

# Probabilistic Grammatical Evolution

Jessica Mégane, Nuno Lourenço, and Penousal Machado

CISUC, Department of Informatics Engineering, University of Coimbra,  
Polo II - Pinhal de Marrocos, 3030 Coimbra, Portugal  
jessicac@student.dei.uc.pt, {nam1, penousal}@dei.uc.pt

**Abstract.** Grammatical Evolution (GE) is one of the most popular Genetic Programming (GP) variants, and it has been used with success in several problem domains. Since the original proposal, many enhancements have been proposed to GE in order to address some of its main issues and improve its performance.

In this paper we propose Probabilistic Grammatical Evolution (PGE), which introduces a new genotypic representation and new mapping mechanism for GE. Specifically, we resort to a Probabilistic Context-Free Grammar (PCFG) where its probabilities are adapted during the evolutionary process, taking into account the productions chosen to construct the fittest individual. The genotype is a list of real values, where each value represents the likelihood of selecting a derivation rule. We evaluate the performance of PGE in two regression problems and compare it with GE and Structured Grammatical Evolution (SGE).

The results show that PGE has a better performance than GE, with statistically significant differences, and achieved similar performance when comparing with SGE.

**Keywords:** Genetic Programming, Grammatical Evolution, Probabilistic Context-Free Grammar, Probabilistic Grammatical Evolution, Genotype-to-Phenotype Mapping

## 1 Introduction

Evolutionary Algorithms (EAs) are loosely inspired by the ideas of natural evolution, where a population of individuals evolves through the application of selection, variation (such as crossover and mutation) and reproduction operators. The evolution of these individuals is guided by a fitness function, which measures the quality of the solutions that each individual represents to the problem at hand. The application of these elements is repeated for several iterations and it is expected that, over time, the quality of individuals improves.

Genetic Programming (GP) [1] is an EA that is used to evolve programs. Over the years many variants of GP have been proposed, namely concerned with how individuals (i.e., computer programs) are represented. Some of these variants make use of grammars to enforce syntactic constraints on the individual solutions. The most well known grammar-based GP approaches are Context-free Grammar Genetic Programming (CFG-GP), introduced by Whigham in

[2], and Grammatical Evolution (GE) introduced by Ryan *et al.* [3]. The main distinction between the two approaches is the representation of the individual’s solution (genotype) in the search space. CFG-GP uses a derivation-tree based representation, and the mapping is made by reading the terminal symbols (tree leaves), starting from the left leaf to the right. In GE there is a distinction between the genotype, a variable length string of integers, and the phenotype of the individual. The mapping between the genotype and the phenotype is performed through a Context-Free Grammar (CFG).

GE is one of the most popular GP variants, in spite of the debate in the literature [4] concerning its relative performance when compared to other grammar-based variants. To address some of the main criticisms of GE, several improvements have been proposed in the literature related to the population initialisation [5], grammar design [6] and the representation of individuals [7,8,9,10].

In this paper we introduce a new representation to GE. In concrete, we proposed a new probabilistic mapping mechanism to GE, called Probabilistic Grammatical Evolution (PGE). In PGE the genotype is a list of probabilities and the mapping is made using a Probabilistic Context-Free Grammar (PCFG) to choose the productions of the individual’s phenotype. All derivation rules in the grammar start with the same chance of being selected, but over the evolutionary process, the probabilities are updated considering the derivation rules that were selected to build the fittest individual. To evaluate the performance of PGE, we compare its performance with GE and SGE [11] in two benchmark problems. PGE showed statistically significant improvements when compared with GE and obtained similar performance when compared to SGE.

The remainder of the paper is structured as follows: Section 2 introduces Grammatical Evolution and related work. Section 3 describes our approach called Probabilistic Grammatical Evolution (PGE), Section 4 details the experimental framework used to study the performance of PGE, and Section 5 describes the main results. Finally, Section 6 gathers the main conclusions and provides some insights towards future work.

## 2 Grammatical Evolution

GE [3] is a Grammar-based GP approach where the individuals are presented as a variable length string of integers. To create an executable program, the genotype (i.e., the string of integers) is mapped to the phenotype (program) via the productions rules defined in a Context-Free Grammar (CFG). A grammar is a tuple  $G = (NT, T, S, P)$  where  $NT$  and  $T$  represent the non-empty set of *Non-Terminal* (NT) and *Terminal* (T) symbols,  $S$  is an element of  $NT$  called the axiom and  $P$  is the set of production rules. The rules in  $P$  are in the form  $A ::= \alpha$ , with  $A \in NT$  and  $\alpha \in (NT \cup T)^*$ . The  $NT$  and  $T$  sets are disjoint. Each grammar defines a language  $L(G) = \{w : S \xrightarrow{*} w, w \in T^*\}$ , that is the set of all sequences of terminal symbols that can be derived from the axiom.

The genotype-phenotype mapping is the key issue in GE, and it is performed in several successive steps. To select which derivation rule should be selected to

replace a NT, the mapping relies on the modulo operator. An example of this process is shown in Fig. 1. The genotype is composed of integers values randomly generated between  $[0,255]$ . The mapping starts with the axiom  $\langle start \rangle$ . In this case, there is only one derivation possible, and we rewrite the axiom with  $\langle expr \rangle$ . Then we proceed the expansion of  $\langle expr \rangle$ . Since this NT has two possible expansion rules available, we need to select which one will be used. We start by taking the first unused value of the genotype, which is 54, and divide it by the number of possible options. The remainder of this operation indicates the option that should be used. In our example,  $54 \bmod(2) = 0$ , which results in the first production being selected. This process is performed until there are no more NT symbols to expand or there are no more integers to read from the genotype.

In this last case and if we still have NT to expand, a wrapping mechanism can be used, where the genotype will be re-used, until it generates a valid individual or the predefined number of wraps is over. If after all the wraps we still have not mapped all the NT, the mapping process stops, and the individual will be considered invalid.

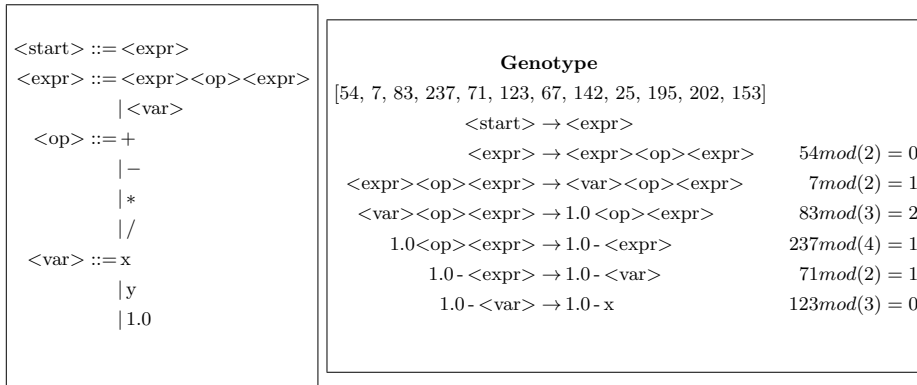


Fig. 1: Example of GE mapping

Even though GE has been applied to several problem domains, there is a debate in the literature concerned with its overall performance [4,12]. GE has been criticised for having high redundancy and low locality [13,14]. A representation has high redundancy when several different genotypes correspond to one phenotype. Locality is concerned with how changes in the genotype are reflected on the phenotype. These criticisms have triggered many researchers into looking how GE could be improved [15,7,16,11].

## 2.1 Representation Variants

O'Neill *et al.* [7] proposed Position Independent GE ( $\pi$ GE), introducing a different mapping mechanism that removes the positional dependency that exists

in GE. In  $\pi$ GE each codon is composed of two values (nont, rule), where nont is used to select the next non-terminal to be expanded and the rule selects which production rule should be applied from the selected non-terminal. In Fagan *et al.* [17] several different mapping mechanisms were compared, and  $\pi$ GE showed better performance over GE, with statistical differences. Another attempt to make GE position independent is Chorus [15]. In this variant, each gene encodes a specific grammar production rule, not taking into consideration the position. This proposal did not show significant differences when comparing with standard GE.

Structured Grammatical Evolution is a recent proposal to address the locality and redundancy problems of GE [11,10]. SGE proposes a one-to-one mapping between the genotype and the non-terminal symbols of the grammar. Each position in the genotype of SGE is a list of integers, where each element of this list is used to select the next derivation rule. This genotype structure, ensures that the modification of a codon does not affect the chosen productions of other non-terminals, reducing the overall changes that can occur at the phenotypical level, which results in a higher locality.

In [11] different grammar-based GP approaches were compared, and the authors showed that SGE achieved a good performance when compared with several grammar-based GP representations. These results were in line with Fagan *et al.* [17], which showed that different genotype-phenotype mapping can improve the performance of grammar-based GP.

Some probabilistic methods have been proposed to try to understand more about the distribution of fitter individuals and have been effective in solving complex problems [18]. Despite its potential, few attempts have been made to use probabilities in GE.

In [8], was implemented a PCFG (Figure 2) to do the mapping process of GE, where the genotype of the individual is a vector of probabilities used to choose the productions rules. This approach also implements Estimation of Distribution Algorithms (EDA) [19], a probabilistic technique that replaces the mutation and crossover operators, by sampling the probability distribution of the best individuals, to generate the new population, each generation. The probabilities start all equal and are updated each generation, based on the frequency of the chosen rules of the individuals with higher fitness. The experiments were inclusive, since the proposed approach had a similar performance to GE.

Kim *et al.* [9] proposed Probabilistic Model Building Grammatical Evolution (PMBGE), which utilises a Conditional Dependency Tree (CDT) that represents the relationships between production rules used to calculate the new probabilities. Similar to [8], the EDA mechanism was implemented instead of the genetic operators. The results showed no statistical differences between GE and the proposed approach.

### 3 Probabilistic Grammatical Evolution

Probabilistic Grammatical Evolution (PGE) is a new representation for Grammatical Evolution. In PGE we rely on a Probabilistic Context-Free Grammar (PCFG) to perform the genotype-phenotype mapping. A PCFG is a quintuple  $PG = (NT, T, S, P, Probs)$  where  $NT$  and  $T$  represent the non-empty set of *Non-Terminal* (NT) and *Terminal* (T) symbols, respectively,  $S$  is an element of  $NT$  called the axiom,  $P$  is the set of production rules, and  $Probs$  is a set of probabilities associated with each production rule. The genotype in PGE is a vector of floats, where each element corresponds to the probability of choosing a certain derivation rule. The overall mapping procedure is shown in Alg. 1 and Fig. 2 shows an example of the application of the PGE mapping.

The panel on the left shows a PCFG, where each derivation rule has a probability associated. The set of NT is composed of  $\langle start \rangle$ ,  $\langle expr \rangle$ ,  $\langle op \rangle$  and  $\langle var \rangle$ . The right panel of Fig. 2 shows how the mapping procedure works.

$\langle start \rangle ::= \langle expr \rangle$ (1.0) $\langle expr \rangle ::= \langle expr \rangle \langle op \rangle \langle expr \rangle$ (0.5) $\quad   \langle var \rangle$ (0.5) $\langle op \rangle ::= +$ (0.33) $\quad   *$ (0.33) $\quad   -$ (0.33) $\langle var \rangle ::= x$ (0.5) $\quad   1.0$ (0.5)	<p style="text-align: center;"><b>Genotype</b></p> <p style="text-align: center;">[0.8, 0.2, 0.98, 0.45, 0.62, 0.37, 0.19]</p> $\langle start \rangle \rightarrow \langle expr \rangle$ (0.8) $\langle expr \rangle \rightarrow \langle expr \rangle \langle op \rangle \langle expr \rangle$ (0.2) $\langle expr \rangle \langle op \rangle \langle expr \rangle \rightarrow \langle var \rangle \langle op \rangle \langle expr \rangle$ (0.98) $\langle var \rangle \langle op \rangle \langle expr \rangle \rightarrow x \langle op \rangle \langle expr \rangle$ (0.45) $x \langle op \rangle \langle expr \rangle \rightarrow x * \langle expr \rangle$ (0.62) $x * \langle expr \rangle \rightarrow x * \langle var \rangle$ (0.37) $x * \langle var \rangle \rightarrow x * x$ (0.19)
---	---

Fig. 2: Example of mapping with PCFG

It begins with the axiom,  $\langle start \rangle$ . We start by taking the first value of the genotype, which is 0.8, and since there is only one expansion available, the non-terminal  $\langle expr \rangle$  will be chosen. Next, we need to rewrite  $\langle expr \rangle$ , which has two derivation options. We take the second value of the genotype, 0.2, and compare it to the probability associated with the first derivation option ( $\langle expr \rangle \langle op \rangle \langle expr \rangle$ ). Since  $0.2 < 0.5$ , we select this derivation option to rewrite  $\langle expr \rangle$ . This process is repeated until there are no more non-terminals left to expand, or no probabilities left in the genotype. When this last situation occurs, we use a wrapping mechanism similar to the standard GE, where the genotype will be reused a certain number of times. If after the wrapping we still have not mapped the individual completely, the mapping process stops, and the individual will be considered invalid.

---

**Algorithm 1** Mapping with PCFG

---

```
1: procedure GENERATEINDIVIDUAL(genotype, pcfg)
2:   start = pcfg.getStart()
3:   phenotype = [start]
4:   for codon in genotype do
5:     symbol = phenotype.getNextNT()
6:     productions = pcfg.getRulesNT(symbol)
7:     cum_prob = 0.0 ▷ Cumulative Sum of Probabilities
8:     for prod in productions do
9:       cum_prob = cum_prob + prod.getProb()
10:      if codon < cum_prob then
11:        selected_rule = prod
12:        break
13:      end if
14:    end for
15:    phenotype.replace(symbol, selected_rule)
16:    if phenotype.isValid() then
17:      break
18:    end if
19:  end for
20: end procedure
```

---

In PGE, the probabilities are updated each generation after evaluating the population, based on how many times each derivation rule has been selected by the best individual of the current generation or the best individual overall. When a derivation rule is used to create one of these individuals, its probability should be increased, otherwise if a derivation rule is not used, we should decrease it. Alternating between these two bests helps us to avoid using the same individual in consecutive generations to adjust the probabilities of the PCFG, balancing global exploration with local exploitation. All the adjustments are performed using a parameter  $\lambda$  called *learning factor* which smooths the transitions on the search space. The lambda value should be between 0% and 100%. At each generation, each individual is mapped using an updated version of grammar.

To update the probabilities in the grammar, we use Alg. 2, where  $j$  is the number of productions of a non-terminal symbol of the grammar,  $i$  is the index of the production probability that is being updated and  $\lambda$  is the learning factor.

The probabilities are updated based on two different rules. The first rule increases the probability of a derivation rule, taken into account the frequency that it was selected by the best individual (Alg. 2 line 5). The second rule decreases the probability of the derivation options that are never used to expand a non-terminal (Alg. 2 line 7). The *min* operator ensures that when we update the probabilities they are not greater than 1.

After the update of the probabilities for each derivation rule, we make sure that the sum of the probabilities of all derivation rules, for a non-terminal, is 1. If the sum surpasses this value, the excess is proportionally subtracted from the

---

**Algorithm 2** Probabilistic Grammatical Evolution

---

```
1: procedure UPDATEPROBABILITIES(best)
2:   counter = best.getCounter()      ▷ list with times each rule was expanded
3:   for each production rule i of each NT do
4:     if counteri > 0 then
5:        $prob_i = \min(prob_i + \frac{\lambda * counter_i}{\sum_{k=1}^j counter_k}, 1.0)$ 
6:     else
7:        $prob_i = prob_i - \lambda * prob_i$ 
8:     end if
9:   end for
10:  while  $\sum_{k=1}^j prob_i \neq 1.0$  do
11:     $extra = (1.0 - \sum_{k=1}^j prob_i) / j$ 
12:    for each production rule i do
13:       $prob_i = prob_i + extra$ 
14:    end for
15:  end while
16: end procedure
```

---

derivation options for a non-terminal. If the sum is smaller than one, the missing amount is added equally to the production rules of the non-terminal.

We are going to use the example of grammar and individual presented Fig. 2, to show how the probabilities are updated. On the left panel of Fig. 3 we can see the derivation tree of the individual that was used to update the probabilities of the PCFG on the right. The learning factor used was 0.01 (1%).

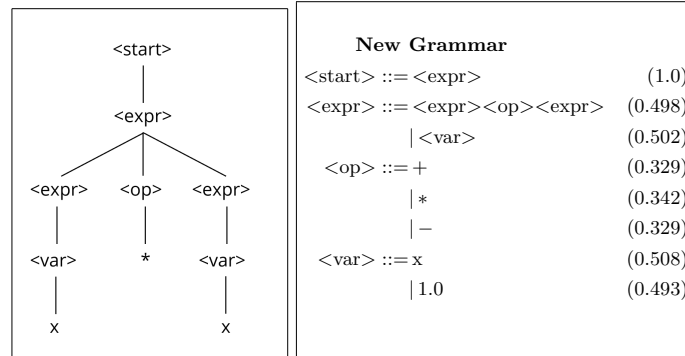


Fig. 3: Example of probability updating in PGE

Since the symbol *< start >* has only one expansion rule, the probability of choosing *< expr >* stays always 1. Looking at the derivation tree, we can see that the first derivation rule (*< expr >< op >< expr >*) of the non-terminal

$\langle expr \rangle$  was selected once, and the second derivation rule ( $\langle var \rangle$ ) was selected twice.

Using the Alg. 2, the new probabilities for the first derivation option is  $\min(0.5 + \frac{0.01*1}{3}, 1)$  equals 0.5033(3), and for the second rule,  $\min(0.5 + \frac{0.01*2}{3}, 1)$  equals 0.5066(6). As the sum of the probabilities surpass 1, the excess  $((1 - 0.5066 + 0.5033)/2 = 0.00495)$  is subtracted from both probabilities and the value is rounded to 3 decimal places. Then we distribute the excess for the derivations rules: the first rule is updated to 0.498 and the second rule is updated to 0.502. This process is applied to the other symbols. For the non-terminal  $\langle op \rangle$  the rule  $+$  and  $-$  were never chosen so they will update equally  $((0.33 - 0.01 * 0.33)$  that equals 0.3267), and the rule  $*$  was chosen once ( $\min(0.33 + \frac{0.01*1}{1}, 1)$  that equals 0.34). As the sum of the three probabilities is smaller than one, 0.0022 must be added to the three probabilities, and the result rounded, staying with 0.329 for the  $+$  and  $-$  symbols, and 0.342 for the  $*$ . The non-terminal  $\langle var \rangle$  was expanded twice for the terminal  $x$  and never expanded for the terminal 1.0. By applying the algorithm, the first rule ( $x$ ) should be updated to 0.51 ( $\min(0.5 + \frac{0.01*2}{2}, 1)$ ) and the second to 0.495 ( $0.5 - 0.01 * 0.5$ ), being the sum of the two probabilities different than one, the adjustment and rounding should be done and they are updated to 0.508 and 0.493, respectively.

## 4 Experimental Analysis

Over the years, several genotype-phenotype mapping alternatives have been proposed to increase the performance of GE. One that has obtained promising results is Structured Grammatical Evolution (SGE) [10,11]. To evaluate the performance of PGE, we considered the standard GE and SGE algorithms in two different benchmark problems. These benchmark problems were selected taking into account the comparative analysis followed by [11] and the recommendations presented in [20]. For our experimental analysis we selected the Pagie Polynomial and the Boston Housing prediction problem.

### 4.1 Problem Description

**Pagie Polynomial** Popular benchmark problem for testing Genetic Programming algorithms, with the objective of finding the mathematical expression for the following problem:

$$\frac{1}{1 + x^{-4}} + \frac{1}{1 + y^{-4}}. \quad (1)$$

The function is sampled in the interval  $[-5, 5.4]$  with a step of 0.4, and the grammar used is:

```

 $\langle start \rangle ::= \langle expr \rangle$ 
 $\langle expr \rangle ::= \langle expr \rangle \langle op \rangle \langle expr \rangle \mid ( \langle expr \rangle \langle op \rangle \langle expr \rangle )$ 
            $\mid \langle pre\_op \rangle ( \langle expr \rangle ) \mid \langle var \rangle$ 

```



$\langle op \rangle ::= + \mid - \mid * \mid /$   
 $\langle pre\_op \rangle ::= \sin \mid \cos \mid \exp \mid \log \mid \text{inv}$   
 $\langle var \rangle ::= x \mid y \mid 1.0$

The division and logarithm functions are protected, i.e.,  $1/0 = 1$  and  $\log(f(x)) = 0$  if  $f(x) \leq 0$ .

**Boston Housing** This is a famous Machine Learning problem to predict the prices of Boston Houses. The dataset comes from the StatLib Library [21] and has 506 entries, with 13 features. It was divided in 90% for training and 10% for testing. The grammar used for the Boston Housing regression problem is as follows:

$\langle start \rangle ::= \langle expr \rangle$   
 $\langle expr \rangle ::= \langle expr \rangle \langle op \rangle \langle expr \rangle \mid ( \langle expr \rangle \langle op \rangle \langle expr \rangle )$   
 $\quad \mid \langle pre\_op \rangle ( \langle expr \rangle ) \mid \langle var \rangle$   
 $\langle op \rangle ::= + \mid - \mid * \mid /$   
 $\langle pre\_op \rangle ::= \sin \mid \cos \mid \exp \mid \log \mid \text{inv}$   
 $\langle var \rangle ::= x[1] \mid \dots \mid x[13] \mid 1.0$

## 4.2 Parameters

For all the experiments reported, the fitness function is computed using the Root Relative Squared Error (RRSE), where lower values indicate a better fitness. The parameters are presented in Table 1. These parameters were selected in order to make the comparisons between all the approaches fair. Additionally, and to avoid side effects, the wrapping mechanism was removed from GE and PGE. Concerning the variation operators for PGE, we used the standard one-point crossover, and float mutation which generates a new random float between  $[0,1]$ . Additionally, PGE uses a learning factor of  $\lambda = 1.0\%$ .

## 5 Results

The experimental results in this section will be presented in terms of the mean best fitness, which results from the execution of 100 independent runs. To compare all approaches we performed a statistical study. Since the results did not follow any distribution, and the populations were independently initialised, we employed the Kruskal-Wallis non-parametric test to check if there were meaningful differences between the different groups of approaches. When this happened we used the Mann-Whitney *post-hoc* test with Bonferroni correction. For all the statistical tests we considered a significance level  $\alpha = 0.05$ .

Table 1: Parameters used in the experimental analysis for GE, PGE and SGE

Parameters	Value		
	GE	PGE	SGE
Population Size	1000		
Generations	50		
Elitism	10%		
Mutation Probability	5%		
Crossover Probability	90%		
Tournament	3		
Max Number of Wraps	0	-	
Size of Genotype	128	-	
Max. Initialisation depth	-	6	
Max. Tree Depth	-	17	

Fig. 4 depicts the results for the Pagie Polynomial. It is possible to see that all the methods start from similar fitness values, but as the evolutionary process progresses, differences between the approaches emerge. The fitness of the solutions being evolved by SGE rapidly decrease in the early generations, but slows down after a certain number of generations ( $\approx 20$ ). For GE, the fitness decreases slowly through the generations. This results are in line with the ones presented in [10].

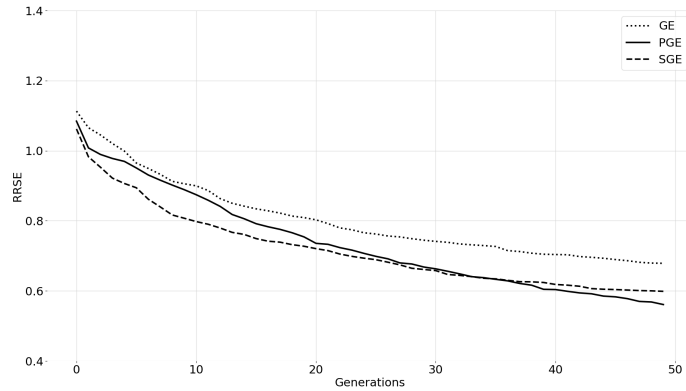


Fig. 4: Results for the Pagie Polynomial

The PGE performance is different from both SGE and GE. In the early generations PGE decreases slowly (following the same trend as GE), but then, around the 10th generation, the fitness of PGE starts to rapidly decrease. Around generation 35 it surpasses the quality of SGE. The first row of Table 2 shows

the Mean Best Fitness and the Standard Deviation for each approach. We can see that for the Pagie Polynomial problem PGE obtains the lowest error.

Table 2: Mean Best Fitness and Standard Deviation for all the methods used in the comparison. Results are averages of 100 independent runs.

Problem	PGE	GE	SGE
Pagie Polynomial	<b>0.56±0.16</b>	0.68±0.17	0.59±0.13
Boston Housing Train	0.82±0.12	0.88±0.14	<b>0.78±0.13</b>
Boston Housing Test	0.84±0.13	0.90±0.15	<b>0.79±0.12</b>

In terms of statistical significant differences, the Kruskal-Wallis showed meaningful differences between the approaches. The *post-hoc* results are depicted in Table 3. Looking at the results it is possible to see that both PGE and SGE are better than GE with statistical significant differences. When comparing PGE and SGE we only found marginal differences (p-value = 0.04) on the Boston Housing Training.

Table 3: Results of the Mann-Whitney *post-hoc* statistical tests. Bonferroni correction is used and the significance level  $\alpha = 0.05$  is considered.

	PGE - GE	PGE - SGE
Pagie Polynomial	<b>0.00</b>	0.24
Boston Housing Training	<b>0.00</b>	0.04
Boston Housing Test	<b>0.00</b>	0.10

Fig. 5 shows the results for the Boston Housing problem. Looking at the training results (Fig. 5 (a)), it is possible to see that the fitness of SGE individuals rapidly decrease, and continue too over the entire evolutionary process. The performance of GE is in line with what we observed previously, i.e., a slow decrease on the fitness. Even though the training results are important to understand the behaviour of the methods, the testing results are more relevant, because they allow us to evaluate the generalisation ability of the models evolved by each approach. Looking at the test results (Fig. 5(b)) we can see that SGE and PGE are building models that can generalise better to unseen data.

Once again we applied a statistical analysis to check whether there were differences between the approaches (Table 3). The results, for both training and test, show that PGE is statistically significant than GE, but there are no differences between PGE and SGE.

Finally we present an analysis on how the probabilities of certain derivation rules progress over the generations. This analysis will give us insights into what are the rules that are more relevant.

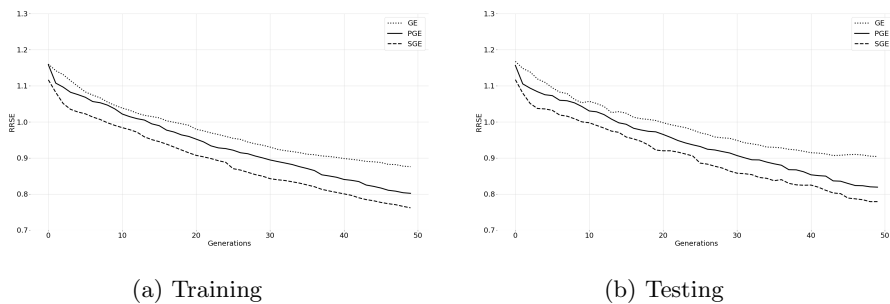


Fig. 5: Results for the Boston Housing problem

For the Page Polynomial, Fig. 6 presents the evolution of the PCFG’s probabilities for the non-terminal  $\langle op \rangle$  over the generations. As one would expect, the probabilities associated with the symbols that are required to solve the problem are higher, namely the ones associated with the terminal symbols  $+$  and  $/$ .

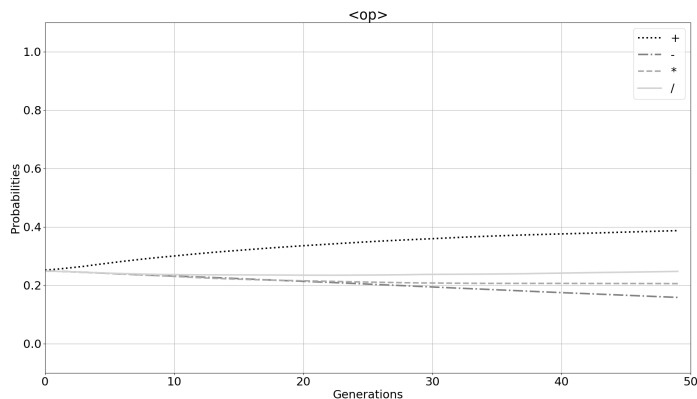


Fig. 6: Evolution of grammar probabilities of non-terminal  $\langle op \rangle$  with the Page Polynomial. Results are averages of 100 runs.

Concerning the Boston Housing, the progression of the probabilities are depicted in Fig. 7 and in Table 4. Concretely, the results show the probabilities for the derivation options of the non-terminal  $\langle var \rangle$ . This symbol was selected because it contains the features that describe the problem. Looking at the evolution of the probabilities associated with each production, we can understand which of these features are more relevant to accurately predict the price of houses. Looking at the results (Fig. 7), one can see that PTRATIO stands out

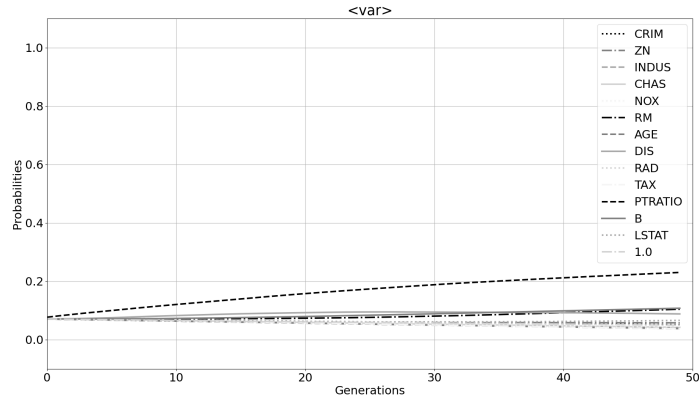


Fig. 7: Evolution of grammar probabilities of non-terminal  $\langle var \rangle$  with the Boston Housing problem. Results are averages of 100 runs.

in terms of the probability of being selected. PTRATIO represents the pupil-teacher ratio by town. This is in line with the results reported by [22]. Another interesting result is to see that the feature RM (the third most important feature in [22]), which is the average number of rooms per dwelling, is also on the top three of our results. These results confirm not only the relevance of these features to the Boston Housing problem, but also allow us to perform feature selection and provide an explanation to the results achieved. This means that at the end of the evolutionary process one can look at the final distribution of the probabilities in the grammar, and analyse the relative importance of each production and derivation rules, and see how they are used to create the best models.

Table 4: Probabilities of the Boston Housing Dataset’s productions of the non-terminal  $\langle var \rangle$  at the end of the evolutionary process. Results are averages of 100 runs.

Production	Probability
PTRATIO	0.23
B	0.11
RM	0.1
DIS	0.09
LSTAT	0.07
ZN	0.06
NOX, CRIM, 1.0	0.05
RAD, TAX, RADIUS, CHAS, AGE	0.04

## 6 Conclusion

Grammatical Evolution (GE) has attracted the attention of many researchers and practitioners. Since its proposal in the late 1990s, it has been applied with success to many problem domains. However, it has been shown that it suffers from some issues.

In this paper we proposed a new mapping mechanism and genotypic representation called Probabilistic Grammatical Evolution (PGE). In concrete, in PGE the genotype of an individual is a variable length sequence of floats, and the genotype-phenotype mapping is performed using a Probabilistic Context-Free Grammar (PCFG). Each derivation rule has a probability associated, and they are updated with taking into account the number of times that the derivation rules were selected by the best individuals. In order to maintain a balance between global and local exploration we alternate between the best overall individual and the best individual of generation, respectively.

PGE was compared with standard GE and SGE in two different benchmarks. The results show that for both problems PGE is statistically better than GE and has a similar performance when compared to SGE. We also analyse how the probabilities associated with the different productions progress over time, and it was possible to see that the production rules that are more relevant to the problem at hand have higher probabilities of being selected.

In terms of future work one needs to consider alternative mechanisms to adjust probabilities of the production rules. Another line of work that needs to be conducted is concerned with the analysis of the locality and redundancy in PGE.

## Acknowledgements

This work is partially funded by the project grant DSAIPA/DS/0022/2018 (GADgET), by national funds through the FCT - Foundation for Science and Technology, I.P., within the scope of the project CISUC - UID/CEC/00326/2020 and by European Social Fund, through the Regional Operational Program Centro 2020. We also thank the NVIDIA Corporation for the hardware granted to this research.

## References

1. Koza, J.R.: Genetic Programming: On the Programming of Computers by Means of Natural Selection. MIT Press, Cambridge, MA, USA (1992)
2. Whigham, P.A.: Grammatically-based genetic programming. In: Proceedings of the Workshop on Genetic Programming: From Theory to Real-World Applications, vol. 16. pp. 33–41 (1995)
3. Ryan, C., Collins, J.J., O’Neill, M.: Grammatical evolution: Evolving programs for an arbitrary language. In: Lecture Notes in Computer Science, pp. 83–96. Springer Berlin Heidelberg (1998), <https://doi.org/10.1007/bfb0055930>

4. Whigham, P.A., Dick, G., Maclaurin, J., Owen, C.A.: Examining the best of both worlds of grammatical evolution. In: Proceedings of the 2015 Annual Conference on Genetic and Evolutionary Computation. pp. 1111–1118. ACM (Jul 2015), <https://doi.org/10.1145/2739480.2754784>
5. Nicolau, M.: Understanding grammatical evolution: initialisation. Genetic Programming and Evolvable Machines 18(4), 467–507 (Jul 2017), <https://doi.org/10.1007/s10710-017-9309-9>
6. Nicolau, M., Agapitos, A.: Understanding grammatical evolution: Grammar design. In: Handbook of Grammatical Evolution, pp. 23–53. Springer International Publishing (2018), [https://doi.org/10.1007/978-3-319-78717-6\\_2](https://doi.org/10.1007/978-3-319-78717-6_2)
7. O’Neill, M., Brabazon, A., Nicolau, M., McGarraghy, M., Keenan, P.:  $\pi$ grammatical evolution. In: Genetic and Evolutionary Computation – GECCO 2004, pp. 617–629. Springer Berlin Heidelberg (2004), [https://doi.org/10.1007/978-3-540-24855-2\\_70](https://doi.org/10.1007/978-3-540-24855-2_70)
8. Kim, H.T., Ahn, C.W.: A new grammatical evolution based on probabilistic context-free grammar. In: Proceedings in Adaptation, Learning and Optimization, pp. 1–12. Springer International Publishing (2015), [https://doi.org/10.1007/978-3-319-13356-0\\_1](https://doi.org/10.1007/978-3-319-13356-0_1)
9. Kim, H.T., Kang, H.K., Ahn, C.W.: A conditional dependency based probabilistic model building grammatical evolution. IEICE Transactions on Information and Systems E99.D(7), 1937–1940 (2016), <https://doi.org/10.1587/transinf.2016ed18004>
10. Lourenço, N., Pereira, F.B., Costa, E.: Unveiling the properties of structured grammatical evolution. Genetic Programming and Evolvable Machines 17(3), 251–289 (Feb 2016), <https://doi.org/10.1007/s10710-015-9262-4>
11. Lourenço, N., Ferrer, J., Pereira, F.B., Costa, E.: A comparative study of different grammar-based genetic programming approaches. In: Lecture Notes in Computer Science, pp. 311–325. Springer International Publishing (2017), [https://doi.org/10.1007/978-3-319-55696-3\\_20](https://doi.org/10.1007/978-3-319-55696-3_20)
12. Ryan, C.: A rebuttal to whigham, dick, and maclaurin by one of the inventors of grammatical evolution: Commentary on on the mapping of genotype to phenotype in evolutionary algorithms by peter a. whigham, grant dick, and james maclaurin. Genetic Programming and Evolvable Machines 18(3), 385–389 (Feb 2017), <https://doi.org/10.1007/s10710-017-9294-z>
13. Keijzer, M., O’Neill, M., Ryan, C., Cattolico, M.: Grammatical evolution rules: The mod and the bucket rule. In: Lecture Notes in Computer Science, pp. 123–130. Springer Berlin Heidelberg (2002), [https://doi.org/10.1007/3-540-45984-7\\_12](https://doi.org/10.1007/3-540-45984-7_12)
14. Rothlauf, F., Oetzel, M.: On the locality of grammatical evolution. In: Lecture Notes in Computer Science, pp. 320–330. Springer Berlin Heidelberg (2006), [https://doi.org/10.1007/11729976\\_29](https://doi.org/10.1007/11729976_29)
15. Ryan, C., Azad, A., Sheahan, A., O’Neill, M.: No coercion and no prohibition, a position independent encoding scheme for evolutionary algorithms – the chorus system. In: Lecture Notes in Computer Science, pp. 131–141. Springer Berlin Heidelberg (2002), [https://doi.org/10.1007/3-540-45984-7\\_13](https://doi.org/10.1007/3-540-45984-7_13)
16. Bartoli, A., Castelli, M., Medvet, E.: Weighted hierarchical grammatical evolution. IEEE Transactions on Cybernetics 50(2), 476–488 (Nov 2018), <https://doi.org/10.1109/tcyb.2018.2876563>
17. Fagan, D., O’Neill, M., Galván-López, E., Brabazon, A., McGarraghy, S.: An analysis of genotype-phenotype maps in grammatical evolution. In: Lecture Notes in Computer Science, pp. 62–73. Springer Berlin Heidelberg (2010), [https://doi.org/10.1007/978-3-642-12148-7\\_6](https://doi.org/10.1007/978-3-642-12148-7_6)

18. Kim, K., Shan, Y., Hoai, N.X., McKay, R.I.: Probabilistic model building in genetic programming: a critical review. *Genetic Programming and Evolvable Machines* 15(2), 115–167 (Sep 2013), <https://doi.org/10.1007/s10710-013-9205-x>
19. Larrañaga, P., Lozano, J.A. (eds.): *Estimation of Distribution Algorithms: A New Tool for Evolutionary Computation*. Springer US (2002), <https://doi.org/10.1007/978-1-4615-1539-5>
20. McDermott, J., De Jong, K., O'Reilly, U., White, D.R., Luke, S., Manzoni, L., Castelli, M., Vanneschi, L., Jaskowski, W., Krawiec, K., Harper, R.: Genetic programming needs better benchmarks. In: *Proceedings of the fourteenth international conference on Genetic and evolutionary computation conference - GECCO'12*. ACM Press (2012), <https://doi.org/10.1145/2330163.2330273>
21. Harrison, D., Rubinfeld, D.: Boston Housing Data. <http://lib.stat.cmu.edu/datasets/boston> (1993), [Online; accessed 27-December-2020]
22. Che, J., Yang, Y., Li, L., Bai, X., Zhang, S., Deng, C.: Maximum relevance minimum common redundancy feature selection for nonlinear data. *Information Sciences* 409-410, 68–86 (Oct 2017), <https://doi.org/10.1016/j.ins.2017.05.013>