

Gene Expression Programming*

Jason Brownlee
jasonb@CleverAlgorithms.com
The Clever Algorithms Project
<http://www.CleverAlgorithms.com>

April 4, 2010
Technical Report: CA-TR-20100404-1

Abstract

The Clever Algorithms project aims to describe a large number of Artificial Intelligence algorithms in a complete, consistent, and centralized manner, to improve their general accessibility. The project makes use of a standardized algorithm description template that uses well-defined topics that motivate the collection of specific and useful information about each algorithm described. This report describes the Gene Expression Programming algorithm using the standardized template.

Keywords: Clever, Algorithms, Description, Optimization, Gene, Expression, Programming

1 Introduction

The Clever Algorithms project aims to describe a large number of algorithms from the fields of Computational Intelligence, Biologically Inspired Computation, and Metaheuristics in a complete, consistent and centralized manner [1]. The project requires all algorithms to be described using a standardized template that includes a fixed number of sections, each of which is motivated by the presentation of specific information about the technique [2]. This report describes the Gene Expression Programming algorithm using the standardized template.

2 Name

Gene Expression Programming, GEP

3 Taxonomy

Gene Expression Programming is a Global Optimization algorithm and an Automatic Programming technique, and it is an instance of an Evolutionary Algorithm from the field of Evolutionary Computation. It is a sibling of other Evolutionary Algorithms such as the Genetic Algorithm as well as other Evolutionary Automatic Programming techniques such as Genetic Programming and Grammatical Evolution.

*© Copyright 2010 Jason Brownlee. Some Rights Reserved. This work is licensed under a Creative Commons Attribution-Noncommercial-Share Alike 2.5 Australia License.

4 Inspiration

Gene Expression Programming is inspired by the replication and expression of the DNA molecule, specifically at the gene level. The expression of a gene involves the transcription of its DNA to RNA which in turn forms amino acids that make up proteins in the phenotype of an organism. The DNA building blocks are subjected to mechanisms of variation (mutations such as coping errors) as well as recombination during sexual reproduction.

5 Metaphor

Gene Expression Programming uses a linear genome as the basis for genetic operators such as mutation, recombination, inversion, and transposition. The genome is comprised of chromosomes and each chromosome is comprised of genes that are translated into an expression tree to solve a given problem. The robust gene definition means that genetic operators can be applied to the sub-symbolic representation without concern for the structure of the resultant gene expression, providing separation of genotype and phenotype.

6 Strategy

The objective of the Gene Expression Programming algorithm is to improve the adaptive fit of an expressed program in the context of a problem specific cost function. This is achieved through the use of an evolutionary process that operates on a sub-symbolic representation of candidate solutions using surrogates for the processes (descent with modification) and mechanisms (genetic recombination, mutation, inversion, transposition, and gene expression) of evolution.

7 Procedure

A candidate solution is represented as a linear string of symbols called Karva notation or a K-expression, where each symbol maps to a function or terminal node. The linear representation is mapped to an expression tree in a breadth-first manner. A K-expression is fixed length and is comprised of one or more sub-expressions (genes), which are also defined with a fixed length. A gene is comprised of two sections, a head which may contain any function or terminal symbols, and a tail section that may only contain terminal symbols. Each gene will always translate to a syntactically correct expression tree, where the tail portion of the gene provides a genetic buffer which ensures closure of the genes expression.

Algorithm 1 provides a pseudo-code listing of the Gene Expression Programming algorithm for minimizing a cost function.

8 Heuristics

- The length of a chromosome is defined by the number of genes, where a gene length is defined by $h + t$. The h is a user defined parameter (such as 10), where as t is defined as $t = h(n - 1) + 1$. The n represents the maximum arity of functional nodes in the expression (such as 2 if the arithmetic functions $\times, \div, -, +$ are used).
- The mutation operate substituted expressions along the genome, although must respect the gene rules such that function and terminal nodes are mutated in the head of genes, whereas only terminal nodes are substituted in the tail of genes.

Algorithm 1: Pseudo Code for the Gene Expression Programming algorithm.

Input: Grammar, $Population_{size}$, $Head_{length}$, $Tail_{length}$, $P_{crossover}$, $P_{mutation}$
Output: S_{best}

```
1 Population  $\leftarrow$  InitializePopulation( $Population_{size}$ , Grammar,  $Head_{length}$ ,  $Tail_{length}$ );
2 foreach  $S_i \in Population$  do
3    $S_{i_{program}} \leftarrow$  DecodeBreadthFirst( $S_{i_{genome}}$ , Grammar);
4    $S_{i_{cost}} \leftarrow$  Execute( $S_{i_{program}}$ );
5 end
6  $S_{best} \leftarrow$  GetBestSolution(Population);
7 while  $\neg$ StopCondition() do
8   Parents  $\leftarrow$  SelectParents(Population,  $Population_{size}$ );
9   Children  $\leftarrow$  0;
10  foreach  $Parent_1, Parent_2 \in Parents$  do
11     $S_{i_{genome}} \leftarrow$  Crossover( $Parent_1, Parent_2, P_{crossover}$ );
12     $S_{i_{genome}} \leftarrow$  Mutate( $S_{i_{genome}}, P_{mutation}$ );
13    Children  $\leftarrow S_i$ ;
14  end
15  foreach  $S_i \in Population$  do
16     $S_{i_{program}} \leftarrow$  DecodeBreadthFirst( $S_{i_{genome}}$ , Grammar);
17     $S_{i_{cost}} \leftarrow$  Execute( $S_{i_{program}}$ );
18  end
19   $S_{best} \leftarrow$  GetBestSolution(Children);
20  Population  $\leftarrow$  Replace(Population, Children);
21 end
22 return  $S_{best}$ ;
```

- Crossover occurs between two selected parents from the population and can occur based on a one-point cross, two point cross, or finally a gene-based approach where genes are selected from the parents with uniform probability.
- An inversion operator may be used with a low probability that reverses a small sequence of symbols (1-3) within a section of a gene (tail or head).
- A transposition operator may be used that has a number of different modes, including: duplicate a small sequences (1-3) from somewhere on a gene to the head, small sequences on a gene to the root of the gene, and moving of entire genes on the chromosome. In the case of intra-gene transpositions, the sequence in the head of the gene is moved down to accommodate the copied sequence and the length of the head is truncated to maintain consistent gene sizes.
- A '?' is included in the terminal set that represents a numeric constant from an array that are evolved on the end of the genome. The constants are read from the end of the genome and are substituted for '?' as the expression tree is created (in breadth first order). Finally the numeric constants are used as array indices in yet another chromosome of numerical values which are substituted into the expression tree.
- Mutation is low (such as $\frac{1}{L}$), selection can be any of the classical approaches (such as roulette wheel or tournament), and crossover rates are typically high (0.7 of offspring)
- Use multiple sub-expressions linked together on hard problems when one gene does not get much progress. The sub-expressions are linked using link expressions which are function

nodes that are either statically defined (such as a conjunction) or evolved on the genome with the genes.

9 Code Listing

Listing 1 provides an example of the Gene Expression Programming algorithm implemented in the Ruby Programming Language based on the seminal version proposed by Ferreira [3]. The demonstration problem is an instance of symbolic regression $f(x) = x^4 + x^3 + x^2 + x$, where $x \in [-1, 1]$. The grammar used in this problem is: Functions: $F = \{+, -, \div, \times, \}$ and Terminals: $T = \{x\}$. The algorithm uses binary tournament selection, uniform crossover and point mutations. The K-expression is decoded to an expression tree in a breadth-first manner, which is then parsed depth first as an ruby expression string for display and direct evaluation.

```

1 def binary_tournament(population)
2   s1, s2 = population[rand(population.size)], population[rand(population.size)]
3   return (s1[:fitness] > s2[:fitness]) ? s1 : s2
4 end
5
6 def point_mutation(grammar, genome, p_mutation, head_length)
7   child, i = "", 0
8   genome.each_char do |v|
9     if rand() < p_mutation
10      if (i < head_length)
11        child << grammar["FUNC"][rand(grammar["FUNC"].length)]
12      else
13        child << grammar["TERM"][rand(grammar["TERM"].length)]
14      end
15    else
16      child << v
17    end
18    i += 1
19  end
20  return child
21 end
22
23 def uniform_crossover(parent1, parent2, p_crossover)
24   return "" + parent1 if rand() >= p_crossover
25   child = ""
26   parent1.length.times do |i|
27     child << ((rand() < 0.5) ? parent1[i] : parent2[i])
28   end
29   return child
30 end
31
32 def reproduce(grammar, selected, pop_size, p_crossover, p_mutation, head_length)
33   children = []
34   selected.each_with_index do |p1, i|
35     p2 = (i.even?) ? selected[i+1] : selected[i-1]
36     child = {}
37     child[:genome] = uniform_crossover(p1[:genome], p2[:genome], p_crossover)
38     child[:genome] = point_mutation(grammar, child[:genome], p_mutation, head_length)
39     children << child
40   end
41   return children
42 end
43
44 def random_genome(grammar, head_length, tail_length)
45   s = ""
46   head_length.times { s << grammar["FUNC"][rand(grammar["FUNC"].length)] }
47   tail_length.times { s << grammar["TERM"][rand(grammar["TERM"].length)] }
48   return s

```

```

49 end
50
51 def target_function(x)
52   x**4.0 + x**3.0 + x**2.0 + x
53 end
54
55 def cost(program, bounds)
56   errors = 0.0
57   10.times do
58     x = bounds[0] + ((bounds[1] - bounds[0]) * rand())
59     expression = program.gsub("x", x.to_s)
60     target = target_function(x)
61     begin score = eval(expression) rescue score = 0.0/0.0 end
62     errors += (((score.nan? or score.infinite?) ? 0.0 : score) - target).abs
63   end
64   return errors
65 end
66
67 def breadth_first_mapping(genome, grammar)
68   off, queue = 0, Array.new
69   root = {}
70   root[:node] = genome[off].chr; off+=1
71   queue.push(root)
72   while !queue.empty? do
73     current = queue.shift
74     if grammar["FUNC"].include?(current[:node])
75       current[:left] = {}
76       current[:left][:node] = genome[off].chr; off+=1
77       queue.push(current[:left])
78       current[:right] = {}
79       current[:right][:node] = genome[off].chr; off+=1
80       queue.push(current[:right])
81     end
82   end
83   return root
84 end
85
86 def tree_to_string(exp)
87   return exp[:node] if (exp[:left].nil? and exp[:right].nil?)
88   left = tree_to_string(exp[:left])
89   right = tree_to_string(exp[:right])
90   return "(#{left} #{exp[:node]} #{right})"
91 end
92
93 def evaluate(candidate, grammar, bounds)
94   candidate[:expression] = breadth_first_mapping(candidate[:genome], grammar)
95   candidate[:program] = tree_to_string(candidate[:expression])
96   candidate[:fitness] = cost(candidate[:program], bounds)
97 end
98
99 def search(grammar, bounds, head_length, tail_length, generations, pop_size, p_crossover,
100           p_mutation)
101   pop = Array.new(pop_size) do
102     {genome=>random_genome(grammar, head_length, tail_length)}
103   end
104   pop.each{|c| evaluate(c, grammar, bounds)}
105   gen, best = 0, pop.sort{|x,y| y[:fitness] <=> x[:fitness]}.first
106   generations.times do |gen|
107     selected = Array.new(pop){|i| binary_tournament(pop)}
108     children = reproduce(grammar, selected, pop_size, p_crossover, p_mutation, head_length)
109     children.each{|c| evaluate(c, grammar, bounds)}
110     children.sort{|x,y| y[:fitness] <=> x[:fitness]}
111     best = children.first if children.first[:fitness] >= best[:fitness]

```

```

111     pop = children
112     gen += 1
113     puts " > gen=#{gen}, f=#{best[:fitness]}, g=#{best[:genome]}"
114 end
115 return best
116 end
117
118 grammar = {"FUNC"=>["+", "-", "*", "/"], "TERM"=>["x"]}
119 bounds = [-1, 1]
120 head_length = 24
121 tail_length = head_length * (2-1) + 1
122 generations = 150
123 pop_size = 100
124 p_crossover = 0.70
125 p_mutation = 2.0/(head_length+tail_length).to_f
126
127 best = search(grammar, bounds, head_length, tail_length, generations, pop_size, p_crossover,
128               p_mutation)
129 puts "done! Solution: f=#{best[:fitness]}, g=#{best[:genome]}, b=#{best[:program]}"

```

Listing 1: Gene Expression Programming algorithm in the Ruby Programming Language

10 References

10.1 Primary Sources

The Gene Expression Programming algorithm was proposed by Ferreira in a paper that detailed the approach, provided a careful walkthrough of the process and operators and demonstrated the the algorithm on a number of benchmark problem instances such as symbolic regression [3].

10.2 Learn More

Ferreira provided an early and detailed introduction and overview of the approach as book chapter, providing a step-by-step walkthrough of the procedure and sample applications [4]. A similar more contemporary and detailed introduction is provided in a second book chapter [5]. Ferreira published a book on the approach in 2002 covering background, the algorithm, and demonstration applications which is now in its second edition [6].

11 Conclusions

This report described the Gene Expression Programming algorithm as an evolutionary automatic programming technique.

12 Contribute

Found a typo in the content or a bug in the source code? Are you an expert in this technique and know some facts that could improve the algorithm description for all? Do you want to get that warm feeling from contributing to an open source project? Do you want to see your name as an acknowledgment in print?

Two pillars of this effort are i) that the best domain experts are people outside of the project, and ii) that this work is (somewhat) wrong by default. Please help to make this work less wrong by emailing the author ‘Jason Brownlee’ at jasonb@CleverAlgorithms.com or visit the project website at <http://www.CleverAlgorithms.com>.

References

- [1] Jason Brownlee. The clever algorithms project: Overview. Technical Report CA-TR-20100105-1, The Clever Algorithms Project <http://www.CleverAlgorithms.com>, January 2010.
- [2] Jason Brownlee. A template for standardized algorithm descriptions. Technical Report CA-TR-20100107-1, The Clever Algorithms Project <http://www.CleverAlgorithms.com>, January 2010.
- [3] C. Ferreira. Gene expression programming: A new adaptive algorithm for solving problems. *Complex Systems*, 13(2):87–129, 2001.
- [4] C. Ferreira. *Soft Computing and Industry: Recent Applications*, chapter Gene Expression Programming in Problem Solving, pages 635–654. Springer-Verlag, 2002.
- [5] C. Ferreira. *Recent Developments in Biologically Inspired Computing*, chapter Gene Expression Programming and the Evolution of computer programs, pages 82–103. Idea Group Publishing, 2005.
- [6] Candida Ferreira. *Gene expression programming: mathematical modeling by an artificial intelligence*. Springer-Verlag, second edition, 2006.