friedman synthetic problems, and Figure 10 reports the percentage each algorithm was ranked top 1 in the real-world, non-Friedman problems. This represents a relative goodness-of-fit performance when pairing each instance of results, and higher values are better. We did not report the top 1 percentage for Friedman datasets because no FEAT variation reached more than 4% of wins, implying that these problems are still challenging for FEAT.

While FEAT achieves a rank of 5 or less approximately 53% of the time, our proposed method, FEAT D-Split, can increase this percentage to over 60%. FEAT with S-Split selection is among the top 5 algorithms in less than 10% of the runs. The D-Split performs better than the other variants, with a total of 10% of wins among the other algorithms. This represents a significant improvement, as the original FEAT performs poorly in this set of problems. We also noticed that S-Split performed better for the non-Friedman datasets than the original FEAT. Given that S-Split was better at finding models with smaller size and complexity, this implies that it found solutions with a better trade-off between $R^2$ and complexity.

Overall, these histograms show that FEAT D-Split increased its rank in the Friedman benchmarks by a small margin, but as already noticed, not enough to be among the top-1 ranks. D-Split and S-Split significantly increase solutions at the top-1, with D-Split ranking in the top-3 algorithms for the non-Friedman benchmarks. Again, as noted in [4], the non-Friedman median $R^2$ does not have much influence on the overall rank. Thus, we could not observe much difference in Fig.8. Notice that this plot discretizes the results presented in the previous plots. The fact that it was not ranked first in a particular instance does not imply that the returned result is not close to the best result. This is better highlighted on the error bar plots.

Finally, we plot the rank for model size and $R^2$ for all algorithms in SRBench in Figure 11. For both model size and $R^2$, smaller values are better. The lines in the plot represent the Pareto fronts for each rank.

We can see that all variations are within the same Pareto front rank, with different levels of trade-off between model size and accuracy. While showing equivalent performance in the synthetic benchmarks, we saw that the lexicase selection achieved a better result in real-world problems.

## 6.4 Scalability

By plotting the execution time versus number of samples or features for the datasets in SRBench, we can visualize how different implementations scale. We re-run the entire FEAT algorithm with SRBench datasets to remove hardware-related differences. Figure 12 reports the training time for different numbers of features in the dataset, and Figure 13 reports the training time for different numbers of samples. We did not compare execution time for other models as they are hardware-dependent. Thus, we could get a biased estimate of performance.
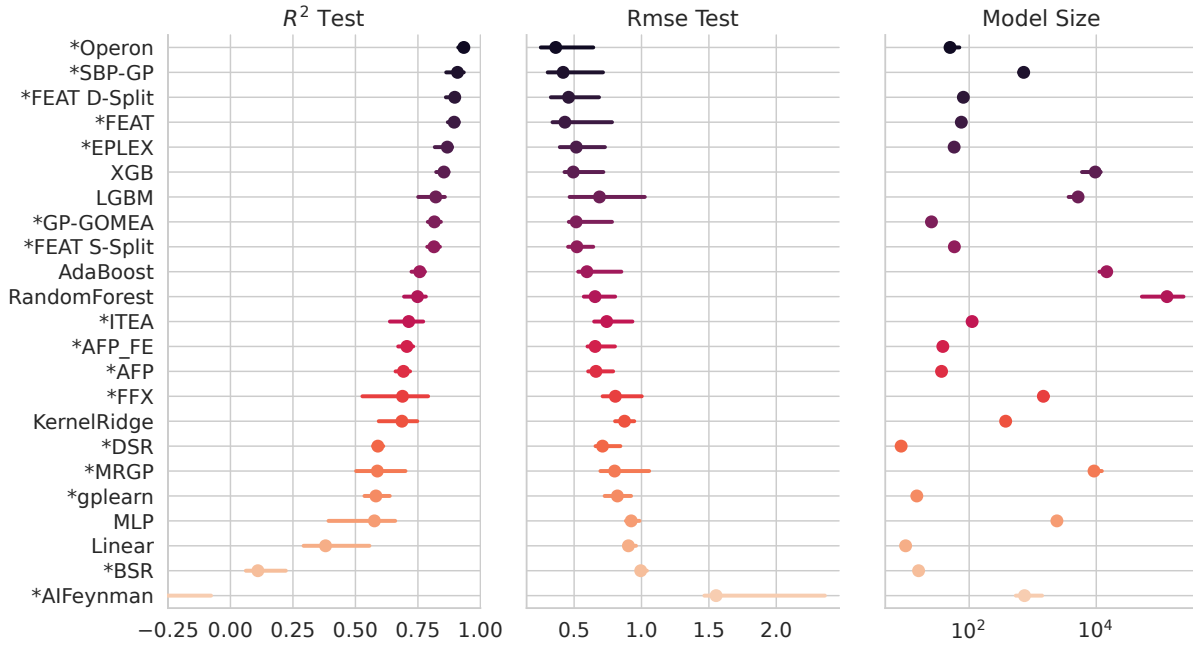
**Figure 8: Comparison of our proposed algorithms with SRBench results. Names with an asterisk indicate that the model is a symbolic regression algorithm. Bars show the** 95% **confidence interval. Size comparisons are made by the number of nodes to keep models comparable.**
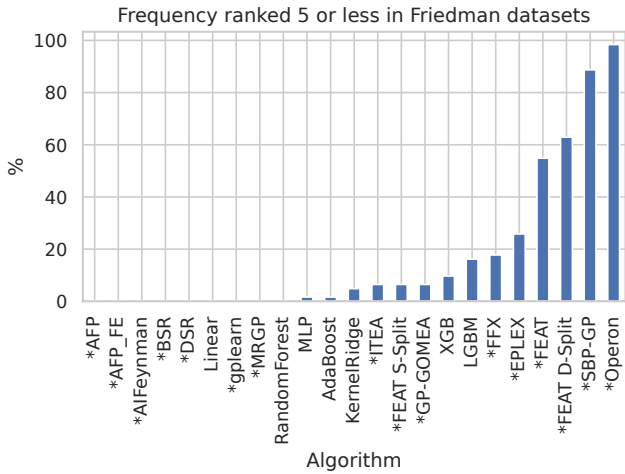


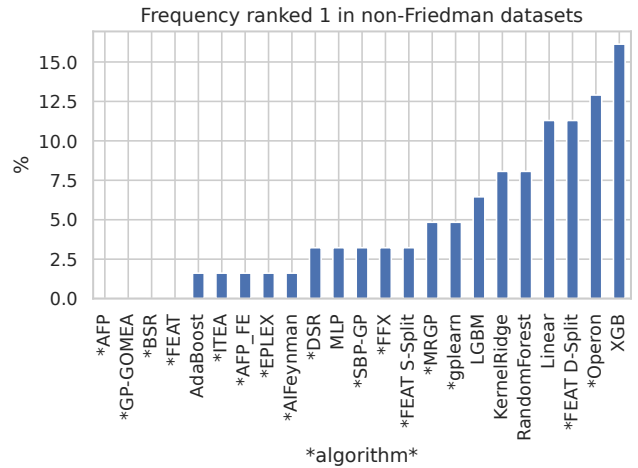**Figure 9: Number of times each algorithm was ranked top 5 or less for each run for the Friedman datasets.**



**Figure 10: Number of times each algorithm was ranked top 1 or less for each run for the non-Friedman datasets.**

All three variants of FEAT with different parent selections have similar curves but different offsets. The results imply that scalability is not correlated with the number of features but with the number of samples. This may be explained because it increases the number of test cases, and our MVT selection uses a higher number of test cases than the original $\epsilon$-lexicase.

As for the number of features, the size of the individuals is limited by the maximum size of symbolic regression trees, and, as

with many symbolic regression algorithms, they have to perform the feature selection implicitly due to size constraints. While the number of features does not impact training time, the maximum expression sizes do. However, we achieved a good performance on SRBench with relatively small models, showing no need to increase the maximum allowed size of the models and implying FEAT could perform well training with the current settings.
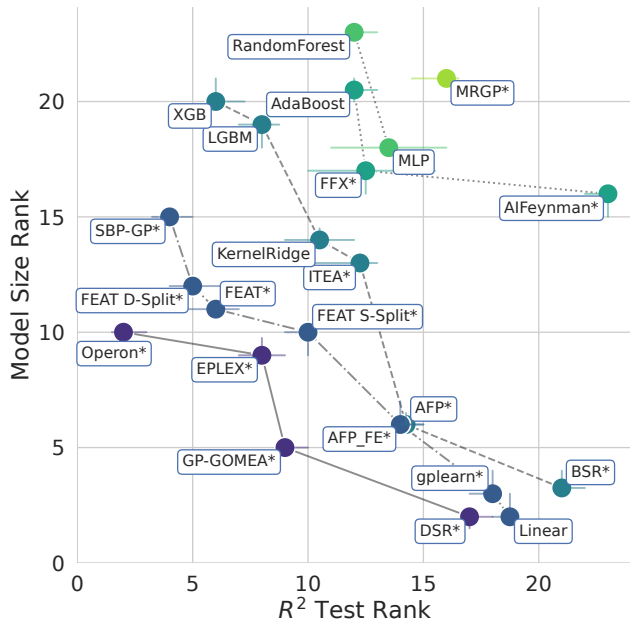
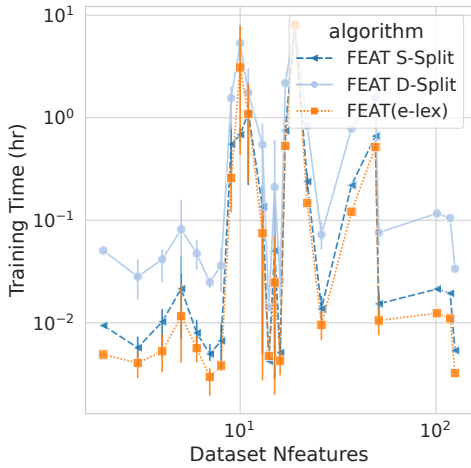**Figure 11: Pareto front of benchmarked algorithms**



**Figure 12: Training time versus number of features in the dataset**

One interesting result is that S-Split uses more tests than D-Split but still has a faster execution. We hypothesize that this is because D-Split needs to calculate the threshold for each test case until one parent is selected, while S-Split needs to calculate it only once.

## 7 CONCLUSIONS

In this paper, we proposed a novel way of estimating the threshold to stay in the pool for $\epsilon$-lexicase selection, an adaptation to improve lexicase for regression problems. We then evaluated these new threshold criteria with 6 small problems and 122 regression
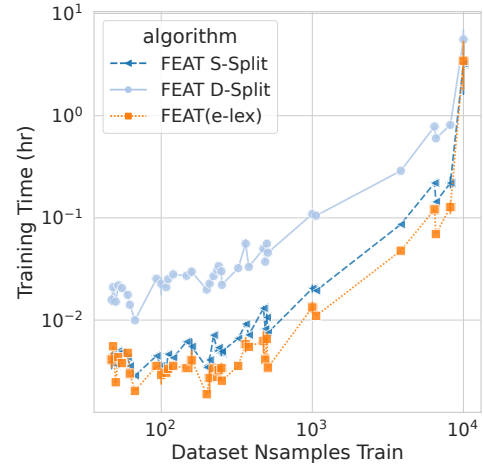


**Figure 13: Training time versus number of samples in the training partition of the dataset**

problems, comparing the results with other 21 algorithms. We achieve competitive results with the previous version of $\epsilon$-lexicase while showing improvements in real-world datasets. The performance in synthetic problems was preserved. The downside of our implementation is that it requires more test cases to run, leading to higher execution times. Nevertheless, using additional resources helps achieve better results, as shown by the final performance comparisons.

In future work, we plan to study how these new criteria change the distribution probability of individuals in the population. We will also investigate how down-sampling can decrease the number of test cases used. Finally, these criteria can also be implemented into other symbolic regression algorithms, and we can assess whether they can be improved with the minimum variance split criteria.

## ACKNOWLEDGMENTS

## REFERENCES

[1] Ryan Boldi, Martin Briesch, Dominik Sobania, Alexander Lalejini, Thomas Helmuth, Franz Rothlauf, Charles Ofria, and Lee Spector. 2024. Informed Down-Sampled Lexicase Selection: Identifying productive training cases for efficient problem solving. *Evolutionary computation* (01 2024), 1–32. https://doi.org/10.1162/evco_a_00346

[2] L. Breiman, J. Friedman, C.J. Stone, and R.A. Olshen. 1984. *Classification and Regression Trees.* Taylor & Francis. https://books.google.com.br/books?id=JwQx-WOmSyQC

[3] Anne Brindle. 1980. Genetic algorithms for function optimization. (1980).

[4] Fabrício Olivetti de França. 2023. Transformation-Interaction-Rational representation for Symbolic Regression: a detailed analysis of SRBench results. *ACM Transactions on Evolutionary Learning* (2023).

[5] K. Deb, A. Pratap, S. Agarwal, and T. Meyarivan. 2002. A fast and elitist multiobjective genetic algorithm: NSGA-II. *IEEE Transactions on Evolutionary Computation* 6, 2 (2002), 182–197. https://doi.org/10.1109/4235.996017

[6] Li Ding, Edward Pantridge, and Lee Spector. 2023. Probabilistic Lexicase Selection. In *Proceedings of the Genetic and Evolutionary Computation Conference*. 1073–1081. https://doi.org/10.1145/3583131.3590375 arXiv:2305.11681 [cs].

[7] Jerome H. Friedman. 2001. Greedy Function Approximation: A Gradient Boosting Machine. *The Annals of Statistics* 29, 5 (2001), 1189–1232. http://www.jstor.org/stable/2699986

[8] Alina Geiger, Dominik Sobania, and Franz Rothlauf. 2023. Down-Sampled Epsilon-Lexicase Selection for Real-World Symbolic Regression Problems. 1109–1117. https://doi.org/10.1145/3583131.3590400

[9] Thomas Helmuth, Nicholas Freitag McPhee, and Lee Spector. 2016. The Impact of Hyperselection on Lexicase Selection. In *Proceedings of the Genetic and Evolutionary Computation Conference 2016* (Denver, Colorado, USA) *(GECCO '16)*. Association for Computing Machinery, New York, NY, USA, 717–724. https://doi.org/10.1145/2908812.2908851

[10] Jose Guadalupe Hernandez, Alexander Lalejini, Emily Dolson, and Charles Ofria. 2019. Random subsampling improves performance in lexicase selection. In *Proceedings of the Genetic and Evolutionary Computation Conference Companion* (Prague, Czech Republic) *(GECCO '19)*. Association for Computing Machinery, New York, NY, USA, 2028–2031. https://doi.org/10.1145/3319619.3326900

[11] Arthur E. Hoerl and Robert W. Kennard. 1970. Ridge Regression: Biased Estimation for Nonorthogonal Problems. *Technometrics* 12, 1 (1970), 55–67. http://www.jstor.org/stable/1267351

[12] Hassan Ismail Fawaz, Germain Forestier, Jonathan Weber, Lhassane Idoumghar, and Pierre-Alain Muller. 2019. Deep learning for time series classification: a review. *Data Mining and Knowledge Discovery* 33, 4 (2019), 917–963.

[13] Michael Kommenda, Bogdan Burlacu, Gabriel Kronberger, and Michael Affenzeller. 2019. Parameter identification for symbolic regression using nonlinear least squares. *Genetic Programming and Evolvable Machines* 21, 3 (Dec. 2019), 471–501. https://doi.org/10.1007/s10710-019-09371-3

[14] John R Koza. 1992. *Genetic Programming: On the Means of Programming Computers by Means of Natural Selection*. MIT Press.

[15] John R Koza. 1994. Genetic programming as a means for programming computers by natural selection. *Statistics and computing* 4 (1994), 87–112.

[16] Krzysztof Krawiec and Una-May O'Reilly. 2014. Behavioral programming: a broader and more detailed take on semantic GP. In *Proceedings of the 2014 Annual Conference on Genetic and Evolutionary Computation* (Vancouver, BC, Canada) *(GECCO '14)*. Association for Computing Machinery, New York, NY, USA, 935–942. https://doi.org/10.1145/2576768.2598288

[17] William La Cava, Thomas Helmuth, Lee Spector, and Jason H. Moore. 2019. A Probabilistic and Multi-Objective Analysis of Lexicase Selection and $\epsilon$-Lexicase Selection. *Evolutionary Computation* 27, 3 (09 2019), 377–402. https://doi.org/10.1162/evco_a_00224 arXiv:https://direct.mit.edu/evco/article-pdf/27/3/377/1858632/evco_a_00224.pdf

[18] William La Cava, Patryk Orzechowski, Bogdan Burlacu, Fabricio de Franca, Marco Virgolin, Ying Jin, Michael Kommenda, and Jason Moore. 2021. Contemporary Symbolic Regression Methods and their Relative Performance. In *Proceedings of the Neural Information Processing Systems Track on Datasets and Benchmarks*, J. Vanschoren and S. Yeung (Eds.), Vol. 1. Curran. https://datasets-benchmarks-proceedings.neurips.cc/paper_files/paper/2021/file/c0c7c76d30bd3dcaefc96f40275bdc0a-Paper-round1.pdf

[19] William La Cava, Tilak Raj Singh, James Taggart, Srinivas Suri, and Jason H. Moore. 2019. Learning concise representations for regression by evolving networks of trees. http://arxiv.org/abs/1807.00981 arXiv:1807.00981 [cs].

[20] William La Cava, Lee Spector, and Kourosh Danai. 2016. Epsilon-Lexicase Selection for Regression. In *Proceedings of the Genetic and Evolutionary Computation Conference 2016*. 741–748. https://doi.org/10.1145/2908812.2908898 arXiv:1905.13266 [cs].

[21] William G. La Cava, Paul C. Lee, Imran Ajmal, Xiruo Ding, Priyanka Solanki, Jordana B. Cohen, Jason H. Moore, and Daniel S. Herman. 2023. A flexible symbolic regression method for constructing interpretable clinical prediction models. *npj Digital Medicine* 6, 1 (June 2023), 107. https://doi.org/10.1038/s41746-023-00833-8

[22] Brad L Miller, David E Goldberg, et al. 1995. Genetic algorithms, tournament selection, and the effects of noise. *Complex systems* 9, 3 (1995), 193–212.

[23] T. Pham-Gia and T.L. Hung. 2001. The mean and median absolute deviations. *Mathematical and Computer Modelling* 34, 7 (2001), 921–936. https://doi.org/10.1016/S0895-7177(01)00109-1

[24] Joseph D Romano, Trang T Le, William La Cava, John T Gregg, Daniel J Goldberg, Praneel Chakraborty, Natasha L Ray, Daniel Himmelstein, Weixuan Fu, and Jason H Moore. 2021. PMLB v1.0: an open-source dataset collection for benchmarking machine learning methods. *Bioinformatics* 38, 3 (Oct. 2021), 878–880. https://doi.org/10.1093/bioinformatics/btab727

[25] Lee Spector. 2012. Assessment of problem modality by differential performance of lexicase selection in genetic programming: a preliminary report. In *Proceedings of the 14th Annual Conference Companion on Genetic and Evolutionary Computation* (Philadelphia, Pennsylvania, USA) *(GECCO '12)*. Association for Computing Machinery, New York, NY, USA, 401–408. https://doi.org/10.1145/2330784.2330846

[26] Marco Virgolin, Tanja Alderliesten, and Peter A. N. Bosman. 2019. Linear scaling with and within semantic backpropagation-based genetic programming for symbolic regression. In *Proceedings of the Genetic and Evolutionary Computation Conference* (Prague, Czech Republic) *(GECCO '19)*. Association for Computing Machinery, New York, NY, USA, 1084–1092. https://doi.org/10.1145/3321707.3321758