

Convergence of distributed symbolic regression using metaheuristics

Ben Cardoen
ben.cardoen@student.uantwerpen.be

May 8, 2017

1 Abstract

Symbolic regression (SR) fits a symbolic expression to a set of expected values. Amongst its advantages over other techniques is the ability for a practitioner to interpret the resulting expression, determine important features by their usage in the expression, and an insight into the behavior of the resulting model (e.g. continuity, derivatives, extrema). SR combines a discrete combinatoric problem (combining base functions) with a continuous optimization problem (selecting and mutating constants). One of the main algorithms used in SR is Genetic Programming (GP). The convergence characteristics of SR using GP is still an open issue. The continuous aspect of the problem has traditionally been an issue in GP based symbolic regression. This paper will study convergence of a GP-SR implementation on selected use cases known for bad convergence. We introduce modifications to the classical mutation and crossover operators and observe their effects on convergence. The constant optimization problem is studied using a two phase approach. We apply a variation on constant folding in the GP algorithm and evaluate its effects. The hybridization of GP with 3 metaheuristics (Differential Evolution, Artificial Bee Colony, Particle Swarm Optimization) is studied.

2 Introduction

2.1 Overview

2.2 Symbolic Regression

2.2.1 Compared to other techniques

2.2.2 Interpreting results

2.2.3 Applications

2.3 Convergence

2.4 Metaheuristics

2.5 Genetic Programming

2.6 Constant optimization problem

2.7 Distributed

3 Design

In this section we will detail the design of our tool, the algorithm and its parameters.

3.1 Algorithm

3.1.1 Input and output

The algorithm accepts a matrix $X = n \times k$ of input data, and a vector $Y = 1 \times k$ of expected data. It will evolve expressions that result, when evaluated on X , in an $1 \times k$ vector Y' that approximates Y . N is the number of features, or parameters, the expression can use. K is the number of datapoints. The algorithm makes no assumptions on the domain or range of the expressions or data set. While domain specific knowledge can be of great value, in real world situations such data is not always known. We aim to make a tool that operates under this uncertainty.

3.1.3 Entities

We will describe briefly the important entities in our design. In Figure 2 the architecture of our tool is shown. Functionality such as IO, statistics and plotting is left out here for clarity. The interested reader is referred to the source code documentation.

Algorithm The algorithm hierarchy extends each instance with new behavior, only overriding those functions needed. Code reuse is maximized here by the usage of hook functions, which in the superclasses result in a noop operation. An example of this is the evolve function, by and large the most complex function of the algorithm. By

3.1.2 Control flow

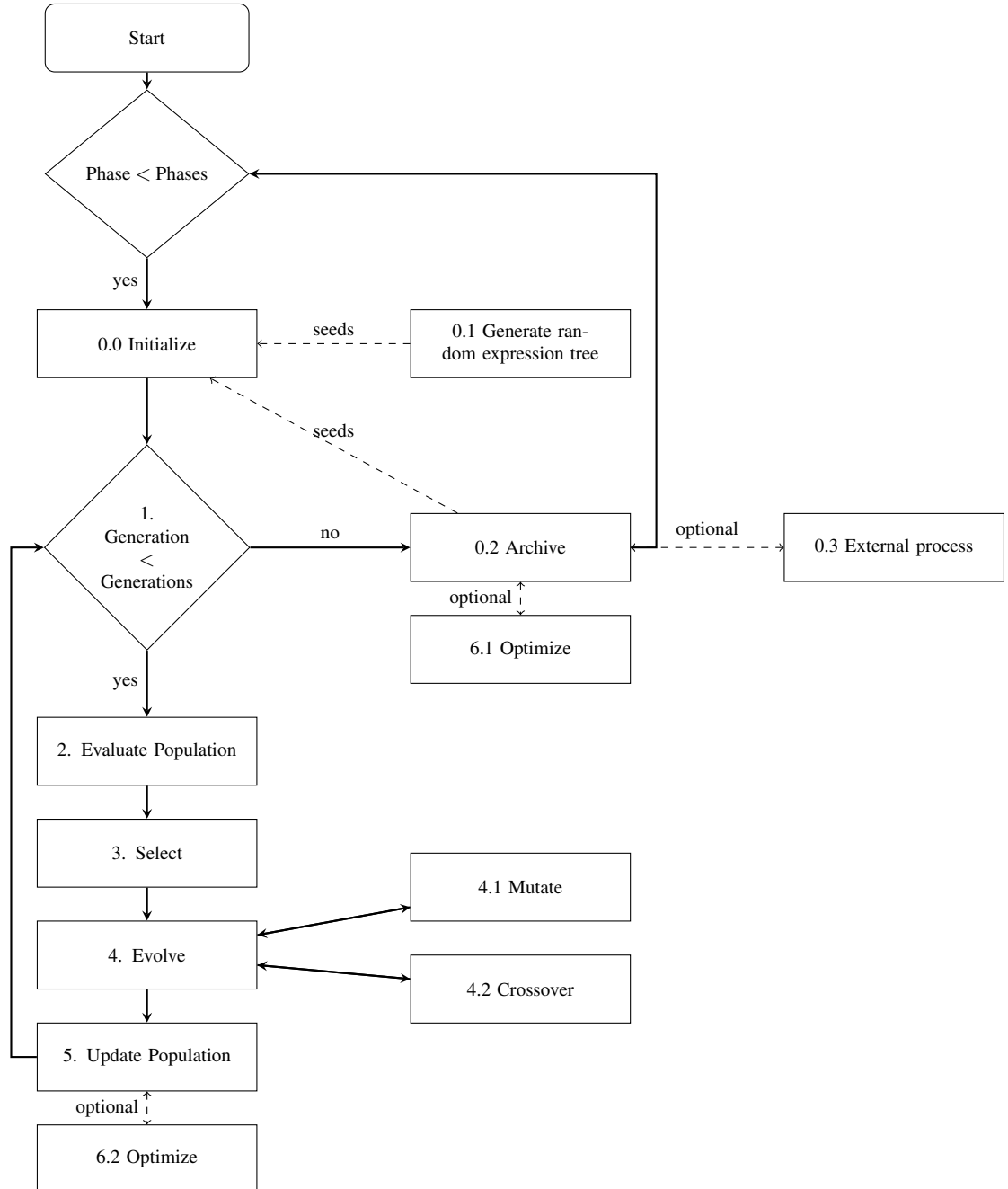


Figure 1: CSRM control flow.

using hooks such as *requireMutation()* this function can remain in the superclass, and the subclass need only implement the hook.

Tree The Tree and Node classes represent the main data structure of CSRM’s population. The Tree is comprised of Node objects, each of which holds an optional constant weight. Leaves are either features (VariableNode) or constants. In 3.8 we go in depth into the design choices made for this representation. The Tree class holds a few utility functions, for example to create trees from an expression using a parser in the tool module. Functionality needed by the mutation and crossover operators is present here as well. Of note are both construction functions, which allow for efficient generation of random trees. In section 3.3.1 we will cover the problem these functions solve.

Optimizer CSRM provides an interface to continuous optimizers. In section 5 we will go deeper into the specific algorithms used. The population data structure is shared between the subclasses. The instances in the population can be subclassed if specific functionality is needed by an optimizer. This approach allows us to reuse a large part of the code between algorithms. For convenience a stub optimizer is added. The optimization algorithms share configuration parameters, and encapsulate those that are distinct to their specific nature.

3.1.4 Implementation

Our tool is implemented in Python. The language offers portability, access to rich libraries and fast development cycles. The disadvantages are speed and memory usage compared with compiled languages (e.g. C++) or newer scripting languages (e.g. Julia). Furthermore, Python’s usage of a global interpreter lock makes shared memory parallelism not feasible. Distributed programming is possible using MPI.

3.2 Fitness

In this work we interpret optimization as a fitness minimization process. In the following we will discuss the fitness function and its behavior.

3.2.1 Distance function

The goal of the algorithm is to find f^* such that

$$Y = f(X)$$

$$Y' = f'(X)$$

$$\text{dist}(Y, Y') = e$$

results in e minimal.

Not all distance functions are equally well suited for this purpose. A simple root mean squared error (RMSE) function has the issue of scale, the range of this function is $[0, +\infty)$, which makes comparison problematic, especially if we want to combine it

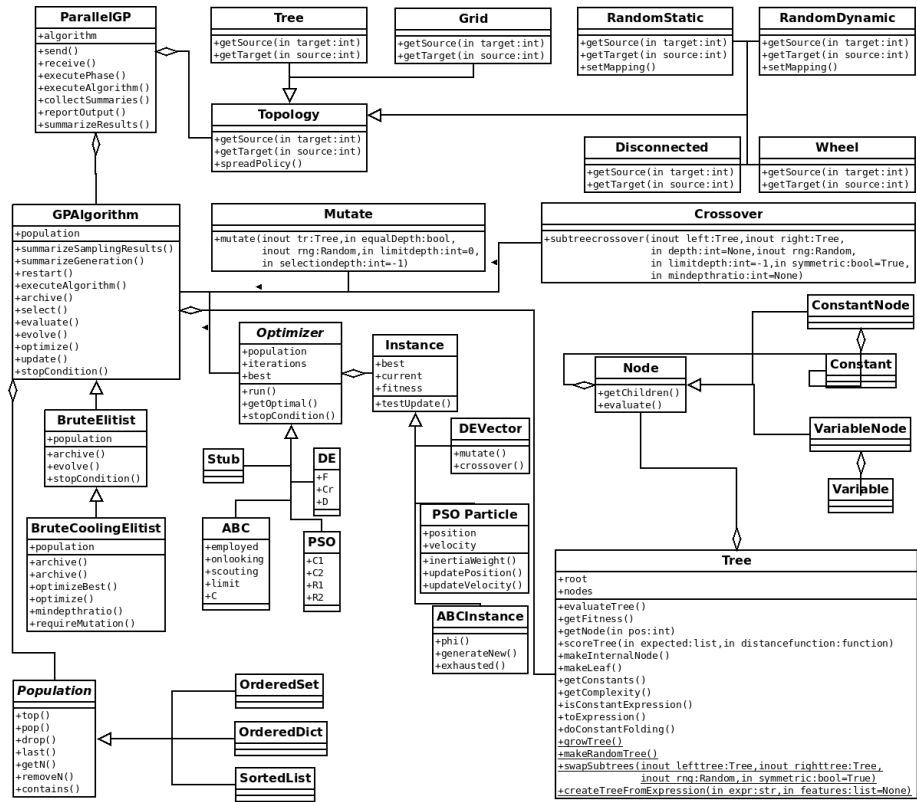


Figure 2: UML of CSRM.

with other objective functions. A simple linear weighted sum requires that all terms use the same scale. Normalization of RMSE is an option, however there is no single recommended approach to obtain this NRMSE.

In this work we use a distance function based on the Pearson Correlation Coefficient. Specifically, we define

$$dist_p(Y, Y') = 1 - |r|$$

with

$$r = \frac{\sum_{i=0}^n (y_i - E[Y]) * (y'_i - E[Y'])}{\sqrt{\sum_{j=0}^n (y_j - E[Y])^2 * \sum_{k=0}^n (y'_k - E[Y'])^2}}$$

R has a range of [-1, 1] where 1, -1 indicate linear and negative linear correlation and 0 no correlation. This function has a range [0,1] which facilitates comparison across domains and allows combining it with other objective functions. The function reflects the aim of the algorithm. We not only want to assign a good (i.e. minimal) fitness value to a model that has a minimal distance, we also want to consider linearity between Y and Y'.

3.2.2 Diversity

Diversity, the concept of maintaining semantically different specimens, is an important aspect in metaheuristics. Concepts such as niching and partitioning are often used to promote this behavior, amongst other reasons to prevent premature convergence or even to promote multimodality. Our tool uses a simple measure that mimics the more advanced techniques stated above. It should be clear that for any combination of input and output data, there are a huge set of expressions with an identical fitness value. Such expressions can lead to premature convergence, consider a population of 5 where the 3 best samples all have fitness scores of 0.134. Our tool will aim to prevent retaining expressions that have identical fitness values. In contrast, in some metaheuristics [21] allowing replacement of solutions with identical fitness values (not duplication, but replacement), can actually help avoiding local minima.

3.2.3 Predictive behavior

The algorithm evaluates expressions based on training data X_t . X_t is an $n \times rk$ matrix based on the original input matrix X ($n \times k$). R is the sampling ratio, the ratio between training and test data. After completion of the algorithm the population is ordered based on minimized fitness values calculated on the training data. In real world applications practitioners are also interested in the predictive qualities of the generated expression. In other words, how well do the expressions score on unknown data. In the final evaluation we score each expression on the full data to obtain this measure. While this gives us some information on how good an expression is on the full data set, we are also interested in how the convergence to this value progresses. If we add 10 more generations, or increase the population by a factor 1.5, what do we gain or lose in predictive quality of the expressions? To define this we use a correlation measure between the fitness values using the training data and those of the full data. This measure quantifies the predictive value of the final results. Finally, we calculate a correlation

trend between the training fitness values at the end of each phase, and the final fitness values calculated on the full data. This trend describes the convergence process of the algorithm over time, specifically directed at the predictive value of the solutions found. It is important to note that the entire final population is considered in these calculations, not only the best. While ideal, there is no guarantee that the expression that has the lowest fitness value on the training data will score best on the full data set. We would like these measures to give us a descriptive capability for the entire process.

3.2.4 Convergence limit

As a stopcondition our tool uses a preset number of iterations. The user specifies the number of generations (g) and phases (r), and the algorithm will at most execute $g \times r$ iterations. Convergence stalling is implemented by keeping track of the number of successful operations (mutation or crossover). If this value crosses a threshold value convergence has stalled and the algorithm is terminated.

3.3 Initialization

Initialization is done using the 'full' method [17]. The algorithm has a parameter initial depth, each new expression in the population is created using that depth. Unless the maximum depth is equal to the initial depth, the algorithm will quickly vary in depth, evolving an optimal depth.

3.3.1 Invalid expressions

Generating a random expression is done by generating random subtrees and merging them. An important observation here is that randomly generated expressions can be invalid for their domain. The ubiquitous example here is division by zero. Several approaches to solve this problem exist. One can define 'safe' variants of the functions, in case of division by zero, returning a predefined value that will still result in a valid tree. The downside to this approach is that the division function's semantics is altered, a practitioner, given a resulting expression, would have to know that some functions are no longer corresponding entirely to their mathematical equivalents, and what 'safe' values the implementation uses. The other option is assigning maximum fitness to an invalid expression. While simple, this approach needs a careful implementation. From a practical standpoint wrapping functions in exception handling code will quickly deteriorate performance. This last problem is solved by a quick domain check for each calculation, avoiding exceptions.

Invalidity probability We define the probability that a randomly generated tree of depth d , with n possible features, k possible base functions, and j data points as q . With more complex problems d will have to be increased. GP is also susceptible to bloat [8], increasing d even further. This issue will affect generation of subtrees in mutation. With more datapoints the probability of at least one leading to an invalid evaluation will increase. An increase in d will lead to an exponential increase in the number of nodes in the tree. A node can be either a basefunction or a leaf (parameter

or constant). For each additional node the probability q increases. We can conclude that q , while irrelevant for small problems and depths, becomes a major constraint for larger problem statements.

Bottom up versus top down There are two methods to generate a tree based expression : bottom up and top down. The top down approach is conceptually simpler, we select a random node and add randomly selected child nodes until the depth constraint is satisfied. The problem with this approach is that the expression can only be evaluated at completion, early detection of an invalid subtree is impossible. In contrast in a bottom up construction we generate small subtrees, evaluate them and if and only if valid merge them into a larger tree. This allows for early detection of invalid expressions. A downside of this approach is the repeated evaluation of the subtrees, which can be mitigated by caching the result of the subtree.

Disallowing invalid expressions in initialization We can generate random expressions and let the evolution stage of the algorithm filter them out based on their fitness value, or we can enforce the constraint that no invalid expressions are introduced. The last option is computationally more expensive at first sight, since the algorithm is capable by definition of eliminating unfit expressions from the population. This can lead to unwanted behavior in the algorithm itself. For high q values we can have a significant part of the population that is at any one time invalid. This can lead to premature convergence, similar to a scenario where the population is artificially small or dominated by a set of highly fit individuals. Another observation to make is that the algorithm will waste operations (mutation, crossover) on expressions that are improbable to contribute to a good solution. While more expensive computationally, we therefore prohibit generation of invalid expressions in the initialization.

3.4 Evaluation and cost

Evaluating a tree requires a traversal, once for each datapoint. In order to compare optimization algorithms across implementations practitioners can use as a measure the number of evaluations required to reach a certain fitness level. If this measure is used, one should take into account that not all evaluations are equal. With the trees varying in depth and density the evaluation cost varies significantly. Furthermore, evaluating " $1 + 2$ " is computationally far less expensive than " $\log(3, 4)$ ". A simple count of evaluation functions executed does not really reflect the true computation cost, especially when we consider the variable depth of the trees. Increasing the depth leads to an exponential increase in the number of nodes, and thus in the evaluation cost. Our tool uses a complexity measure that takes into account the density of the tree and which functions are used. A tree comprised of complex functions will score higher in cost than a corresponding tree using simple multiplications and additions. Although this is an option, we do not use this measure in the objective function. We would like to observe the effect of the cost, but not directly influence the algorithm. There is no direct link between more complex functions and an optimal solution. In certain domains the argument can

be made that more complex functions are more likely to lead to overfitting, or more likely to lead to invalid trees due to smaller domain.

3.5 Evolution

In the evolution stage we apply two operators on (a selection of) the population. The operators are configured to constrain the depth of the modified trees, enforcing the maximum depth parameter. We trace the application of the operators during the execution using an effectiveness measure. Each time an operator application is able to lower the fitness of a tree, this measure increases. Using this measure we can gain insight when and how certain optimizations and modifications work inside the algorithm instead of simply observing the algorithm as a black box.

3.5.1 Mutation

Mutation works by replacing a randomly chosen subexpression with a newly generated. In our implementation this means selecting a node in the tree and replacing it with a new subtree. Our mutation operator can be configured to replace the subtree with one of equal depth or a randomly chosen depth. The insertion point can be made depth-sensitive. A shallow node, a node with a low depth, is the root of a subtree with significant depth. Replacing such a subtree is therefore more likely to have a significant effect on the fitness value. Replacing a deep node will have on average a smaller effect. We will investigate this assumption in our experiments. In our implementation the mutation operator is applied using a cooling schedule, or on the entire population. The cooling schedule uses the current fitness and the current generation of a tree to decide if the tree is likely to benefit from mutation. This is similar to the approach in simulated annealing. Mutation introduces new information into the optimization process (a new subtree). This process can be constructive (lowering fitness) or destructive (increasing fitness). The idea behind the cooling schedule is that for fitter trees the mutation operation will not be able to improve fitness (decrease it), while for less fit trees it has a higher probability. This probability is estimated using the current generation and fitness value (relative to the population), and using this information a random choice is made whether or not to apply mutation. In the experiments section we will investigate if this assumption holds, and if it leads to gains in fitness and or effectiveness. The mutation operator has to generate new subtrees, and therefore the same issues seen in the initialization process which we discussed in section 3.3.1 apply here as well. The mutation operator will generate subtrees until a valid one is found. If the depth sensitive operation proves to generate equal or better fitness values, this could significantly reduce the computation cost of the operator.

3.5.2 Crossover

Crossover operates on 2 trees in the population. It selects subtrees from both, and swaps them between each other. The resulting set of 4 trees is scored and the 2 fittest replace the original trees. As with mutation crossover can be configured to operate on a set depth, or pick a random depth. It can also work symmetric, picking an equal depth

insertion point in both trees. Crossover in our implementation can be applied in sequence to the entire population, with pairwise mating in order of fitness. Alternatively a random selection can be used. A variant of both, alternating between the two schedules based on a random choice is a third option. This is similar to the roulette wheel selection method, although without replacement. Crossover, unlike mutation, does not introduce new information in the sense that no new subtrees are generated. Crossover can be configured to work with a decreasing depth, operating only on deeper nodes at the later stages of the algorithm. The assumption for this mode of operation is similar to that made for mutation.

3.6 Selection

After evolution a decision has to be made on which expressions will be retained in the population. In our tool the population is fixed, so a replacement is in order. We use an elitist approach. If an expression has a lower (better) fitness value after mutation, it is replaced. In crossover we combine two expressions r , and t , resulting in two offspring s , u . From these four expressions the two with minimal fitness survive to the next generation.

3.7 Archiving

The algorithm holds an archive that functions as memory for best solutions obtained from the best expressions at the end of a phase., from seeding, or from other processes. At the end of a phase we store the j best expressions out of the population. J is a parameter ranging from 1 to the population size n . With $j == n$ we risk premature convergence, with $j == 1$ the risks exists that we lose good expressions from phases which will have to be rediscovered. While there are numerous archiving strategies described in literature, we use a simple elitist approach. This means that there is no guarantee that the best j samples of phase i are retained, if they have fitness values lower than those present in the archive and the archive has no more empty slots, they will be ignored. This leads us to the size of the archive. While no exact optimal value for this exists, in order to function as memory between phases it should be similar in size to the amount of phases. The j parameter will influence this choice as well.

3.8 Representation and data structures

3.8.1 Expression

We use a tree representation for an arithmetic expression, where each internal node represent a base function (unary or binary), and each leaf either a feature or a constant.

Tree representation We use a hybrid representation of a tree structure. The tree is a set of nodes, where each node holds a list of children nodes. This representation favors recursive traversal algorithms. In addition to this representation, the nodes of a tree are stored in a dictionary keyed on the position of a node. This allows for $O(1)$ access needed by the mutation and crossover operators. The overhead of this extra

representation is minimal, since the keys are integers and only a reference is stored. A list representation would be faster in (indexed) access, but would waste memory for sparse trees. With a mix of unary and binary operators the average nodecount of a tree with depth d is $\frac{2^{d+1}-1+d}{2} = O(2^d)$ resulting in savings on the order of 2^d , with d the depth of the tree. The algorithm has the choice which mode of access to use depending on the usage pattern. As an example, selecting a random node in the tree is $O(1)$, selecting a node with a given depth is equally $O(1)$. Splicing subtrees is equally an $O(1)$ operation.

Base functions The set of base functions is determined by the problem domain. For symbolic regression we use the following set:

+, -, %, /, max, min, abs, tanh, sin, cos, tan, log, exp, power

A problem with this set is the varying domain of each, why may or may not correspond with the domain of the features. A possible extension is to use orthogonal base functions such as Chebyshev polynomials. In our solution the functions listed correspond with their mathematical equivalent, e.g. division by zero is undefined. Other constraints are based on floating point representation (e.g overflow).

Constants Constants are fixed values in leaves in the tree. The representation also allows for multiplicative weights for each node, which can be used to optimize the final solution. In contrast to the base functions with a limited set to choose from, constants have the entire floating point range at their disposal. The probability of selecting a 'right' value is extremely small, and domain information is lacking for these constants, it depends on the entire subtree holding the constant. We select a constant from a small range and allow the algorithm to recombine these constants later to larger values. This limiting of the search space can lead to the algorithm converging faster to a suboptimal solution. The constants are reference objects in the tree, and so can be accessed and modified in place by an optimizer.

Features Features are represented as an object holding a reference to a set of input values, and always occur as leaves. The choice for a random leaf is evenly distributed between constants and features.

3.8.2 Population

The population of trees is kept in a sorted set, where the key is the fitness value and the value is a reference to the tree representing the expression. Sorted datastructures are notably lacking from Python, we use the sortedcontainers [9] module which provides sorted datastructures that are designed with performance in mind. This representation allows for fast access in selecting a (sub)set of fittest trees. It also allows for an $O(1)$ check if we already have an equivalent solution, a tree with a fitness score which we already have in the population. This allows our diversity approach mentioned in 3.2.2. In Figure 2 we see that the actual data structure used as population is interchangeable. If duplicate fitness values are allowed, and membership testing is no longer needed, a sorted list could be used instead.

3.9 Parameters

In this section we briefly list the main parameters of the algorithm and their effects on convergence, runtime and complexity.

3.9.1 Depth

The depth of trees can vary between an initial and maximum value. If we know in advance an optimal solution uses 13 features, the tree should have at least 13 leaves in order to use each feature at least once. This requires a depth of at least 4. In practice the depth will need to be greater due to the use of unary functions, and bloat. Bloat is a known issue in GP [8] where the algorithm, unless constrained by limits or an objective function that penalizes depth, will tend to evolve deep trees without gaining much in fitness. In the worst case entire subtrees can be evolved that do not contribute to the fitness value, mimicking introns in genetics. These can still have a valid purpose, serving as genetic memory. Their disadvantage is also clear: a computational overhead without clear effect on the fitness value. A similar problem arises with the generation of constants. Suppose we would like to generate the constant 3.14, the probability of generating this by a single random call is extremely small. The algorithm will try to build an expression fitting 3.14, for example $1 + (3 * 1) + 28/200$. The problem with this is that this process is highly inefficient. It uses iterations that, if a more efficient approach exists, could be used to optimize the fitness value of the tree further. The constant expression wastes nodes that could be used to improve the tree. One approach is folding such expressions into a single constant, which mitigates some of the effects mentioned but does not prevent such subtrees from forming. Without a solution for this issue, the user would have to take this into account and increase the depth parameter. From 3.3.1 we know that the initialization process and the mutation operator will become more costly exponentially with the increase in depth. A similar argument can be made for the time and space complexity of operating on deeper trees, both increase exponentially.

3.9.2 Population size

The population size is directly related to convergence speed. The catch is the quality of the solution, a very small population will lead to premature convergence, the population is lacking in diversity (or information viewed from a different perspective). A large population on the other leads to a large increase in runtime, unless the operators only work on a subset of the population. The optimal value is problem specific, but values in the range of 20-50 give a good balance for our implementation.

3.9.3 Phases and generations

The execution of the algorithms comprises of g generations and r phases, resulting in a maximum of $g \times r$ iterations. For both parameters a too small value will hinder convergence, while a high value can lead to overfitting. The optimal value is domain specific and dependent on the iterations required to approximate the expected data. This is by virtue of the problem statement unknown. There is a subtle difference between these

parameters. Each phase a reseed of the algorithm is done using the archive. This archive holds the best results from previous phases, external seeds and in a distributed setting the best solutions from other processes. The remainder of the population is initialized with random tree. These trees introduce new information into the optimization process, and while expensive and with a low probability of improving fitness, nonetheless will help avoid premature convergence. If the archive size is less than the population size, new trees will always be introduced. The rate at which the archive fills is dependent on the number of phases and a parameter which determines how many trees are archived. In a distributed setting this process accelerates using the input of other processes. If the archive is full after i phases, and the archive size is equal to the population, further phases will no longer have the benefit of newly generated trees. This does not prevent convergence, but could reduce the convergence rate in some scenarios. Increasing g and r both can lead to overfitting, but in order to decide on r we also have to take into account the archiving strategy. In order to find a good starting value for g one can look at the population size. Diffusion, where the information from each expression is shared with the others, will require at least the same amount of generations as there are expressions, depending on the operators and the individual fitness of the expressions. Concentration, where we only look at maximizing the few existing fittest expressions, requires far less generations, but risks premature convergence.

3.9.4 Samples

The user can provide the algorithm with input data, or can specify a range from which to sample the input data. In addition, the ratio between training and testing data can be adjusted. Care should be taken in tuning this value. A low ratio will increase the probability that evolved solutions are invalid on the test data, while a high ratio will lead to overfitting.

3.9.5 Domain

Domain specific knowledge can significantly reduce the time needed to converge to a solution. In our implementation we assume no domain knowledge. While this increases the complexity of obtaining a good solution, it also makes for a fairer evaluation of the algorithm and optimizations used.

3.10 Incremental support

Our tool supports incremental operation. The user can provide seeds, expressions that are known or assumed to be good initial approximations. The algorithm writes out the best solutions obtained in an identical format. By combining this the user can create a Design Of Experiment (DOE) setup using the algorithm. As a use case, suppose the user wants to apply symbolic regression on the output of an expensive simulation. The simulator has k features or parameters, each with a different domain and j data-points. The user wants insights into the correlation of the parameters. A naive approach would be to generate output data for a full factorial experiment, and feed this into the CSRM tool. For both the simulator and the SR tool the cost of this approach would be

prohibitive. It is likely that some features are even irrelevant, leading to unnecessary computation. Instead we can opt for a combined DOE approach. We start with a subset of features $k' < k$ and datapoints $j' < j$. The simulation results Q of this subset are then given to the CSRM algorithm. It generates a solution, optimized for this subset of the problem. The user then adds more parameters, $k' < k'' < k$ and/or datapoints $j' < j'' < j$, runs the simulator again. The resulting output is given the CSRM tool, with Q as seed. This seed is used as memory of the previous run on the smaller input set. Unless there is no correlation between the incrementing datasets the CSRM tool can use the knowledge gained from the previous run to obtain faster convergence for this new dataset. In addition, by inspecting Q the user can already analyze a (partial) result regarding the initial parameter set. Suppose $k = 5$, and only 3 parameters are used in Q . Then the user can exclude the missing two parameters from the remainder of the experiment, as these are unlikely to contribute to the output. By chaining the simulator and CSRM tool together in such a way, an efficient DOE approach can potentially save both simulation and regression time, or result in increased quality of solutions. The advantages are clear :

- The SR tool can start in parallel with the simulator, instead of having to wait until the simulator has completed all configurations.
- Reducing irrelevant parameters detected by the SR tool can save in configuration size for the simulator.
- The SR tool can reuse previous results as seeds, and tackles and slowly increasing search space with those initial values instead of starting without knowledge in a far larger search space.

Possible disadvantages are :

- The SR tool is not guaranteed to return the optimal solution, it is possible the process is misguided by suboptimal solutions. This risk exist as well in the full approach.
- Interpretation of the results can be needed in order to make the decision to prune features.

We will investigate this approach in our use case in section 7. The following diagram illustrates a DOE hypercube design using our tool and a simulator.

3.11 Statistics and visualization

Stochastic algorithms are notoriously hard to debug. By virtue of the problem statement we do not know whether the returned solution is what is expected. The size of the search space makes detecting all edge cases infeasible. In addition the algorithm functions as a black box, where output is presented to the user without a clear trace indicating how or why this output was obtained. Both for developer and practitioner insight into the algorithms inner workings is vital. Our tool provides a wealth of

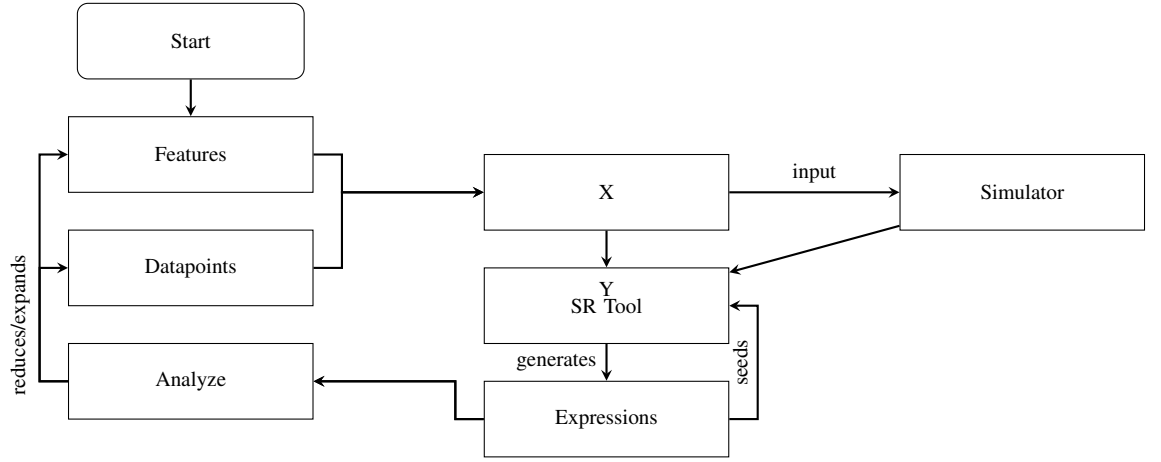


Figure 3: Incremental DOE using CSRM and a simulator.

runtime statistics ranging from fitness values per expression for each point in the execution, convergence behavior over time, depth and complexity of the expressions, cost of evaluations, effectiveness of operators, correlation between training and test fitness values and more. These statistics can be saved to disk for later analysis, or displayed in plots in a browser. Using this information the user can tune the parameters of the algorithm (e.g. overfitting due to an excessively large number of phases), the developer can look at how effective new modifications are (e.g. new mutation operator) and so on. It is even possible to trace the entire run of the algorithm step by step by saving the expressions in tree form, displayed in an SVG image rendered by Graphviz [7]. In Figure 4 a selection of the collected statistics on a testfunction is shown. The third of our set of testproblems was used with a depth $\in [4,10]$, 30 generations, populationsize 30, and 4 phases.

3.12 Conclusion

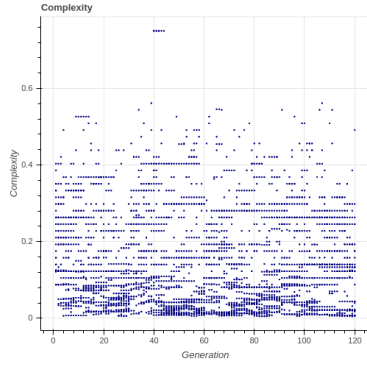
In this section we have covered in detail the design of the CSRM tool, highlighting the choices made.

4 Distributed SR

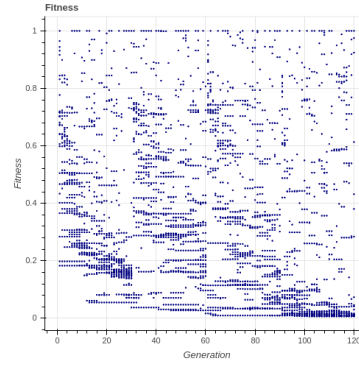
In this section we will cover the parallelization of symbolic regression, and explain our design choices.

4.1 Approaches

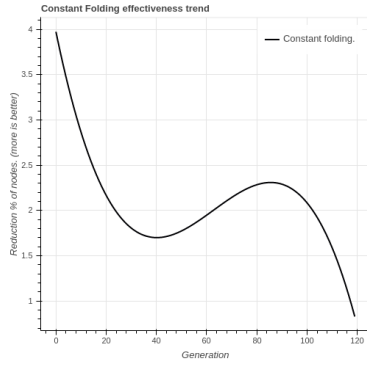
We distinguish between fine grained parallelism and course grained parallelism.



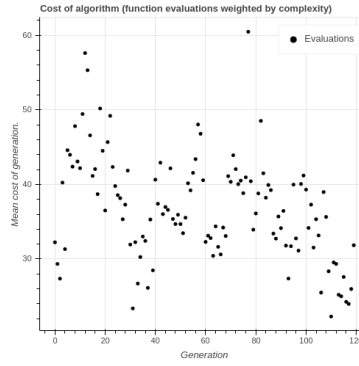
(a) Scaled complexity over generations.



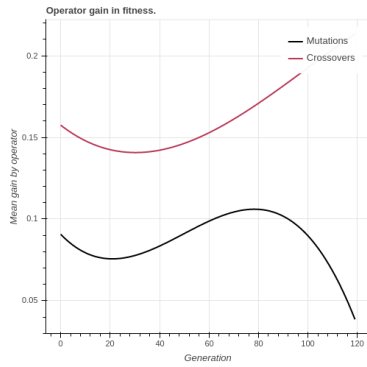
(b) Fitness values over generations.



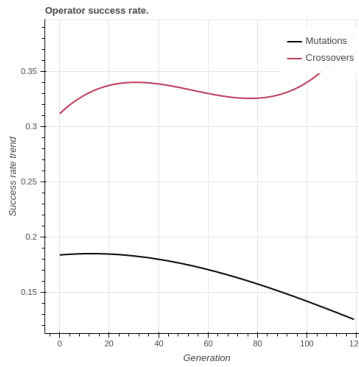
(c) Constant folding savings.



(d) Mean evaluation cost.



(e) Operator gain.



(f) Operator success rate.

Figure 4: Selection of visualizations generated by CSRM.

Fine grained parallelism In fine grained we parallelize a single step in the algorithm, where we execute a number of tasks in parallel that are independent of each other and have a relative short completion time. The evaluation of the fitness function is a prime example of such a workload. The fitness evaluation itself can be quite complex, but in comparison with the entire runtime of the algorithm the computational cost of a single calculation is small. In order to efficiently parallelize such a taskset we have employ a mechanism that has minimum overhead. Shared memory parallelism using threads is ideal for this use case. Optimizations such as dynamically allocated threadpools will reduce the overhead even further, and increase the speedup. Overhead in this case is split over the implementation overhead of starting or assigning a thread and administering it, and the copying of the problem data. Shared memory implies a near zero cost in copying, only a reference is shared. A lock on the shared data is not needed, since the instance the thread operates on is not used by any other part of the program while we execute the fitness function.

Course grained parallelism In course grained parallelism we run a series of tasks in parallel that have a long runtime or complex workload. In our setting an example of course grained parallelism is running the entire algorithm in parallel, with several instances tackling a distinct subset of the problem as separate processes. With a longer runtime of the task, the overhead of parallelizing the problem can be significantly greater without impacting the speedup obtained. Typically we only start and stop such a task once, in contrast with fine grained task which are constantly started and then stopped again. We can use processes or threads for this form of parallelism. Threads have the benefit of being lighter in comparison to processes in terms of overhead. If we allow communication between the tasks threads can elide copying the data, at the cost of introducing locks. Processes typically have to copy data in order to communicate. We can compare both approaches with an interrupt based approach versus message passing. Whichever we choose, communication between tasks requires synchronization of some form. This introduces not only memory and time overhead, but also significantly complicates the implementation. Invariants that hold in sequential implementations are no longer guaranteed in parallel, and a careful implementation is required.

Our approach

Evaluating fine grained parallelism Python has poor support for shared memory parallelism. While threading is available, it will not offer a speedup except for IO bound tasks due to the presence of global lock in the interpreter (GIL). It is possible to use compiled extensions in C to work around this issue, at a high cost in code complexity. A prototype implementation that parallelized the fitness function calculation proved that both threading and processes are unable to offer any speedup, often even running slower than sequential. A large part of this cost is due the overhead of copying Python objects, which are reference counted and thus require a graph traversal in order to create full copies. In order to evaluate a population of n expression tree of depth between 5 and 10 in parallel, we have to first copy the n expressions to each of the n

processes, then evaluate the tree, and then copy the result back. The copying operation is far more expensive than the evaluation function, and both scale exponentially with the depth of the tree. A prototype using threads elided the copy, but the GIL ensured that performance was lower than the sequential approach. This ruled out fine grained parallelism in CSRM.

Evaluating course grained parallelism We create k processes, and give each an instance of the algorithm to execute in parallel. Communication is done via message passing. We use the MPI framework, which offers Python bindings, to allow processes to communicate and synchronize. Copying is still costly, but now a speedup is possible. Threading did not work in this approach due to the GIL. CSRM uses a distributed set of communicating processes in order to solve the SR problem.

4.2 Distributed SR

We will now discuss the benefits of the distributed (course grained) application to our problem. A parallel metaheuristic has the obvious advantage of speedup in comparison with a sequential implementation. By dividing the problem over k processes we can, in ideal circumstances, obtain a speedup linear in k . With speedup we then refer to the time, or number of evaluations, needed to reach a certain threshold of fitness. The advantage of parallelization do not stop with this speedup. We can view the parallel processes as a swarm in itself, where each instance communicates with the others using a predefined pattern. Instead of k standalone processes we now have a set of k cooperating processes. Each process can now use the information of others to improve its own search. Using this approach a superlinear speedup is possible. There are, however, downsides. Communication implies overhead, not only in memory needed but also in synchronization. Without communication a linear speedup is only dependent on the implementation. With communication we introduce time constraints which can bottleneck a subset of the processes, or in the worst case serialize the entire group. We would like to exchange the most valuable information with the least amount of overhead. This balance is problem specific, the cost of copying depends on what exactly is being copied when and to whom it is sent. Even if we restrict us to 1-1 link between two processes, and only exchange the fittest expressions between the two processes, we do not know in advance how large the expression we copy will be, as the depth and sparseness of the tree representing the expression will vary. While each process is given an equal sized subproblem, there are no guarantees that the actual workload of the different processes will be equal. We know that the evaluation of the fitness function has variable computational load. The stochastic nature of the metaheuristics used in the SR implementation compound this issue. By virtue of the problem statement we do not know the optimal solution to our problem, and with different starting points convergence between the different processes is likely to differ. We will address each of these issues.

4.3 Topology

A topology in this context is the set of communication links between the processes. The topology determines the convergence behavior of the entire group. In a disconnected topology, there is no diffusion of information. If a process discovers a highly fit expression, that knowledge has to be rediscovered by the other processes. The only edge case where this is an advantage is if we see the group of processes as a multiple restart version of the sequential algorithm. If the risk for premature convergence due to local optima is high, we can try to avoid those optima by executing the algorithm in k instances, without communication. The implementation should offer the process an efficient way to lookup both processes from which it will receive information (sources) and processes it has to send to (targets). Source lookup is needed in order to solve the synchronization problem. Any topology that contains a cycle between processes can introduce a potential serialization effect at runtime.

Diffusion and concentration In general would like the processes to share information in order to accelerate the search process. With a topology we introduce two concepts : concentration and diffusion. Concentration refers to processes that share little or no information, and focus on their own subset of the search space. Like exploitation concentration is sensitive to overfitting and premature convergence. It is warranted when the process is nearing the global optimum. Diffusion, in this work, is the process of sharing optimal solutions with other processes. Diffusion accelerates the optimization process of the entire group. It is not without risk, however. If diffusion is instant, a single suboptimal solution can dominate the other processes, leading to premature convergence. The topology will determine the synchronization between processes, and the balance between diffusion and concentration.

Synchronization Synchronization between the processes will play an important role. While it won't directly influence convergence, it will constrain the runtime performance of the entire group. If we denote S_i and T_i as the set of processes that are sources and targets respectively for process i , we would like S_i minimal. The message processing code will have to wait for the slowest processes in S_i before continuing. Without asynchronous communication i would have to wait even longer, with the slowest process blocking the receipt of messages from the other processes. Synchronization implies that i cannot send until its receiving stage has completed. If $S_i \cap T_i \neq \emptyset$ we have a cyclic dependency which can introduce deadlock in the implementation. By extension, if there is a cycle in the topology between any i, j processes this deadlock is a real risk. Dealing with this in the communication code is non trivial. We will show in section 4.4 how CSRМ is able to deal efficiently with cycles in the topologies that it implements. Even with deadlock solved or prevented, cycles will introduce a tendency for the the processes to serialize on the slowest process in the node. The time spent waiting in the communication stage is lost to the convergence process. We will show how this is mitigated in our implementation.

4.3.1 Grid

This topology arranges a set of k processes in a square 2D grid. Each process is connected with the neighbouring processes, respectively up, down, left and right. Some variations include diagonal links, or create 3 or higher dimensional meshes. The general idea behind the grid is unchanged in those configurations. A grid connects all processes, allowing for diffusion of the best solution to each individual process. The key observation here is that diffusion is gradual. While a process has an immediate neighbourhood, reaching all processes takes a variable amount of communication links. Diffusion of a dominating solution will take time, and if the solution is suboptimal this time allows the other processes to evolve their own optimal solutions thereby preventing premature convergence. This risk is only mitigated by gradual diffusion, not eliminated. Elimination is only possible in the extreme case by an absence of communication or in a more advanced configuration the usage of cliques. Nodes on the borders of the grid communicate with their counterparts at the mirrored side of the grid. This ensures that the communication process is symmetric. CSRM supports both square grids, where k is square of natural number, or incomplete squares. If $\sqrt{k} \neq n$ for some $n \in \mathbb{N}$ we create a grid that would fit j processes where $\min(j) > k \forall j \in \mathbb{N} \wedge \sqrt{j} = n$. This grid is filled row by row with k processes. This configuration should be avoided, as the communication pattern will not always be symmetrical. It is implemented to allow for a fair comparison with other topologies where the processcount is not a square.

Cost The communication cost for k processes in a single iteration, with m messages sent per exchange, is in our configuration $4km$. The synchronization constraints are high, all processes are interdependent (directly or at most \sqrt{k} links removed). Lookup of targets is statically determined, at configuration each parallel process is given a simple integer indexed mapping, and the symmetry of the mapping simplifies the lookup code.

4.3.2 Wheel

A wheel, or circle topology connects all processes with a single link shared between each source and target. Diffusion is slower than with a grid, with k processes it takes $k-1$ messages to reach all processes. Some variations introduce a 'hub', with spokes reaching out the circle itself, completing the wheel analogy. CSRM implements this topology as a simple circle. A spoked wheel topology, if the hub has bidirectional communication with all processes, has a maximum distance of 2 between all processes offering fast diffusion. Without a hub but bidirectional links the maximum distance is $k/2$.

Cost For k processes with m messages sent per exchange, the communication cost is km . This is a static configuration, with symmetric lookup. A circle forms a cycle, resulting in high synchronization constraints. The variant with hub and spokes introduces even higher synchronization constraints. At a doubling of message cost the doubly linked circle has an significantly higher synchronization cost than the singly linked variant.

4.3.3 Random

In a random topology a process picks a random target for its messages. We can either configure it to do so statically, such that the target remains invariant at runtime, or select a new target each iteration. The number of targets is variable as well. CSRM implements all variants, allowing for both dynamic and static random topologies with a variable amount of targets per sending process. The idea behind a random topology is avoiding patterns that are present in the structured topologies. If such a pattern leads to premature convergence, poor synchronization or fails to gain from the exchange of information between the processes there exists a possibility that a random approach can work. The downside is that we do not know what the maximum distance is between two processes, or even if they are connected. There can be cliques in the topology, or cycles. Avoiding these, or detecting them requires more complex code than a simple random assignment of targets. By increasing the number of targets we decrease the possibility of cliques, and will minimize eventually the distance between two processes, at the cost of increased overhead. CSRM does not enforce constraints such as clique or cycle forming in its random topologies.

Cost The cost for a singly linked random static topology of k processes with m processed per exchange is clearly km . Synchronization constraints are unknown, and would have to be resolved by the processes. Cyclic dependencies between the processes are likely. The lookup code for a static configuration is simple, symmetric and fast. For a dynamic configuration the lookup is simple, but only symmetric at single points in time. This observation is important because the virtual time of a process will diverge from that of the other processes. By allowing a dynamic configuration we introduce a new type of synchronization constraint. In CSRM processes have a copy of the global topology, but this is no longer immutable. In order to resolve sources the process has to know the virtual time of the other processes. The topology remains deterministic, the seed used is identical between each process. Synchronization is highly complex in a random dynamic topology as cycles and cliques vary over time.

4.3.4 Tree

We introduce a binary tree topology. The links are unidirectional, with a single root and a given depth. The processcount should ideally be $k = 2^{d+1} - 1$ for some $d > 0 \in \mathbb{N}$ to create full binary tree. This is not a hard constraint, with values not satisfying we construct a tree of depth $d = \lfloor \log_2 k \rfloor$ and fill the last level left to right until all processes are assigned a position. Communication is unidirectional, from the root to the leaves. This topology is free from cliques or cycles. The leaves act as sinks, while the root is a source. Lookup is fast, static and symmetric. Each process, except the leaves has at most two targets, and one source (except the root). The tree topology offers a structured balance between diffusion and concentration. The distance between two processes range from 1 to $\log_2 k$. In the case of leaves the distance is infinite. Note that depending on the communication strategy, each subtree behaves as a sink, and will not diffuse information. A variant on this topology is a singly linked list, where the maximum distance is k . The problem with this topology is that its diffusion scales

linearly with k .

Cost For k processes and m messages per exchange, we have km messages exchanged. Synchronization overhead is minimal, there are no cycles and a process is influenced at most by $\log_2 k$ other processes and influences at most $k-1$. With the exception of a disconnected topology the tree topology allows for a high speedup compared with grid, random and circle.

4.3.5 Disconnected

The disconnected topology has zero diffusion and an infinite distance between processes. The only applications of this topology are in cases where premature convergence due to diffusion is known in advance and can for some reason not be mitigated, and as a comparison with the other topologies in order to measure memory and synchronization cost.

Cost Cost is near zero, no messages are sent nor is any synchronization needed. In practice the collecting of all results will still have to be done by either an elected process or a process statically assigned as collector. This holds true for all topologies.

4.4 Asynchronous communication

We have to tackle deadlock and synchronization delay in our parallel implementation. We will first describe the interaction between the processes and using this context demonstrate our solution.

Control flow A parallel process in our implementation executes a single phase of the algorithm, then collects at most m of its fittest expressions from its archive. The set of target processes is resolved using the topology, and then the m messages are sent to the targets. After the sending stage the process looks up its sources using the topology, and waits for messages from those sources. After receiving the messages are decoded into solutions which are passed to the algorithm for its next phase.

Waiting for messages If the process does not wait for messages from the other processes the convergence behavior of the entire group becomes extremely hard to predict. In most platforms there are no hard guarantees about inter process scheduling. We can end up in extreme cases with a single process only sending and never consuming messages. This would break the intention of any structured topology. While this is an extreme example, even in average causes an extra level of non determinism is introduced without there being an explicit need for it. In addition, not waiting for messages requires that messages are collected by the parallel framework, and this can lead to internal buffers varying strongly in size. The sending code would have to be asynchronous as well. We will show another approach where the deterministic execution of the parallel metaheuristic is retained. This does not imply that the process itself is no

longer stochastic, only that we are guaranteed that our designed communication pattern is strictly adhered to.

Deadlock From our previous discussion we know that deadlock is a risk if there are cycles in the topology. Since each process has full knowledge of the topology (except in a random dynamic topology) it is possible to write an implementation that solves deadlock by ordering the sending and receiving of the messages. Suppose A waits for B, and B for A in the most simple example. A simple solution of A sending to B, B receiving from A, then B sending to A and A receiving from B would break the deadlock. This simple interleaving solution is not unique, but the priority is not important here. As long as the deadlock is broken the processes can continue. There are 2 serious issues with this approach. First, this approach serializes concurrent processes. This is detrimental to any speedup we were hoping to achieve. Second, the implementation becomes more complex. The order of calls is no longer statically defined, we need an algorithm that acts as a central coordinator between any group of processes in a cycle, computes a solution, and executes that solution in order. We cannot directly invoke operations on other processes, so the only solution is to execute the coordinating algorithm in parallel on all processes if they detect that they are a member of a cycle. Each process will then know when it can send and receive based on the computed order. While this solves deadlock, we still end up with a serialized execution and the communication code becomes far more complex than a simple sequential sending and receiving call.

Asynchronous solution CSRM instead opts for an asynchronous sending of messages in order to resolve deadlock and mitigate the serialization. Cycle detection is no longer needed. A process executes a phase, then looks up its targets for messages. Instead of sending the messages with a blocking call to the framework, the sending process now allocates a buffer and send each set of messages to its target with an asynchronous call. It stores future object that can later be checked to verify if the receiving process has executed its mirror receive call. The order of sending is no longer relevant, the sending calls return immediately for all targets. Next the process calls receive for each of its sources. This is a blocking call, but this does not introduce a risk for deadlock. In the worst case our process waits until the source has stopped its own sending calls, but due to its asynchronous nature this process is quite fast. The order of sources is again no longer a risk for deadlock. After receiving all messages the process continues its execution as before. Each process still has the allocated buffers and callback objects from its send operation. The buffers cannot be deallocated as long as the receiving process has not acknowledged receipt of the contents. In order to check this we invoke the future objects, but care must be taken here. Calling the futures for each receiving process results in a blocking call. If we time this operation incorrectly we simply have deadlock all over again. The latest a process can wait is that point in time where the send buffers are needed again, which is in the next iteration. As soon as a process enters the sending phase, the first operation it executes is waiting on the future objects. After each returning call the corresponding buffer is cleared for reuse. It should be clear that this approach does not simply delay deadlock but prevents it from

occurring. The asynchronous approach allows the processes to interleave in any order that breaks the deadlock. Instead of an explicit algorithm, we simply let each cycle of processes resolve the deadlock in their own optimal sequence. Finally note that the receive call returns immediately once the framework has registered the corresponding send call, it is not directly waiting on the sending process.

Synchronization delay We know from our discussion that the execution time of a single phase will vary per process, and even between phases for a single process. The receiving of messages is a synchronization point between processes, but not a strict one as the sending process will only wait for receipt in the next phase. This means that a process will only be waiting on any other process after completion of a phase. If we introduce drift as the virtual time difference between processes, in order for processes to delay each other the drift would have to be greater than the execution time of a single phase. A strict serialization effect is therefore not possible. What we introduce with this implementation is latency buffer equal to the runtime of the next phase. The time to execute a phase, while variable, will still be on average far greater than the time needed to communicate results. One of our aims in distributing SR is obtaining a speedup by, amongst other things, finding a balance between phase runtime and communication time.

Future improvements If we allocate a buffer per communication phase, we can wait even longer before invoking future objects. This means we can leave a configurable amount of phases before we wait on the receipt of messages. The downside of this approach is that the memory constraints for each process now significantly increase from m to pm where p is the number of phases we opt to wait. The receive call is executed in sequence for the list of sources in the target process. In order to minimize delay even further we could execute this call asynchronously as well. Extra code would need to be introduced to handle the possible non deterministic behavior this introduces. The order in which the messages are received can influence the algorithm. When the archive is seeded with these external expressions we drop any expressions with identical fitness values. There is thus a risk that by varying the order we introduce non determinism in our results. This can be resolved in the receiving code by preallocating the buffer for all sources and after all asynchronous receive calls storing the messages in the same order of the sources list. Asynchronous receiving is not implemented in CSRM.

4.5 Communication strategies

If the process has j targets, and the algorithm is configured to distribute the m best solutions, we can use several approaches to send those to their targets. We can spread m over j , using a random, interleaving, slicing or any other type of sampling technique. An alternative is copying m j times so that each target receives m messages. CSRM has a `Spreadpolicy` interface that hides the implementation of this behavior for the processes. CSRM defaults to a slicing policy. If we have m messages to distribute over j targets, we will sequentially assign each of the j process $\lfloor m/j \rfloor$ messages. This policy is efficient as it avoids copying but has a direct effect on the diffusion between

the processes. First of all the value of m should be chosen with the topology in mind. If $m = 2$ in a grid topology, each process will only send and receive 2 messages, not 4. This will alter the diffusion pattern. For a tree topology m should be 2 as well, else we create a symmetric series of cliques, which, while elegant, will unlikely be intended. If we use a copy policy the value of m can be as low as 1. For random topologies m depends on the parameter determining the number of targets. In our circle topology $m = 1$ is sufficient for both policies.

4.6 K fold sampling

4.7 Conclusion

We have covered our distributed design and implementation and highlighted the issues that drove our choices. With a wealth of topologies and configurations we also increase the dimensions of the parameter optimization problem. In the experiments we will evaluate some, but not all parameter choices. Our implementation can easily be extended with more topologies or policies without risking deadlock or serialization.

5 Constant optimization

5.1 Constant Optimization

Selecting base functions is a combinatoric, discrete problem. Selecting the right constants to use is a continuous problem, that GP tries to solve with a combinatoric approach. There are several aspects to this issue, we will go over each individually and show our approach.

5.1.1 Initialization revisited

During initialization it is possible that a generated tree represents a constant expression. Such an expression is only a valid approximation if no feature has an influence on Y , which is an unlikely edge case. A constant expression is of no use to the algorithm, but without measures to prevent or remove these they will still be formed. Ideally such an expression will have a worse fitness value than non constant expressions, and will be eventually filtered out. Since detecting a constant expression is feasible in worst case $O(n)$ with n the number of nodes in the tree, a more efficient approach is preventing constant expressions from being generated.

Constant expression detection An expression is a constant expression if all its children are constant expressions. As a base case, a leaf node is a constant expression if it is not a feature. This problem statement allows us to define a recursive algorithm to detect constant expressions. It should be noted that its complexity is $O(n)$ only in the worst case, when the tree is a constant expression. Upon detecting a non constant subtree, the algorithm returns early without a full traversal.

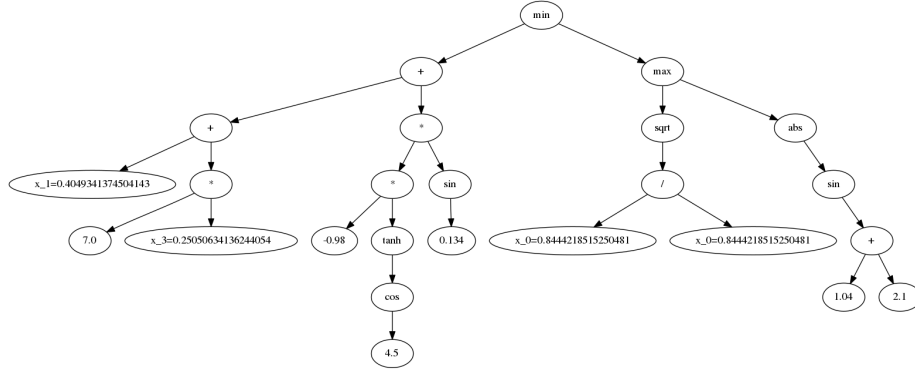


Figure 5: Tree before subtree folding.

Preventing constant expressions Using the checking procedure in the initialization, a tree marked as a constant expression is not allowed in the initialization procedure. It is still possible to create constant expressions by applying mutation and crossover. If the left subtree of a node is a constant expression, and the right is not, and this right is replaced by either mutation or crossover with a constant expression then the tree becomes a constant expression. The mutation operator will not generate constant subtrees, so this leaves only crossover. Our tool does not prevent constant expressions from forming in this way, the evaluation following crossover will filter out the constant expressions using evolutionary pressure.

5.1.2 Folding

Constant subtree problem A tree can contain subtrees that represent constant expressions. This is an immediate effect of the GP algorithm trying to evolve the correct constant. This can lead to large subtrees that can be represented by a single node. Nodes used in such a constant subtree are not available for base functions. They waste memory and evaluation cost, without their presence the tree could become a fitter instance.

Constant subtree folding It is possible to use a depth sensitive objective function to try to mitigate this effect, but a more direct approach is replacing the subtrees. Using the previous constant expression detection technique we can collect all constant subtrees from a tree. We evaluate each subtree, and replace it with the constant value it represents. This leads to savings in nodes and possibly a reduction in depth. These savings can have an effect on the convergence as well. Mutation and crossover will no longer operate on constant subtrees, and the iterations and place in the tree that becomes available can be used to improve the fitness of the tree. Constant folding requires an $O(n)$ detection check, since the entire tree needs to be traversed. The folding operation itself is at worst $O(n)$, if the entire tree is constant.

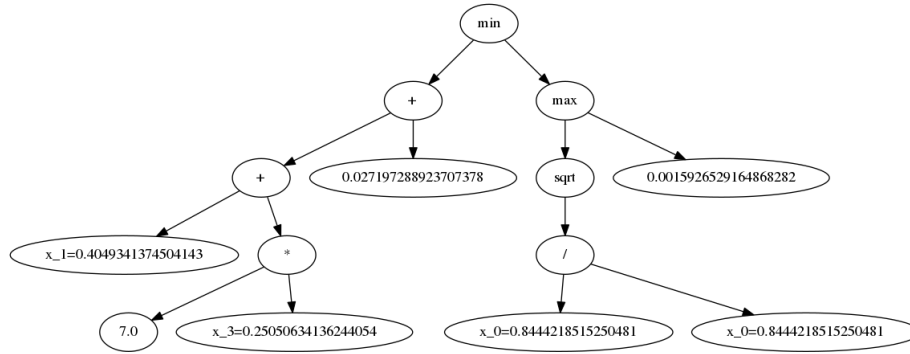


Figure 6: Tree after subtree folding.

Edge cases In Figure 6 we see the effect of applying the constant subtree folding. There are subtle edge cases that can't be detected using the above method. Consider the tree representing the expression

$$f(x) = \max(x, 4 + \sin(0.3))$$

The subexpression $4 + \sin(0.3)$ is detected and folded into the constant 4.26. We should not stop there, since f can still be a constant expression if $\forall i : x_i < 4.26$. To verify this we need to evaluate f on all i datapoints. In contrast, the constant subtree detection code needs only 1 evaluation. In Figure 6 we see a similar case in the right subtree where the first value of x_0 the right subtree is indeed constant. Even if we should evaluate f for all i , this will not guarantee us that f is indeed a constant expression. All this check then proves is that for all i in the training data, f is constant. The testing data holds new values, for which f may or may not be constant. We conclude that this edge case cannot be prevented from occurring. Another case occurs in expressions of the form

$$f(x, y) = \tan(y) * \sqrt{x/x}$$

In this reduced example the node $\sqrt{x/x}$ is simply $\sqrt{1}$. Outside of this example, detecting such cases is non trivial. There is a clear benefit to do so, despite their low occurrence: if the domain of x includes 0 this tree is never generated because it leads to division by zero in a single datapoint. Discarding this expression is however not needed, since $\sqrt{1}$ is a valid subexpression. One way to solve this is to use mathematical software to simplify the expressions, and then convert them back to tree representations. Deciding to apply this step should be based on the cost versus benefit and the frequency of such occurrences.

Conclusion Constant detection and folding mitigate some of the side effects of the constant optimization problem, they do not lead to a more efficient finding of the 'right' constant values. For this we need a continuous optimizer, an algorithm designed specifically to optimize real valued problems.

5.2 Optimizers

Hybridizing GP Using a real valued optimizer in combination with GP is a known solution [19, 6]. Which algorithm to combine is a difficult question. To our knowledge there is no comparison made between optimization algorithms to find out which is a better fit to combine with GP.

Problem statement Given a tree with k constant leaves, with all constant subtrees folded, we would like to find the most optimal values for those constants resulting in a better fitness. It is vital to perform the constant folding step before optimization takes place. Suppose a given tree has on average k constants, which after folding become j with $j \leq k$. Without folding the optimizer has to solve a k dimensional optimization problem, whereas after folding the problem is only j dimensional. The underlying approximation only has j dimensions, so this step is a necessity.

Initialization In most optimization problems the initial solution is not known, only the problem domain and an objective function. The problem we face here does have an initial solution, namely that generated by GP. Instead of choosing random points in the search space, we therefore opt by perturbing this initial solution. An immediate problem here is that the optimizer may simply converge on the initial solution. This risk can be high, given that GP already has evolved it as a suboptimal solution. The problem faced by the population initialization step 3.3.1 reappears here. We could pick random values in the search space, but these are likely to generate invalid trees. A balance between exploration and exploitation is once again called for.

Domain Each tree instance has a completely different domain for each constant. We cannot therefore guide the optimizer with domain specific knowledge. This also reflects in the choice of parameters for the optimizer, which will have to be suboptimal for specific problem instances, as we want them to work on the largest set of problems.

Comparison and configuration In order to keep the comparison between the algorithms fair they are each configured as similar as is possible. Since they are all population based, and given a maximum number of iterations, all three share these same parameter values. Each application of the algorithm has a significant cost in comparison with the main GP algorithm. In a single iteration with population p , the GP algorithm is likely to perform $2n$ evaluations (mutation and crossover). If we give the optimizer a population of m and q iterations, it will execute at most in the order of $m \times q$ evaluations. Based on optimal value for PSO [12] we use a default population of 50, and 50 iterations. This means the cost of the optimizer will quickly dominate that of the main GP algorithm, depending on when and on what we apply it. We can apply the optimizer on each expression each iteration, only on a selection of the best expressions at each iteration, or only at the end of phase on all, a selection, or only the best expressions. In our experiments we will show the effect of these choices on convergence.

5.2.1 ABC

Artificial Bee Colony [10] is a relative new nature inspired optimization algorithm. It is not limited to continuous optimization, and has even been used for symbolic regression itself [11]. One of its key advantages over other algorithms is a lower parameter count. Optimal values for an optimization algorithm have a large impact on its convergence behavior, so much so that other optimizers can be required to find the parameters of an optimizer. With a small parameter set finding optimal values becomes easier. Finding optimal values for these parameters is quite often related to the problem domain, and as we have seen each instance here will have a new domain. ABC is good at exploration (thanks to the scouting phase) though sometimes lacking in exploitation. To resolve this ABC can be combined with more exploitative algorithms [18].

Algorithm ABC is a population based algorithm, using 3 distinct phases per iteration. We will refrain from using the nature analogy in describing the algorithm as it can be confusing. The algorithm maintains a set of potential solutions and a population of particles. Each particle is either employed, onlooker, or scout. An employed particle perturbs a known solution, and if an improvement in fitness is obtained replaces the old solution. If this fails a preset number of iterations, the solution is discarded and a new one scouted. Scouting in this context is generating a new potential solution. After the employed phase the onlooking particles decide, based on a fitness weighted probability, which solutions should be further optimized. In contrast to the employed particles they swap out solutions each iterations, whereas an employed particle is linked to a single solution. Finally exhausted solutions are replaced by scouted solutions.

Initialization The algorithm initializes its solution set by perturbing it. Perturbation is done by multiplying each constant with a random number $\in [-1, 1]$. Using the modified constants the tree recalculates its fitness value and the global best. Each instance records its best solution. The original source, the expression tree that we wish to optimize, is retained as a solution, and not perturbed. This can lead to premature convergence if the algorithm is configured to focus too much on exploitation versus exploration.

Modification In the employed stage, a solution x at time i is modified by combining it with another randomly chosen solution y . The choice for y is made with replacement, and obviously $y \neq x$. With k random $\in [0, |x|)$

$$\begin{cases} x_{ij} = x_{ij} & j \neq k \\ x_{ij} = y_{ij} & j = k \end{cases}$$

The employed particle tries to improve the current source. If the modification leads to an improved fitness value, the solution is updated. Note that only a strict improvement warrants an update, an equal fitness value will not lead to an update. The onlooker phase is executed next. Each onlooker is assigned a solution using roulette wheel selection. For the set of solutions S we calculate fitness weights using:

$$w_i = \frac{1}{1 + f_i} \forall i \in S$$

Optimization in CSRM is a fitness minimization process, which leads to the fraction. A list is built using a cumulative sum of these weights. From this a uniform random number picks a single value based on these weights. A smaller fitness value will have a relatively larger section of the list compared with a larger fitness value resulting in a fitness bias in selection. Once assigned to an onlooker, the same modification process used by the employed particle is applied. Since selection is done with replacement it is possible that a single highly fit value is modified more than twice in an iteration. After the onlooker phase the algorithm checks in the scouting phase if any sources are exhausted. Exhaustion indicates that a solution can't be improved upon for at least *limit* modifications. Depending on the fitness value this limit can be reached faster for more fit values. There are several edge cases to consider here. Clearly this approach is beneficial if the exhausted solution has a poor fitness value. Improvement was impossible, so replacement is likely to improve the overall quality of the solution set by introducing new information. There is however the risk that algorithm discards highly fit solutions that fail to be improved. Discarding such a solution is warranted if that solution is a local optimum, but by the very nature of the problem statement we cannot know this in advance. We could prohibit discarding the best solution, even if the improvement limit has been exceeded. This does not solve all edge cases. If ABC is used in a multimodal search, it is still possible that valid solutions are discarded. Compounding the problem is that an equality update is not used. In CSRM we implement a strict check, so it is possible that equivalent solutions are discarded. With s scouts and e exhausted solutions $\min(s, e)$ solutions will be replaced with new solutions. These are generated using a Normal distribution. In most optimization problem good starting positions are not known and thus a random point is selected. In our problem we already have a reasonably good solution, so we use this initial value as the mean of the normal distribution and generate values within 2 standard deviations around that mean. A configurable scaling factor is introduced, in our test problems we use 20. This value is trade-off between exploration and an increasing probability for generating invalid solutions. The values chosen have a far greater range than those generated by the initialization process. We do not know the domain of our problem, nor do we have the computational resources to cover the entire floating point range. The initial value is already evolved as a fit solution to our problem. The probability that the true optimum is far beyond the range of our initial value is estimated as low, though we can never be sure of this. The scouted solution replaces the exhausted solution. In our problem statement this can lead to issues. There is no guarantee that the scouted solution is actually valid. As we have seen in the initialization problem 3.3.1 this probability can be quite high. It is possible that an increasing part of the solutions are invalid. These will still be chosen to contribute in the modification step. It is unlikely though not impossible that they contribute to the convergence. In the worst case with enough iterations and a very small domain there it is possible that the entire population is replaced with invalid solutions. In our implementation the values of the threshold regulating exhaustion are chosen such that this is unlikely to occur. One solution here is to apply the same generation loop used in the GP algorithm, which keeps generating solutions until a valid one is found. Given the already high cost in evaluations the optimizer introduces we have chosen not to use this here. This applies for all optimizers implemented in CSRM, not just ABC. Note that the same problem is present in the initialization step, although far

less severe given the small perturbation applied there.

Selection A solution is updated if the fitness value is improved. The new fitness weights for the roulette wheel selection are calculated and the global best is updated. Unlike PSO a solution is only updated if an improvement is measured. In contrast to DE, equal fitness values do not lead to an update. An equality update allows an optimization algorithm to cross zero gradient areas in the fitness landscape. The influence other solutions have on each other in ABC is not as great as in DE. From the modification stage we also observe that at most 2 dimensions per iteration per solution are modified. In PSO the entire position, in all dimensions, is updated. In DE this depends on the Cr parameter which we discuss in 5.2.3. This distinction can have a large impact on convergence. The balance sought here is influenced by the interdependence of the constants. The modifications made by the algorithm can be seen as a process trying to extract information about this dependency. Suppose we have k constants, of which 2 have a large effect on the result. Then it only makes sense to modify those 2 dimensions, modification in the others does not gain anything except noise. The problem becomes more difficult if those two are correlated. Modifying all dimensions will not guide our optimization process as clearly as modifying only those 2. Modifying only a single one of the 2, as in ABC, is too strict as we lose the information about the correlation. We do not know in advance if our problem instances are separable or not. Related to this discussion, PSO can be sensitive to a bias along the axes of the dimensions [20] with improvements suggested in recent work [2].

Cost With a solution set of n, m employed, j scouts, i onlookers and k iterations we now look at the evaluation cost. Initialization requires n evaluations. Each iteration we execute m evaluations in the employed phase, i evaluations in the onlooker phase and at most j evaluations in the scouting phase. With m, j and i all \leq to n we have per iteration at most 3n evaluations. With our configuration, listed below, this value will be at most 2n. This results in an evaluation complexity of $n + k \cdot 2n$, or $O(kn)$.

Configuration

- $\text{limit} = 0.75 * \text{onlookers} / 2$: If a solution can't be improved after this many iterations, it is marked exhausted and will be scouted for a new value. This limit is scaled by the number of dimensions per instance.
- $\text{population} = 50$: This is the solution set, or set of sources.
- $\text{onlookers} = 25$: The number of onlookers, instances that will be assigned solutions to exploit based on fitness values. Setting this value to half that of the employed finds a balance between exploitation and evaluation cost.
- $\text{employed} = 50$: Instances that try to improve an assigned solution. If we use a value lower than the solution set we have to define an assignment procedure, which would mimick the onlooker phase. We therefore set the employed count equal to the size of the solutions set.

- scouts = 25 : This is a maximum value, up to this number are used to scout after a solution is exhausted. A higher scouting value leads to more exploration, a lower value favors exploitation. More exploration would result in the initialization problem dominating the runtime cost of the optimizer.

This configuration is guided by the findings in [10].

5.2.2 PSO

Particle Swarm Optimization [12] is one of the oldest population based metaheuristics. It consists of n particles that share information with each other about the global best solution.

Algorithm

Initialization Each particle is assigned a n dimensional position in the search space. A particle's position is updated using its velocity. This last is influenced by information from the global best and the local best. The concept of inertia is used to prevent velocity explosion [3]. Each particle is given a random location at start. In our application we already have an (sub)optimal solution, the constant values in the tree instances have been evolved by the GP algorithm. Rather than pick random values, we perturb the existing solution. This is a difficult trade-off. If the perturbation is not large enough the optimizer will simply converge on the known solution. If we perturb too greatly the risk for invalid solutions increases, rendering a large selection of the population invalid. We can initialize the population with n perturbed solutions or with $n-1$ with the initial value remaining intact. The $n-1$ solution is useful to test the algorithm, ideally the swarm will converge on the known best value. When applying the optimizer the n -perturbation approach is used, minimizing the risk for premature convergence. CSRM multiplies each constant with a random value in $[0,1]$. Each particle is assigned a small but non-zero velocity. The reason for this is again avoiding premature convergence. Without this velocity all particles are immediately attracted to the first global best. While attraction to this value is desired, it should not dominate the population. The small value of the initial velocity once again reflects an empirical discovered balance between exploration and exploitation. Each particle is assigned an inertia weight. This value is one approach to combat the velocity explosion problem, which we will cover in 5.2.2. Finally the global best is recorded.

Modification The algorithm updates all particles in sequence, then records the new global best. Let d be the dimension of the problem, or in our case the number of constants in the tree to optimize. The velocity v at iteration i of a particle is updated using:

$$v_{i+1j} = w_i * v_{ij} + C_1 * (p_{ij} - g_{ij}) * R_1 + C_2 * (p_{ij} - G_{ij}) * R_2 \forall j \in [0, d]$$

with

- v_i Current velocity

- p_i Current position (set of constant values)
- g_i Local best
- G_i Global best
- C_1 Constant weight influencing the effect the local best has on the velocity.
- C_2 Constant weight influencing the effect the global best has on the velocity.
- w_i Inertia weight simulating physical inertia.
- R_1 Random value perturbing the effect the local best has on the velocity.
- R_2 Random value perturbing the effect the global best has on the velocity.

Without the inertia weight PSO has issues with velocity explosion, the velocity has a tendency to increase to large values. This increases the distance between particles, but more importantly is far more likely to generate positions that are no longer inside the domain of one or more of the dimensions. Inertia weighting will dampen this effect. The position is updated using :

$$x_{i+1,j} = x_{ij} + v_{ij} \forall j \in [0, d)$$

The R_1 and R_2 parameters make the modification stochastic, they introduce perturbations in the calculation. These changes have the benefit that they can break non optimal behavior resulting from the deterministic calculation. If there is a suboptimal best value (local or global) that leads to a too strong attraction and thus forcing premature convergence, we can with a certain probability escape from such a value by perturbing the velocity calculation. The C constants determine how strong the effect is of respectively the local best and the global best. This reflects the balance between exploration and exploitation respectively, where a particle is influenced more by its own information or that of the swarm. After all particles are updated, the new global best is recorded for the next iteration.

Selection An interesting difference with other algorithms is that the position is always updated, whether it improves the fitness or not. The local best records the best known position, but the particle is allowed to visit positions with lower fitness values. This allows it to escape local optima. The global best is obviously only updated with an improved value. The comparison is strict, meaning that only a better value can become the new global best. This may not seem significant, but allowing equality updates can actually benefit an optimization process. To see this it helps to view the fitness domain as a landscape with troughs, valleys and heights. Allowing for equality updates allows the global best to move despite no apparent improvement in fitness. In this analogy we represent a (sub) optimal value with low points in the landscape. This allows it to cross areas where a zero gradient is observed, for example between two low points where the lower one is a local optima.

Cost With a population size of n , the algorithm requires n fitness evaluations per iteration. The computational cost of updating the velocity and position is small given the evaluation cost of an expression tree over several datapoints. However, it is linear in d , the number of dimensions. If the tree increases in size, the number of constants can increase in the worst case exponentially. On average due to the construction algorithm we expect that half the leaves in the tree are constants. Given our discussion of the constant folding algorithm we know that a full binary tree is unlikely, so the number of constant nodes is equally unlikely to increase exponentially, but will nonetheless scale poorly. In our implementation we will halt the algorithm if it cannot improve the global best after $k/2$ consecutive iterations, where k is the maximum number of iterations it can execute. The initialization stage adds another n evaluations, in addition to n per iteration. The total cost in fitness evaluations is therefore $n(k+1)$, resulting in a worst case evaluation complexity of $O(nk)$.

Configuration This overview gives the values of each parameter used in CSRM's PSO implementation.

- $C_1 = 2$
- $C_2 = 2$: Setting both to 2 is recommended as the most generic approach [13].
- $w_i = \frac{1+r}{2}$ with r random in $[0,1]$: In early implementation the inertia weight was kept constant [5] There are a large number of strategies for an inertia weight. Dynamically decreasing inertia may improve convergence significantly. We opt for a random inertia weight as it has been shown [1] to lead to faster convergence. Since our use case requires fast convergence on very limited iterations, this strategy is clearly favored.
- R_1, R_2 r with r random in $[0,1]$
- population = 50 : PSO is not sensitive to populations larger than this value, providing a robust default value. [14]

CSRM's optimizer does not set constraints on the domain of each constant, it is different for each problem instance. Finding the domain of a constant in the expression tree requires a domain analysis of the expression tree. With features in the tree involved, finding the exact domain is infeasible, given that some of the datapoints are unknown. It is therefore possible that a particle obtains values outside the valid domain of one or more constants, resulting in an invalid expression tree. This will result in the particle temporarily no longer contributing to the search process.

5.2.3 DE

Differential Evolution is a vector based optimization algorithm, or rather as the name implies, it operates by computing the difference between particles.

Algorithm The algorithm has a population of n vectors, similar to the other algorithms it holds a linear set of values to optimize, one per dimension.

Initialization Similar to our approach in initialization PSO, we perturb a known (sub) optimal solution. A vector stores its current value, and the best value.

Modification Each iteration the algorithm processes all vectors. For each vector \vec{v} , three distinct randomly selected vectors are selected. From these 3 vectors a new 'mutated' vector is obtained:

$$\vec{v} = \vec{w} + F(\vec{y} - \vec{z})$$

With $\vec{w}, \vec{y}, \vec{z}$ randomly chosen and not equal to \vec{v} .

The selection occurs with replacement. Several selection schemes exists, and the size of the selection is equally configurable. From this step the algorithm lends its name. The F factor influences the effect of the difference. Then we apply a crossover operation, using vectors x and v and probability parameter Cr. We select a random index j with $j \in [0, |x|]$. We then create a new vector u:

$$u_i = k_i \forall i \in [0, |x|]$$

and k_i equal to

$$\begin{cases} v_i & i = j \vee r < Cr \\ x_i & i \neq j \wedge r \geq Cr \end{cases}$$

This is binomial crossover, another frequently used selection operation is exponential crossover.

Selection For a given selected vector \vec{v} and created vector \vec{u} we now test if \vec{v} is a better candidate than \vec{u} , in other words has a lower or equal fitness value. Note the distinction here with PSO, the equality test allows DE vectors to cross areas without a gradient. If $f(\vec{u}) \leq f(\vec{v})$ the vector is replaced with \vec{u} . The global and local best are updated as well. A difference with PSO is that a PSO particle changes regardless of fitness value, whereas in DE the modification is only committed if a better or equal fitness value is obtained. The first approach allows an optimization algorithm to break free from local optima. DE uses the random selection of other vectors to create a similar effect. If we use the landscape analogy, as long as at least one DE vector is outside a depression in the landscape, but all the others are converging to the suboptimal minimum, DE has a probability to escape a local optima. Unlike PSO DE (in our configuration) does not use the global best in its calculations, sharing of information is completely distributed over the vectors.

Cost For each vector 3 other vectors are used, or restated we create 2 new vectors. Similar to PSO these calculations have a complexity linear in d, the dimensionality of the problem. The fitness function is called once per iteration per vector. Compared to PSO we therefore have the exact same evaluation complexity of $O(nk)$.

Configuration CSRM uses a DE/rand/2/bin configuration. The DE/x/y/z notation reflects its main configuration, where x is the vector perturbed, y is the vectors used in the differential step and z is the crossover operation (binomial). This configuration is referenced [21] as one of the most competitive for multimodal problems with good convergence to the global optimum. Since we start from a probable local optimum the choice for a random vector instead of the global best vector also helps avoid premature convergence. This overview gives the values of each parameter used in CSRM’s DE implementation.

- $F = 0.6$: F should be in $[0.4, 1]$ Large F values favor exploration, whereas small F values favor exploitation. The value of 0.6 is reported as good starting value. In our problem domain we already have a (sub) optimal solution which we wish to improve, so the risk of premature convergence is present, hence the small bias for exploration.
- $Cr = 0.1$: The Cr values should be in $[0, 0.2]$ for separable functions, and $[0.9, 1]$ for non separable functions. We cannot assume dependency between the constants, and therefore use a value of 0.1. This results in DE focussing alongside the axes of each dimension in its search trajectory.

Compared to PSO DE has a low parameter count, optimal values for these parameters can be found in literature [4]. The population size should be $t * d$ with t in $[2, 10]$. Since we do not know d in advance, and to keep the comparison fair we set the population at 50, allowing for optimal values for up to 25 dimensions (constants). While DE has a small set of parameters, their effect is still quite pronounced. There exists implementations of DE that use self adapting parameters, but this is beyond our scope. It should be noted that CSRM’s optimizer has a very small optimization budget (in evaluation cost) and each new problem has potentially new characteristics. We therefore chose for the most robust values and configuration.

6 Experiments

6.1 Benchmark problems

Recent work on the convergence of GP-based SR [15, 16] featured a set of benchmark problems that pose convergence problems for SR implementations. We reuse these problems in our work in order to study convergence of CSRM’s implementation.

6.1.1 Problems

These problems use at most five features, when present to CSRM we do not give the algorithm this knowledge. In other words it assumes each problem is a function of 5 features which may or may not influence the expected outcome. This is an extra test in robustness for the algorithm, while also testing the algorithm’s capability as a classifier.

$$1.57 + (24.3 * x_3)$$

$$\begin{aligned}
& 0.23 + 14.2 * \frac{x_3 + x_1}{3.0 * x_4} \\
& -5.41 + 4.9 * \left(\frac{x_3 - x_0 + \frac{x_1}{x_4}}{3 * x_4} \right) \\
& -2.3 + 0.13 * \sin(x_2) \\
& 3.0 + (2.13 * \ln(x_4)) \\
& 1.3 + 0.13 * \sqrt{x_0} \\
& 213.80940889 - 213.80940889 * e^{-0.54723748542 * x_0} \\
& 6.87 + 11 * \sqrt{7.23 * x_0 * x_3 * x_4} \\
& \frac{\sqrt{x_0}}{\ln(x_1)} * \frac{e^{x_2}}{x_3^2} \\
& 0.81 + 24.3 * \frac{2.0 * x_1 + 3.0 * x_2^2}{4.0 * x_3^3 + 5.0 * x_4^4} \\
& 6.87 + 11 * \cos(7.23 * x_0^3) \\
& 2.0 - 2.1 * \cos(9.8 * x_0) * \sin(1.3 * x_4) \\
& 32 - 3.0 * \frac{\tan(x_0)}{\tan(x_1)} * \frac{\tan(x_2)}{\tan(x_3)} \\
& 22 - 4.2 * ((\cos(x_0) - \tan(x_1)) * \frac{\tanh(x_2)}{\sin(x_3)}) \\
& 12.0 - 6.0 * \frac{\tan(x_0)}{e^{x_1}} * (\ln(x_2) - \tan(x_3))
\end{aligned}$$

6.2 Operators

6.3 Constant Folding

6.3.1 Savings

6.4 Constant optimization

We look at the effect constant optimization using different algorithms has on different configurations of the tool. The measures used in the comparison are best fitness on training and test data, mean fitness on training and test data, and optimization cost.

6.4.1 Test problem

To verify our implementation for the optimizers we use a simple test problem and observe for each optimizer if it is able to optimize this instance to a known optimal value.

$$f(x_0, x_1, x_2) = 1 + x_1 * \sin(5 + x_2) * x_0 + (17 + \sin(233 + 9))$$

We give each optimizer a population of 50, 50 iterations and compare the results for 10 runs, displaying best value obtained, mean, and standard deviation of the fitness values compared to the known best value.

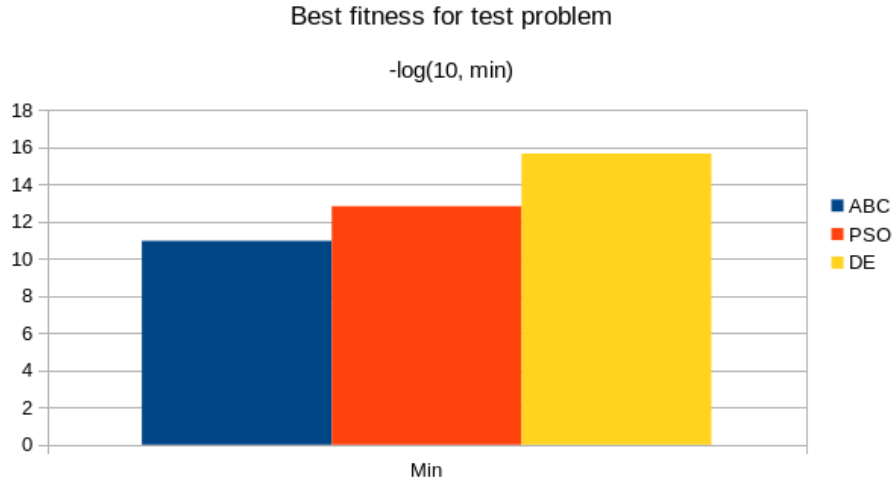


Figure 7: Logarithmic value of best fitness for each optimizer.

Best fitness In Figure 7 we see that DE outperforms PSO and ABC with several orders of magnitude. The best fitness value obtained was 2.22×10^{-16} . As smaller but significant difference is present between PSO and ABC. This result is somewhat surprising given that fact that ABC is allowed to perform more evaluations in its configuration. From our previous discussion 5.2.2,5.2.3,5.2.1 we can conclude that for this test problem DE is clearly preferable as it obtains the best result at minimum cost. ABC has almost double the cost compared to PSO and DE, with PSO and DE having an equal cost in evaluations. The results on this testproblem do not necessarily mean that in the application of the three optimizers the results will be identical. Here we have a known optimal solution and want to observe how fast the optimizers converge to it. When we optimize evolved expressions we do not know what the optimal solution is. The problem statement is different, and so the convergence behavior is likely to differ as well.

Distribution of fitness In Figures 8 and 9 we see that both the mean and standard deviation follow the same pattern as seen for the minimum fitness value with DE leading the others by several orders of magnitude. With all three distributions behaving similarly, this result provides a more solid foundation for our conclusions that for this problem DE is indeed the better optimizer.

6.4.2 Optimizer experiments setup

We test the 15 expressions with the following configuration:

- population : 20
- minimum depth : 4

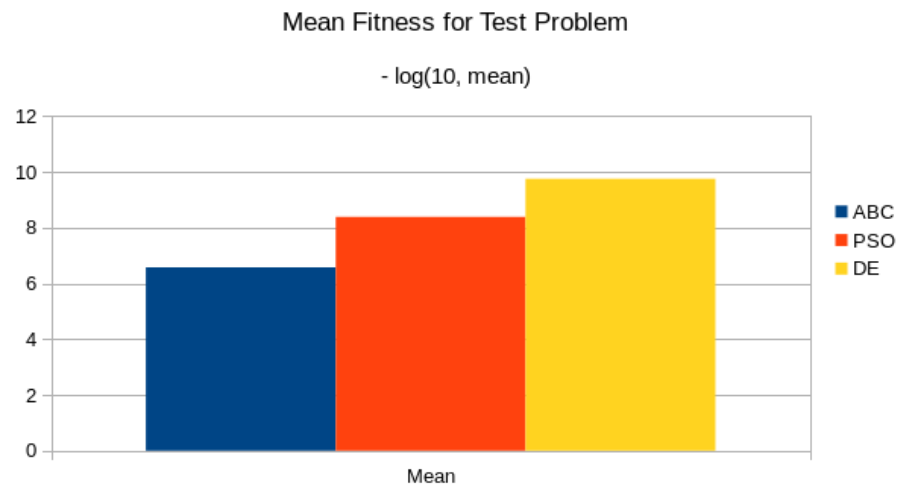


Figure 8: Logarithmic scaled mean fitness for each optimizer.

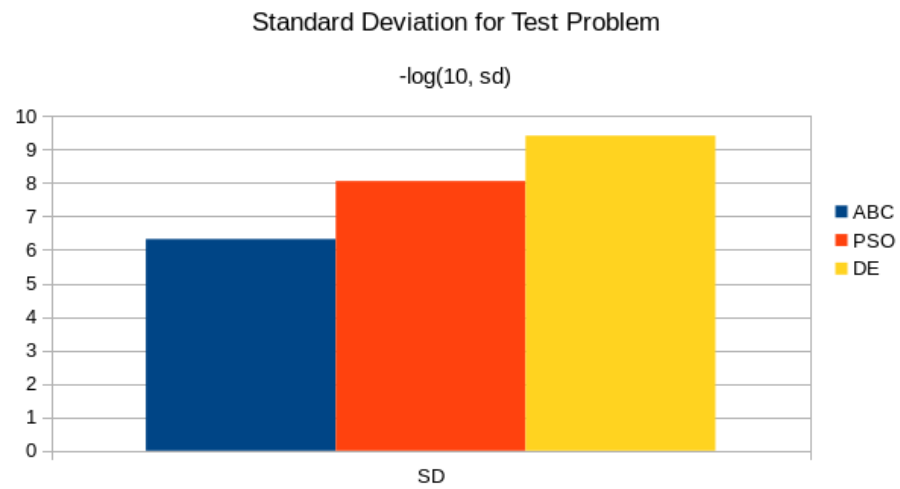


Figure 9: Logarithmic scaled standard deviation fitness for each optimizer.

- maximum depth : 10
- phases : (2, 5, 10)
- generations per phase : 20
- datapoints : 20
- range : [1,5]
- features : 5
- archivesize : 20
- expressions to archive per phase : 4
- optimization strategy : optimize expressions archived at end of phase

6.4.3 Measures

We compare the relative gain in fitness compared to not using an optimizer for all expressions. In other words, if m_n is a measure obtained by the algorithm without the optimizer, and m_a the same measure with the optimizer, we then define the relative gain as :

$$g_{ma} = \frac{m_n}{m_a}$$

If m_a is zero, we use $-\log_{10}(m_n)$ to represent the gain. If both are zero, the gain is obviously 1. A value of $g > 1$ indicates the ratio with which the optimizer improves the result. A g value < 1 indicates a regression. These 15 functions have wildly varying convergence behavior. In order to make sense of the data, we then apply a log scale :

$$gl_{ma} = -\log_{10}(g_{ma})$$

A value of $gl_{ma} > 0$ indicates improvement, with the units transformed to orders of magnitude. A zero value indicates no improvement is registered, and negative values indicate regression. As measures we use the best fitness value on the training data, and the best on the full data set. We take the mean of the fitness of the 5 best expressions on training and the full data as well. This last measure gives us an indication on how the optimization process acts on the 'best' set of the population. Note that in our configuration, the 4 best expressions are always optimized.

6.4.4 2 Phases

In Figure 10 we see the performance of the algorithms on training data. In Table 1 we see that for problems 0, 4, 5 there is no improvement possible (e.g. 0 zero fitness value), which explains the absence of any value in the figure. We see that for the training fitness data the improvements are significant, with ABC scoring an increase of 2.5 orders of magnitude for problem 6. For the other problems the increase is still large, especially given that our fitness function has a range of [0,1]. We also observe

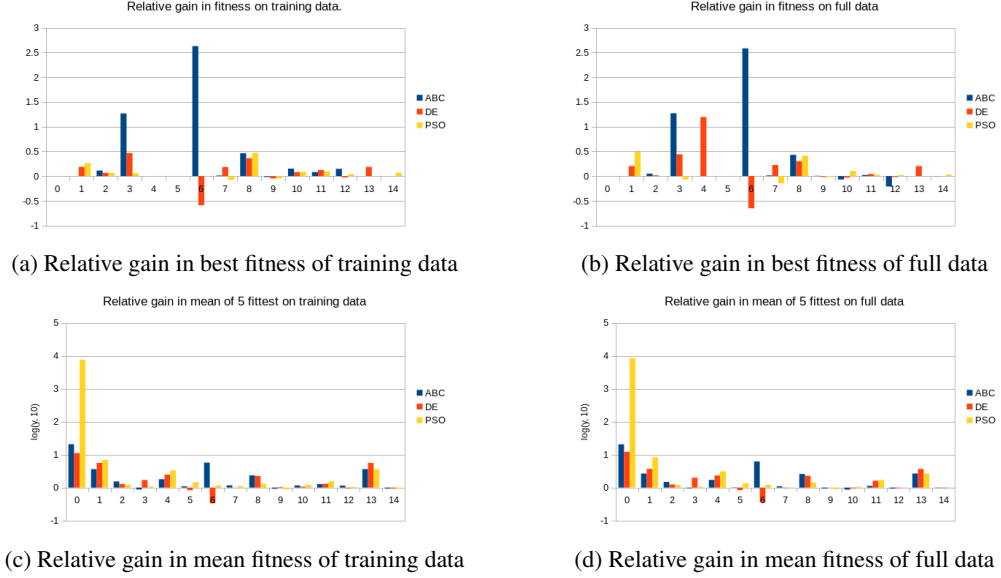


Figure 10: Relative gain of optimizer after 2 phases.

the significant regression for problem 6. This is likely caused due to overfitting. The algorithm in question (DE) optimizes the 4 best candidates of the last phase (1), but it is possible that these optimized expressions actually form a local optimum. By archiving these the convergence of the algorithm is hindered in the next phase. Note that DE allows equality updates, where expressions with the same fitness values are accepted as better. The same behavior occurs in a far less significant effect for expressions 7 and 9. A second explanation is our implementation of the population. The algorithm enforces distinct fitness values for all expressions. In an edge case it is possible that these optimized samples form a barrier, preventing other expressions from evolving past them. The optimized expressions in effect trap the rest of the population, which given our low generation count can explain this behavior. The mean fitness of the 5 best expressions shows significant improvements. Important to observe is the similarity between the two plots, the correlation between fitness values on training and full data is strong. This was a concern in the setup of the experiments. The optimizers could introduce overfitting on the training data. This risk is mitigated by the relatively low number of iterations each optimizer has been allocated. For the minimum fitness on the full data ABC outperforms the others. For the mean evaluation PSO is a good candidate. In this stage of the experiments, there is no single algorithm that is ideal for all problems. This once again confirms the NFL theorem [22].

Table 1: Relative Gain in minimum fitness on training data after 2 phases.

	0	1	2	3	4	5	6	7	8	9	10	11	12	13	14
ABC	1.000e+00	1.000e+00	1.295e+00	1.848e+01	1.000e+00	1.000e+00	4.253e+02	1.030e+00	2.897e+00	9.615e-01	1.415e+00	1.207e+00	1.404e+00	1.000e+00	9.970e-01
DE	1.000e+00	1.536e+00	1.163e+00	2.923e+00	1.000e+00	1.000e+00	2.584e-01	1.526e+00	2.294e+00	8.972e-01	1.210e+00	1.327e+00	9.397e-01	1.536e+00	9.960e-01
PSO	1.000e+00	1.839e+00	1.169e+00	1.150e+00	1.000e+00	1.000e+00	1.000e+00	8.356e-01	2.947e+00	8.971e-01	1.226e+00	1.247e+00	1.096e+00	1.000e+00	1.172e+00

Table 2: Gain in minimum fitness on full data after 2 phases.

	0	1	2	3	4	5	6	7	8	9	10	11	12	13	14
ABC	1.000e+00	1.000e+00	1.122e+00	1.865e+01	1.000e+00	1.000e+00	3.831e+02	1.034e+00	2.691e+00	1.016e+00	8.528e-01	1.051e+00	6.169e-01	1.000e+00	9.938e-01
DE	1.000e+00	1.597e+00	1.039e+00	2.757e+00	1.565e+01	1.000e+00	2.233e-01	1.674e+00	2.004e+00	9.617e-01	9.337e-01	1.103e+00	9.584e-01	1.597e+00	9.970e-01
PSO	1.000e+00	3.098e+00	9.915e-01	8.538e-01	1.000e+00	1.000e+00	1.000e+00	7.145e-01	2.586e+00	9.617e-01	1.274e+00	1.067e+00	1.055e+00	1.000e+00	1.069e+00

Table 3: Relative gain in mean fitness of 5 fittest expressions on training data after 2 phases.

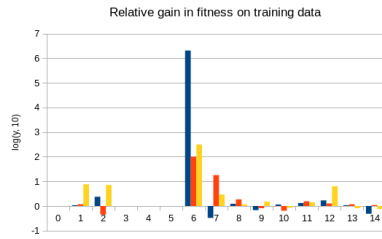
	0	1	2	3	4	5	6	7	8	9	10	11	12	13	14
ABC	1.335e+01	1.113e+00	4.824e+00	5.543e-01	3.479e+02	4.395e+02	7.059e-02	4.477e-01	1.215e+00	9.073e-01	2.214e+00	9.169e-01	1.564e+00	1.113e+00	4.064e-01
DE	3.620e+01	1.299e+00	1.115e+00	1.128e+00	1.204e-01	5.385e+08	3.202e+01	1.338e+00	1.835e+00	1.154e+00	1.268e+00	1.626e+00	1.246e+00	1.299e+00	7.247e-01
PSO	5.046e+02	8.017e+00	2.302e+00	4.220e-01	1.665e+01	2.237e+01	2.144e-01	6.974e-01	1.240e+00	9.647e-01	1.408e+00	1.347e+00	5.933e+00	7.152e-01	8.523e-01

Table 4: Relative gain in mean fitness of 5 fittest expressions on full data after 2 phases.

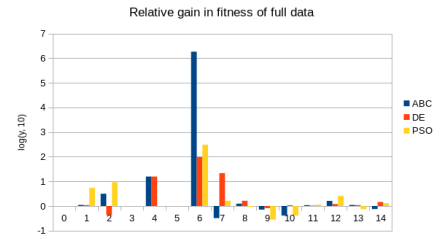
meanfullfitness algorithm	0	1	2	3	4	5	6	7	8	9	10	11	12	13	14
ABC	1.382e+01	2.306e-01	5.950e+00	5.551e-01	3.264e+02	5.091e+02	7.308e-02	4.463e-01	1.188e+00	9.214e-01	4.665e-01	7.197e-01	1.463e+00	2.306e-01	8.096e-01
DE	3.244e+01	1.315e+00	1.068e+00	1.228e+00	1.153e-01	5.392e+08	3.234e+01	1.348e+00	1.685e+00	1.153e+00	1.193e+00	1.441e+00	1.210e+00	1.315e+00	1.199e+00
PSO	4.920e+02	6.794e+00	2.531e+00	4.117e-01	1.721e+01	9.994e+00	2.146e-01	6.486e-01	1.113e+00	7.060e-01	5.830e-01	1.324e+00	2.613e+00	8.126e-01	1.094e+00

6.4.5 5 Phases

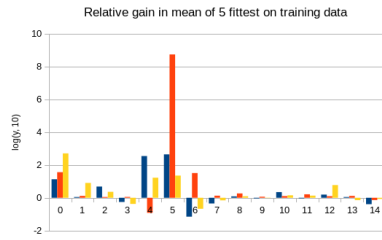
With 5 phases we see in Figure11 a more diverse effect. While ABC scores exceptionally good on problem 6, in sharp contrast with the 2 phase experiment, we see that PSO scores overall better for the training data. These results are logarithmic scaled, an improvement in fitness of factor 10 results in a value of 1 in the plots. When it comes to improving the mean of the best 5 expressions, PSO is a stable choice if we disregard the outlier values for problem 5. The correlation between training and full fitness scores is good for both measures. This demonstrates that the optimizer is not in this experiment introducing overfitting on the training data. The adverse effect of the optimizer on some test problems is still present. For the best fitness values on the full data DE is the better candidate. While ABC scores exceptionally high on problem 6, DE scores better overall. When we look at the mean there is no clear winner.



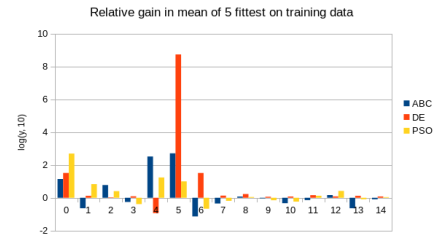
(a) Relative gain in best fitness of training data



(b) Relative gain in best fitness of full data



(c) Relative gain in mean fitness of training data



(d) Relative gain in mean fitness of full data

Figure 11: Relative gain of optimizer after 5 phases.

Table 5: Relative Gain in minimum fitness on training data after 5 phases.

algorithm	0	1	2	3	4	5	6	7	8	9	10	11	12	13	14
ABC	1.000e+00	1.074e+00	2.371e+00	1.000e+00	1.000e+00	1.000e+00	2.050e+06	3.277e-01	1.230e+00	6.779e-01	1.146e+00	1.321e+00	1.690e+00	1.074e+00	4.837e-01
DE	1.000e+00	1.176e+00	4.459e-01	1.000e+00	1.000e+00	1.000e+00	9.987e+01	1.790e+01	1.853e+00	8.129e-01	6.428e-01	1.553e+00	1.266e+00	1.176e+00	1.095e+00
PSO	1.000e+00	7.591e+00	7.112e+00	0.000e+00	1.000e+00	1.000e+00	3.117e+02	2.904e+00	1.159e+00	1.506e+00	8.442e-01	1.396e+00	6.302e+00	8.146e-01	7.435e-01

Table 6: Gain in minimum fitness on full data after 5 phases.

	0	1	2	3	4	5	6	7	8	9	10	11	12	13	14
ABC	1.000e+00	1.103e+00	3.145e+00	1.000e+00	1.565e+01	1.000e+00	1.846e+06	3.220e-01	1.243e+00	7.106e-01	4.082e-01	1.080e+00	1.613e+00	1.103e+00	7.483e-01
DE	1.000e+00	1.087e+00	4.240e-01	1.000e+00	1.565e+01	1.000e+00	9.756e+01	2.158e+01	1.629e+00	8.213e-01	1.079e+00	1.040e+00	1.202e+00	1.087e+00	1.454e+00
PSO	1.000e+00	5.426e+00	9.234e+00	0.000e+00	1.000e+00	1.000e+00	3.070e+02	1.634e+00	9.216e-01	2.786e-01	4.071e-01	1.131e+00	2.563e+00	7.444e-01	1.288e+00

Table 7: Relative gain in mean fitness of 5 fittest expressions on training data after 5 phases.

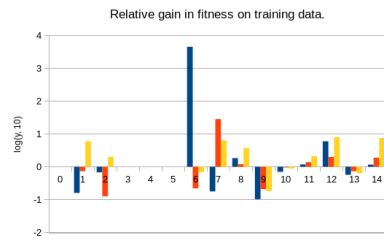
	0	1	2	3	4	5	6	7	8	9	10	11	12	13	14
ABC	2.064e+01	3.632e+00	1.536e+00	8.774e-01	1.788e+00	1.096e+00	5.726e+00	1.173e+00	2.348e+00	9.209e-01	1.163e+00	1.280e+00	1.151e+00	3.632e+00	9.300e-01
DE	1.117e+01	5.624e+00	1.277e+00	1.696e+00	2.468e+00	8.307e-01	3.333e-01	9.953e-01	2.249e+00	1.038e+00	1.067e+00	1.305e+00	9.416e-01	5.624e+00	9.302e-01
PSO	7.494e+03	6.977e+00	1.225e+00	1.091e+00	3.335e+00	1.461e+00	1.154e+00	1.118e+00	1.300e+00	8.794e-01	1.212e+00	1.559e+00	1.041e+00	3.632e+00	9.242e-01

Table 8: Relative gain in mean fitness of 5 fittest expressions on full data after 5 phases.

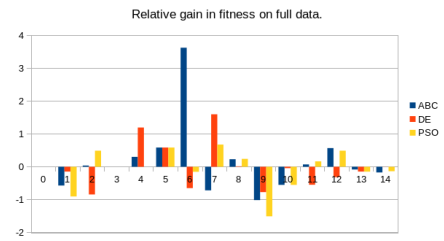
	0	1	2	3	4	5	6	7	8	9	10	11	12	13	14
ABC	2.048e+01	2.672e+00	1.487e+00	9.516e-01	1.709e+00	1.021e+00	6.217e+00	1.093e+00	2.589e+00	9.346e-01	8.686e-01	1.126e+00	9.308e-01	2.672e+00	9.816e-01
DE	1.213e+01	3.687e+00	1.244e+00	1.997e+00	2.324e+00	8.352e-01	3.572e-01	9.731e-01	2.288e+00	1.009e+00	9.496e-01	1.616e+00	9.641e-01	3.687e+00	9.668e-01
PSO	8.382e+03	8.286e+00	1.204e+00	1.071e+00	3.111e+00	1.362e+00	1.219e+00	9.817e-01	1.411e+00	8.939e-01	1.059e+00	1.698e+00	1.018e+00	2.672e+00	9.840e-01

6.4.6 10 Phases

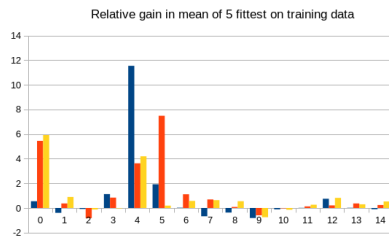
If we observe the convergence after 10 phases we see a more pronounced effect. In Figure 12 we see that for several problems the optimizers are no longer improving w.r.t the unoptimized algorithm. This only holds for the best values, for the mean values the improvements are still significant. It becomes clear that the optimizer can force the algorithm into a local optimum from which it becomes hard to escape. The correlation between fitness results on the training data and full data is starting to weaken as well, in comparison to the experiments with 2 and 5 phases. If we look at the fitness values for the full data DE is the more stable of the three algorithms. When it regresses its losses are smaller than the others, while its gains are strongest on the most problems. For the mean fitness of the full data a similar argument can be made, with the exception of problem 2 where DE fails severely. Another aspect is that after 100 generations the fitness values are extremely small, in the order of $1e-15$. We measure the relative gain with respect to the algorithm without an optimizer, but as the fitness values decrease rounding errors start to influence the calculations more and more. The fitness values are approaching the floating point epsilon values. For our implementation epsilon is set at $2.22e-16$. For problem 0, a minimum fitness value of 0 is found after 2 phases. For others far more iterations are needed. We need to make a trade-off in order to be able to compare all 15 problems. Giving each problem an equal budget in iterations is the more fair approach. Another approach is implementing a stop condition that halts within a certain distance of a desired fitness threshold, but this approach is fraught with issues. There is no guarantee exactly how many iterations are needed. This approach requires knowing the problem 'hardness' in advance, but by the very definition of our problem statement we do not know how hard our problem is. We do not know the optimal value, if there is a singular optimal value. In general the topology of the search space SR tries to traverse is not known. A practitioner with a real world problem faces the same issues. A more robust approach is stating in advance how much resources the algorithm can use in its search, and terminate if that budget is exhausted. The exact definition of resource is nuanced. We can use time, but this depends to a large extent on the implementation. A more solid measure is the number of fitness evaluations. Even this is not a constant measure, not all evaluations are equal in computational complexity.



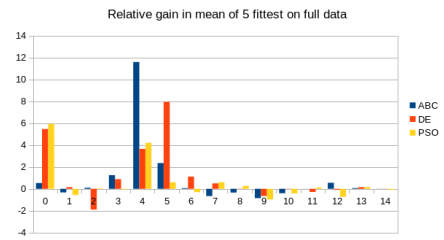
(a) Relative gain in best fitness of training data



(b) Relative gain in best fitness of full data



(c) Relative gain in mean fitness of training data



(d) Relative gain in mean fitness of full data

Figure 12: Relative gain of optimizer after 10 phases.

Table 9: Relative Gain in minimum fitness on training data after 10 phases.

	0	1	2	3	4	5	6	7	8	9	10	11	12	13	14
ABC	1.000e+00	1.591e-01	6.759e-01	1.000e+00	1.000e+00	1.000e+00	4.508e+03	1.771e-01	1.831e+00	1.029e-01	6.983e-01	1.176e+00	5.945e+00	5.674e-01	1.158e+00
DE	1.000e+00	7.304e-01	1.254e-01	1.000e+00	1.000e+00	1.000e+00	2.197e-01	2.816e+01	1.201e+00	2.068e-01	9.480e-01	1.377e+00	1.989e+00	7.304e-01	1.891e+00
PSO	1.000e+00	5.976e+00	1.980e+00	1.000e+00	1.000e+00	1.000e+00	6.857e-01	6.335e+00	3.739e+00	1.822e-01	8.852e-01	2.105e+00	7.945e+00	6.334e-01	7.526e+00

Table 10: Gain in minimum fitness on full data after 10 phases.

	0	1	2	3	4	5	6	7	8	9	10	11	12	13	14
ABC	1.000e+00	2.679e-01	1.087e+00	1.000e+00	2.000e+00	3.846e+00	4.224e+03	1.908e-01	1.694e+00	9.627e-02	2.807e-01	1.178e+00	3.699e+00	8.255e-01	6.709e-01
DE	1.000e+00	7.126e-01	1.432e-01	1.000e+00	1.565e+01	3.846e+00	2.232e-01	3.965e+01	1.031e+00	1.687e-01	9.052e-01	2.827e-01	4.813e-01	7.126e-01	1.005e+00
PSO	1.000e+00	1.246e-01	3.083e+00	1.000e+00	1.000e+00	3.846e+00	7.019e-01	4.738e+00	1.736e+00	3.091e-02	2.797e-01	1.459e+00	3.081e+00	7.083e-01	7.305e-01

Table 11: Relative gain in mean fitness of 5 fittest expressions on training data after 10 phases.

	0	1	2	3	4	5	6	7	8	9	10	11	12	13	14
ABC	3.530e+00	4.032e-01	8.301e-01	1.346e+01	3.503e+11	8.400e+01	1.061e+00	2.112e-01	4.251e-01	1.531e-01	7.788e-01	1.034e+00	5.595e+00	1.041e+00	7.932e-01
DE	2.813e+05	2.333e+00	1.492e-01	6.914e+00	4.242e+03	3.098e+07	1.305e+01	4.957e+00	1.210e+00	2.547e-01	9.384e-01	1.339e+00	1.614e+00	2.333e+00	1.742e+00
PSO	8.346e+05	7.871e+00	7.495e-01	9.256e-01	1.573e+04	1.539e+00	3.743e+00	4.297e+00	3.560e+00	1.781e-01	7.192e-01	1.838e+00	6.612e+00	1.991e+00	3.372e+00

Table 12: Relative gain in mean fitness of 5 fittest expressions on full data after 10 phases.

	0	1	2	3	4	5	6	7	8	9	10	11	12	13	14
ABC	3.389e+00	4.674e-01	1.263e+00	1.779e+01	3.846e+11	2.257e+02	1.163e+00	2.171e-01	4.563e-01	1.408e-01	3.961e-01	1.020e+00	3.568e+00	1.170e+00	1.027e+00
DE	2.901e+05	1.426e+00	1.270e-02	7.677e+00	4.445e+03	8.874e+07	1.310e+01	3.143e+00	1.057e+00	2.200e-01	9.597e-01	5.134e-01	8.172e-01	1.426e+00	1.039e+00
PSO	8.487e+05	2.720e-01	1.113e+00	8.726e-01	1.601e+04	3.942e+00	4.913e-01	4.000e+00	1.855e+00	1.070e-01	3.816e-01	1.327e+00	1.873e-01	1.466e+00	8.262e-01

6.4.7 Cost

The cost in terms of evaluations is linear in the number of phases. For each phase, 4 expressions are optimized with $O(nk)$ complexity. In our configuration $n=50$, $k=50$. The SR algorithm itself requires $O(mg)$ evaluations per phase, where m is the population (20) and g the generations per phase (20). The question whether or not the cost of the optimizer is worth the gain in fitness is in the end one for the practitioner to answer. There is no guarantee that improvement will take place, but from the results on the above hard problems we can still expect significant improvements up to several orders of magnitude.

6.5 Conclusion

6.5.1 Operator cooling schedule

6.5.2 Constant folding

6.5.3 Optimizers

The experiments with the optimizers highlight several issues. There are a wide number of strategies and parameters that influence the effect of the optimizer. We also see that optimizers can hinder the algorithm in its convergence. This is not a general conclusion, but dependent in part on our design choices and the trade-offs made. If we only optimize the final outcome of the algorithm it is obvious that no fitness regression is possible. Only when we apply optimization in the archiving stage are there subtle effects at play that allow for such edge cases. The cost of applying the optimizers is significant. In our implementation the cost is known beforehand, but the gain is not. This holds true in general for this GP SR algorithm. While we can empirically investigate the convergence of a number of problems, there is no known limit to this process. In general, when the optimizers are used the improvements made are far greater than the loss in edge cases. We have tested 3 distinct algorithms as optimizers in order to test which is best. Unfortunately no such algorithm exists. The NFL theorem [22] suggests as much. What we can see is that ABC and DE offer, for our problem set, the best results. Best in this context means the overall highest gain with the lowest losses at an equal cost to the other algorithms. The hardness of the SR problem indicates that, unfortunately, there are an infinite number of problem statements that will have different convergence characteristics. While hybridization of the GP algorithm with other algorithms is a viable strategy, it also substantially increases the number of parameters.

7 Use Case

In this section we will use our tool on a real world use case.

7.1 Problem statement

We want to use symbolic regression on the output of a simulator. The parameter space of the simulator is huge, and simulation instance are computationally expensive. We would like to obtain symbolic expressions relating parameters values to simulator output. These expressions can be used to gain insight into which parameters are correlated, which have a larger effect on output and which are irrelevant.

7.2 Design of Experiment

7.3 Incremental support

We will use the incremental support built in into CSRM to apply SR to the simulator in tandem.

7.4 Results

8 Related Work

9 Conclusion

10 Future work

References

- [1] BANSAL, J. C., SINGH, P., SARASWAT, M., VERMA, A., JADON, S. S., AND ABRAHAM, A. Inertia weight strategies in particle swarm optimization. In *Nature and Biologically Inspired Computing (NaBIC), 2011 Third World Congress on* (2011), IEEE, pp. 633–640.
- [2] BONYADI, M. R., AND MICHALEWICZ, Z. A locally convergent rotationally invariant particle swarm optimization algorithm. *Swarm Intelligence* 8, 3 (2014), 159–198.
- [3] CLERC, M., AND KENNEDY, J. The particle swarm - explosion, stability, and convergence in a multidimensional complex space. *Trans. Evol. Comp* 6, 1 (Feb. 2002), 58–73.
- [4] DAS, S., MULLICK, S. S., AND SUGANTHAN, P. Recent advances in differential evolution an updated survey. *Swarm and Evolutionary Computation* 27 (2016), 1 – 30.
- [5] EBERHART, R. C., AND SHI, Y. Comparing inertia weights and constriction factors in particle swarm optimization. In *Proc. of the 2000 Congress on Evolutionary Computation* (Piscataway, NJ, 2000), IEEE Service Center, pp. 84–88.

- [6] EVETT, M., AND FERNANDEZ, T. Numeric mutation improves the discovery of numeric constants in genetic programming. In *Genetic Programming*. In (1998), Morgan Kaufmann, pp. 66–71.
- [7] GANSNER, E. R., AND NORTH, S. C. An open graph visualization system and its applications to software engineering. *SOFTWARE - PRACTICE AND EXPERIENCE* 30, 11 (2000), 1203–1233.
- [8] GARDNER, M.-A., GAGN, C., AND PARIZEAU, M. Controlling code growth by dynamically shaping the genotype size distribution. *Genetic Programming and Evolvable Machines* 16, 4 (2015), 455–498.
- [9] JENKS, G., ET AL. SortedContainers: Sorted list, dictionary and set types for Python, 2014–. [Online; accessed ;today].
- [10] KARABOGA, D., AKAY, B., AND OZTURK, C. *Artificial Bee Colony (ABC) Optimization Algorithm for Training Feed-Forward Neural Networks*. Springer Berlin Heidelberg, Berlin, Heidelberg, 2007, pp. 318–329.
- [11] KARABOGA, D., OZTURK, C., KARABOGA, N., AND GORKEMLI, B. Artificial bee colony programming for symbolic regression. *Inf. Sci.* 209 (Nov. 2012), 1–15.
- [12] KENNEDY, J., AND EBERHART, R. Particle swarm optimization, 1995.
- [13] KENNEDY, J., AND EBERHART, R. C. Particle swarm optimization. In *Proceedings of the 1995 IEEE International Conference on Neural Networks* (Perth, Australia, IEEE Service Center, Piscataway, NJ, 1995), vol. 4, pp. 1942–1948.
- [14] KENNEDY, J. F., KENNEDY, J., EBERHART, R. C., AND SHI, Y. Swarm intelligence, 2001.
- [15] KORNS, M. F. *Accuracy in Symbolic Regression*. Springer New York, New York, NY, 2011, pp. 129–151.
- [16] KORNS, M. F. *A Baseline Symbolic Regression Algorithm*. Springer New York, New York, NY, 2013, pp. 117–137.
- [17] KOZA, J. R. *Genetic Programming: On the Programming of Computers by Means of Natural Selection*. MIT Press, Cambridge, MA, USA, 1992.
- [18] LI, Z., WANG, W., YAN, Y., AND LI, Z. Psabc: A hybrid algorithm based on particle swarm and artificial bee colony for high-dimensional optimization problems. *Expert Systems with Applications* 42, 22 (2015), 8881 – 8895.
- [19] O’NEILL, M., AND BRABAZON, A. Grammatical differential evolution. In *Proceedings of the 2006 International Conference on Artificial Intelligence, ICAI 2006* (Las Vegas, Nevada, USA, June 26-29 2006), H. R. Arabnia, Ed., vol. 1, CSREA Press, pp. 231–236.
- [20] SPEARS, W. M., GREEN, D. T., AND SPEARS, D. F. Biases in particle swarm optimization. *Int. J. Swarm. Intell. Res.* 1, 2 (Apr. 2010), 34–57.

- [21] STORN, R., AND PRICE, K. Differential evolution – a simple and efficient heuristic for global optimization over continuous spaces. *Journal of Global Optimization* 11, 4 (1997), 341–359.
- [22] WOLPERT, D. H., AND MACREADY, W. G. No free lunch theorems for optimization. *Trans. Evol. Comp* 1, 1 (Apr. 1997), 67–82.