

Studies in Computational Intelligence, Volume 21

Editor-in-chief

Prof. Janusz Kacprzyk
Systems Research Institute
Polish Academy of Sciences
ul. Newelska 6
01-447 Warsaw
Poland
E-mail: kacprzyk@ibspan.waw.pl

Further volumes of this series
can be found on our homepage:
springer.com

Vol. 5. Da Ruan, Guoqing Chen, Etienne E.
Kerre, Geert Wets (Eds.)
Intelligent Data Mining, 2005
ISBN 3-540-26256-3

Vol. 6. Tsau Young Lin, Setsuo Ohsuga,
Churn-Jung Liao, Xiaohua Hu, Shusaku
Tsumoto (Eds.)
Foundations of Data Mining and Knowledge Dis-
covery, 2005
ISBN 3-540-26257-1

Vol. 7. Bruno Apolloni, Ashish Ghosh, Ferda, Al-
paslan, Lakhmi C. Jain, Srikanta
Patnaik (Eds.)
Machine Learning and Robot Perception, 2005
ISBN 3-540-26549-X

Vol. 8. Srikanta Patnaik, Lakhmi C. Jain,
Spyros G. Tzafestas, Germano Resconi,
Amit Konar (Eds.)
Innovations in Robot Mobility and Control,
2005
ISBN 3-540-26892-8

Vol. 9. Tsau Young Lin, Setsuo Ohsuga,
Churn-Jung Liao, Xiaohua Hu (Eds.)
Foundations and Novel Approaches in Data Min-
ing, 2005
ISBN 3-540-28315-3

Vol. 10. Andrzej P. Wierzbicki, Yoshiteru
Nakamori
Creative Space, 2005
ISBN 3-540-28458-3

Vol. 11. Antoni Ligęza
Logical Foundations for Rule-Based
Systems, 2006
ISBN 3-540-29117-2

Vol. 13. Nadia Nedjah, Ajith Abraham,
Luiza de Macedo Mourelle (Eds.)
Genetic Systems Programming, 2006
ISBN 3-540-29849-5

Vol. 14. Spiros Sirmakessis (Ed.)
Adaptive and Personalized Semantic Web, 2006
ISBN 3-540-30605-6

Vol. 15. Lei Zhi Chen, Sing Kiong Nguang,
Xiao Dong Chen
Modelling and Optimization of
Biotechnological Processes, 2006
ISBN 3-540-30634-X

Vol. 16. Yaochu Jin (Ed.)
Multi-Objective Machine Learning, 2006
ISBN 3-540-30676-5

Vol. 17. Te-Ming Huang, Vojislav Kecman,
Ivica Kopriva
Kernel Based Algorithms for Mining Huge
Data Sets, 2006
ISBN 3-540-31681-7

Vol. 18. Chang Wook Ahn
Advances in Evolutionary Algorithms, 2006
ISBN 3-540-31758-9

Vol. 19. Ajita Ichalkaranje, Nikhil
Ichalkaranje, Lakhmi C. Jain (Eds.)
Intelligent Paradigms for Assistive and
Preventive Healthcare, 2006
ISBN 3-540-31762-7

Vol. 20. Wojciech Pecznec, Agata Pórola
Advances in Verification of Time Petri Nets
and Timed Automata, 2006
ISBN 3-540-32869-6

Vol. 21. Cândida Ferreira
Gene Expression Programming: Mathemati-
cal Modeling by an Artificial Intelligence,
2006
ISBN 3-540-32796-7

Cândida Ferreira

Gene Expression Programming

Mathematical Modeling
by an Artificial Intelligence

Second, revised and extended edition

 Springer

Dr. Cândida Ferreira
Chief Scientist
Gepsoft Ltd.
73 Elmtree Drive
Bristol BS13 8NA
United Kingdom
E-mail: candidaf@gepsoft.com

Library of Congress Control Number: 2006921791

ISSN print edition: 1860-949X

ISSN electronic edition: 1860-9503

ISBN-10 3-540-32796-7 Springer Berlin Heidelberg New York

ISBN-13 978-3-540-32796-7 Springer Berlin Heidelberg New York

This work is subject to copyright. All rights are reserved, whether the whole or part of the material is concerned, specifically the rights of translation, reprinting, reuse of illustrations, recitation, broadcasting, reproduction on microfilm or in any other way, and storage in data banks. Duplication of this publication or parts thereof is permitted only under the provisions of the German Copyright Law of September 9, 1965, in its current version, and permission for use must always be obtained from Springer-Verlag. Violations are liable to prosecution under the German Copyright Law.

Springer is a part of Springer Science+Business Media
springer.com

© Springer-Verlag Berlin Heidelberg 2006

Printed in The Netherlands

The use of general descriptive names, registered names, trademarks, etc. in this publication does not imply, even in the absence of a specific statement, that such names are exempt from the relevant protective laws and regulations and therefore free for general use.

Cover design: deblik, Berlin

Typesetting: by the author and TechBooks

Printed on acid-free paper SPIN: 11506591

89/Strasser

5 4 3 2 1 0

To José Simas
For All the Dreams

and

To my Grandfather, Domingos de Carvalho
For His Vision

Preface to the Second Edition

The idea for this second edition came from Janusz Kacprzyk on April 29, 2005, who kindly invited me to his new Springer series, Studies in Computational Intelligence. The initial plan was to correct the usual typos and mistakes but leave the book unchanged, as Janusz thought (and I agreed with him) that it was the proper moment for a second edition. But then there was the problem of the new format and I had to reformat and proofread everything again. And I just thought that while I was at it, I might as well change some things in the book to make it more enjoyable and interesting. Foremost in my thoughts was the restructuring of chapter 4, The Basic GEA in Problem Solving. In that chapter, buried together with a wide variety of problems, were several important new algorithms that I wanted to bring to the forefront. These algorithms include: the GEP-RNC algorithm (the cornerstone of several new other algorithms); automatically defined functions; polynomial induction; and parameter optimization. So I removed all these materials from chapter 4 and gave them the deserved attention by writing four new chapters (chapter 5 Numerical Constants and the GEP-RNC Algorithm, chapter 6 Automatically Defined Functions in Problem Solving, chapter 7 Polynomial Induction and Time Series Prediction, and chapter 8 Parameter Optimization). Then this new structure just begged for me to include one of the last additions to the GEP technique – decision trees – that I was regrettably unable to include in the first edition (I implemented decision trees in August of 2002, two months before sending the manuscript to the printer). So, chapter 9, Decision Tree Induction, is totally new to this second edition and an interesting addition both to the book and to GEP. The last three chapters, chapter 10 Design of Neural Networks, chapter 11 Combinatorial Optimization, and chapter 12 Evolutionary Studies, remain basically unchanged.

With all this restructuring in chapter 4, I was able to develop several new topics, including solving problems with multiple outputs in one go and

designing parsimonious solutions both with parsimony pressure and user defined functions, also interesting new extensions to the GEP technique. Furthermore, the section on Logic Synthesis was totally restructured and an interesting analysis of the most common universal logical systems is presented.

Chapter 3, The Basic Gene Expression Algorithm, with the exception of section 2, Fitness Functions and the Selection Environment, remains practically unchanged. In section 2, however, I introduce several new fitness functions that are then used to explore more efficiently the solution landscapes of the problems solved in the book.

The inversion operator is one of the latest additions to the GEP technique, and you will notice an extra entry for it in chapters 3 and 10. Furthermore, you'll also notice that both IS and RIS transposition were slightly modified so that the transposon sizes were automatically chosen rather than a priori set. But unbeknownst to me, this slightly different implementation had consequences in the performance of these operators, and the new implementation is slightly more unpredictable in terms of performance. This is the reason why in the Evolutionary Studies of chapter 12, the old implementation of these operators is still used. And to be fair, the comparison of the inversion operator with these transposition operators should also use a fixed set of sizes for the inverted sequences. And this is the reason why inversion is not analyzed in chapter 12.

Cândida Ferreira
December 20, 2005

Preface to the First Edition

I developed the basic ideas of gene expression programming (GEP) in September and October of 1999 almost unaware of their uniqueness. I was reading Mitchell's book *An Introduction to Genetic Algorithms* (Mitchell 1996) and meticulously solving all the computer exercises provided at the end of each chapter. Therefore, I implemented my first genetic algorithm and I also implemented what I thought was a genetic programming (GP) system. Like a GP system, this new system could also evolve computer programs of different sizes and shapes but, surprisingly, it surpassed the old GP system by a factor of 100-60,000. So, what happened here? What was responsible for this astounding difference in performance? For an evolutionary biologist, the answer is quite straightforward: this new system – gene expression programming – simply crossed the phenotype threshold. This means that the complex computer programs (the phenotype) evolved by GEP are totally encoded in simple strings of fixed length (the chromosomes or genotype). The separation of the genotype from the phenotype is comparable to opening a Pandora box full of good things or possibilities. Of these good things, perhaps the most important is that there are virtually no restrictions concerning the number or type of genetic operators used. Another important thing is that the creation of higher levels of complexity becomes practically a trivial task. Indeed, it was trivial to create a multigenic system from a unigenic one and a multicellular system from a unicellular one. And each new system creates its own box of new possibilities, which enlarges considerably the scope of this new technique.

In this first book on gene expression programming I describe thoroughly the basic gene expression algorithm and numerous modifications to this new algorithm, providing all the implementation details so that anyone with elementary programming skills (or willing to learn them) will be able to implement it themselves. The first chapter briefly introduces the main players

of biological gene expression in order to show how they relate to the main players of artificial evolutionary systems in general and GEP in particular. The second chapter introduces the players of gene expression programming, showing their structural and functional organization in detail. The language especially created to express the genetic information of GEP chromosomes is also described in this chapter. Chapter 3 gives a detailed description of the basic gene expression algorithm and the basic genetic operators. In addition, a very simple problem is exhaustively dissected, showing all the individual programs created during the discovery process in order to demystify the workings of adaptation and evolution. Chapter 4 describes some of the applications of the basic gene expression algorithm, including a large body of unpublished materials, namely, parameter optimization, evolution of Kolmogorov-Gabor polynomials, time series prediction, classifier systems, evolution of linking functions, multicellularity, automatically defined functions, user defined functions and so forth. The materials of Chapter 5 are also new and show how to simulate complete neural networks with gene expression programming. Two benchmark problems are solved with these GEP-nets, providing an effective measure of their performance. Chapter 6 shows how to do combinatorial optimization with gene expression programming. Multigene families and several combinatorial-specific operators are introduced and their performance evaluated on two scheduling problems. The last chapter discusses some important and controversial evolutionary topics that might be refreshing to both evolutionary computists and evolutionary biologists.

Acknowledgments

The invention of a new paradigm can often create strong resistance, especially if it seems to endanger long established technologies and enterprises. The publication of my work on scientific journals and conferences, which should be forums for discussing and sharing new ideas, became a nightmare and both my work and myself were outright dismissed and treated with scorn. Despite the initial opposition and due to a set of happy circumstances and resources, I was finally able to make my work known and available to all. I am deeply indebted to José Simas, an accomplished graphic and web designer and software developer, for believing in me and in GEP from the beginning and for helping its expansion and promotion on the World Wide

Web. Together we founded Gepsoft and developed software based on gene expression programming which is already helping numerous scientists and engineers worldwide. And thanks to Gepsoft it was possible for me to concentrate fully on the writing of this book and on the development of several new algorithms. Indeed, my work at Gepsoft benefited tremendously from my writing and vice versa.

I am also very grateful to Pedro Carneiro, a talented musician with an avid mind, for reading and editing the first three chapters of the manuscript. José Simas also read several drafts of the manuscript, accompanying the process from the beginning and contributing with valuable discussions and suggestions. He is, in fact, my first reader and I always write with him in my mind.

Finally, I would like to thank José Gabriel, a talented printer and skilled craftsman, for his involvement in the making of this book from the start and for handling its printing with special care.

Cândida Ferreira
October 15, 2002

List of Symbols

ADF	Automatically defined function
ADF-RNC	ADFs with random numerical constants
AND	AND or Boolean function of two arguments #8
APS	Gepsoft Automatic Problem Solver
CA	Cellular automata
Dc	Gene domain for encoding random numerical constants
DT	Decision tree
Dt	Gene domain for encoding the thresholds of neural networks
Dw	Gene domain for encoding the weights of neural networks
EDT	Evolvable decision trees
EDT-RNC	Evolvable decision trees with numeric attributes
ET	Expression tree
FN	False negatives
FP	False positives
GA	Genetic algorithm
GEA	Gene expression algorithm
GEP	Gene expression programming
GEP-ADF	GEP with automatically defined functions
GEP-EDT	GEP for inducing evolvable decision trees
GEP-KGP	GEP for inducing Kolmogorov-Gabor polynomials
GEP-MO	GEP for solving problems with multiple outputs
GEP-NC	GEP with numerical constants
GEP-nets	GEP for inducing neural networks
GEP-NN	GEP for inducing neural networks
GEP-PO	GEP for parameter optimization
GEP-RNC	GEP with random numerical constants
GKL	Gacs-Kurdyumov-Levin rule
GOE	Greater Or Equal or Boolean function of two arguments #13

GP	Genetic programming
GT	Greater Than or Boolean function of two arguments #4
HZero	Parameter optimization algorithm with a head size of zero
IC	Initial configuration
IF	Rule of three arguments #202, If $a = 1$, then b , else c
IS	Insertion sequence elements
LOE	Less Or Equal or Boolean function of two arguments #11
LT	Less Than or Boolean function of two arguments #2
MAJ	Majority function of three arguments or rule #232
MGF	Multigene family
MIN	Minority function of three arguments or rule #23
MUX	3-Multiplexer or rule #172, If $a = 0$, then b , else c
NAND	NAND or Boolean function of two arguments #7
NC	Numerical constants
NLM	NAND-like module
NOR	NOR or Boolean function of two arguments #1
NPV	Negative predictive value
NXOR	NXOR or Boolean function of two arguments #9
OR	OR or Boolean function of two arguments #14
ORF	Open reading frame
PPV	Positive predictive value
RIS	Root insertion sequence elements
RNC	Random numerical constants
TAP	Task assignment problem
TN	True negatives
TP	True positives
TSP	Traveling salesperson problem
UDF	User defined function
ULM	Universal logical module
XOR	Exclusive-OR or Boolean function of two arguments #6

Contents

<i>Preface to the Second Edition</i>	vii
<i>Preface to the First Edition</i>	ix
<i>List of Symbols</i>	xiii
1 Introduction: The Biological Perspective	1
1.1 The Entities of Biological Gene Expression	3
1.1.1 DNA	3
1.1.2 RNA	4
1.1.3 Proteins	6
1.2 Biological Gene Expression	8
1.2.1 Genome Replication	8
1.2.2 Genome Restructuring	8
1.2.2.1 Mutation	10
1.2.2.2 Recombination	12
1.2.2.3 Transposition	13
1.2.2.4 Gene Duplications	14
1.2.3 Transcription	14
1.2.4 Translation and Posttranslational Modifications	16
1.2.4.1 Translation	16
1.2.4.2 Posttranslational Modifications	18
1.3 Adaptation and Evolution	19
1.4 Genetic Algorithms	21
1.5 Genetic Programming	22
1.6 Gene Expression Programming	26
2 The Entities of Gene Expression Programming	29
2.1 The Genome	30
2.1.1 Open Reading Frames and Genes	30
2.1.2 Structural and Functional Organization of Genes	34
2.1.3 Multigenic Chromosomes	37

2.2 Expression Trees and the Phenotype	38
2.2.1 Information Decoding: Translation	39
2.2.2 Posttranslational Interactions and Linking Functions	43
2.3 Cells and the Evolution of Linking Functions	47
2.3.1 Homeotic Genes and the Cellular System	47
2.3.2 Multicellular Systems with Multiple Main Programs	49
2.4 Other Levels of Complexity	51
2.5 Karva Language: The Language of GEP	52
3 The Basic Gene Expression Algorithm	55
3.1 Populations of Individuals	57
3.1.1 Creation of the Initial Population	58
3.1.2 Subsequent Generations and Elitism	62
3.2 Fitness Functions and the Selection Environment	65
3.2.1 The Selection Environment	65
3.2.2 Fitness Functions for Symbolic Regression	66
3.2.2.1 Number of Hits	66
3.2.2.2 Precision and Selection Range	67
3.2.2.3. Mean Squared Error	68
3.2.2.4 R-square	69
3.2.3 Fitness Functions for Classification and Logic Synthesis	69
3.2.3.1 Number of Hits	71
3.2.3.2 Hits with Penalty	71
3.2.3.3 Sensitivity / Specificity	72
3.2.3.4 Positive Predictive Value / Negative Predictive Value	73
3.2.4 Selection Mechanism	73
3.3 Reproduction with Modification	74
3.3.1 Replication and Selection	75
3.3.2 Mutation	77
3.3.3 Inversion	81
3.3.4 Transposition and Insertion Sequence Elements	85
3.3.4.1 Transposition of IS Elements	86
3.3.4.2 Root Transposition	88
3.3.4.3 Gene Transposition	91
3.3.5 Recombination	95
3.3.5.1 One-point Recombination	95
3.3.5.2 Two-point Recombination	99
3.3.5.3 Gene Recombination	102
3.4 Solving a Simple Problem with GEP	106

4 The Basic GEA in Problem Solving	121
4.1 Symbolic Regression	122
4.1.1 Function Finding on a One-dimensional Parameter Space	122
4.1.2 Function Finding on a Five-dimensional Parameter Space	133
4.1.3 Mining Meaningful Information from Noisy Data	135
4.2 Classification Problems	136
4.2.1 Diagnosis of Breast Cancer	137
4.2.2 Credit Screening	141
4.2.3 Fisher's Irises	144
4.2.3.1 Decomposing a Three-class Problem	144
4.2.3.2 Multiple Genes for Multiple Outputs	146
4.3 Logic Synthesis and Parsimonious Solutions	150
4.3.1 Fitness Functions with Parsimony Pressure	151
4.3.2 Universal Logical Systems	152
4.3.2.1 Boolean Logic	156
4.3.2.2 Nand Logic	158
4.3.2.3 Nor Logic	159
4.3.2.4 Reed-Muller Logic	161
4.3.2.5 Mux Logic	162
4.3.3 Using User Defined Functions as Building Blocks	164
4.3.3.1 Odd-parity Functions	166
4.3.3.2 Exactly-one-on Functions	169
4.4 Evolving Cellular Automata Rules for the Density-classification Problem.....	174
4.4.1 The Density-classification Task	175
4.4.2 Two Rules Discovered by GEP	176
5 Numerical Constants and the GEP-RNC Algorithm	181
5.1 Handling Constants in Automatic Programming	181
5.2 Genes with Multiple Domains to Encode RNCs	188
5.3 Multigenic Systems with RNCs	191
5.4 Special Search Operators for Fine-tuning the RNCs	193
5.4.1 Dc-specific Mutation	194
5.4.2 Dc-specific Inversion	196
5.4.3 Dc-specific Transposition	196
5.4.4 Direct Mutation of RNCs	198
5.5 Solving a Simple Problem with GEP-RNC	201
5.6 Three Approaches to the Creation of NCs	210
5.6.1 Problems and Settings	211

5.6.2 Sequence Induction	214
5.6.3 “V” Function	218
5.6.4 Diagnosis of Breast Cancer	225
5.6.5 Analog Circuit Design	228
6 Automatically Defined Functions in Problem Solving	233
6.1 Solving a Simple Modular Function with ADFs	234
6.2 Odd-parity Functions	248
6.3 Kepler’s Third Law	253
6.4 RNCs and the Cellular System	257
6.4.1 Incorporating RNCs in ADFs	257
6.4.2 Designing Analog Circuits with the ADF-RNC Algorithm	258
6.5 Diagnosis of Breast Cancer	264
6.6 Multiple Cells for Multiple Outputs: The Iris Problem	269
7 Polynomial Induction and Time Series Prediction	275
7.1 Evolution of Kolmogorov-Gabor Polynomials	275
7.2 Simulating STROGANOFF in GEP	278
7.3 Evaluating the Performance of STROGANOFF	279
7.3.1 Original and Enhanced STROGANOFF	280
7.3.2 Simpler GEP Systems	286
7.4 Predicting Sunspots with GEP	291
8 Parameter Optimization	297
8.1 The HZero Algorithm	298
8.1.1 The Architecture	298
8.1.2 Optimization of a Simple Function	300
8.2 The GEP-PO Algorithm	312
8.2.1 The Architecture	312
8.2.2 Optimization of a Simple Function	314
8.3 Maximum Seeking with GEP	328
9 Decision Tree Induction	337
9.1 Decision Trees with Nominal Attributes	338
9.1.1 The Architecture	339
9.1.2 A Simple Problem with Nominal Attributes	341
9.2 Decision Trees with Numeric/Mixed Attributes	349
9.2.1 The Architecture	351
9.2.2 A Simple Problem with Mixed Attributes	353
9.3 Solving Problems with GEP Decision Trees	361

9.3.1 Diagnosis of Breast Cancer	362
9.3.2 Classification of Irises	364
9.3.3 The Lymphography Problem	373
9.3.4 The Postoperative Patient Problem	375
9.4 Pruning Trees with Parsimony Pressure	377
10 Design of Neural Networks	381
10.1 Genes with Multiple Domains for NN Simulation	382
10.2 Special Genetic Operators	384
10.2.1 Domain-specific Inversion	386
10.2.2 Domain-specific Transposition	386
10.2.3 Intragenic Two-point Recombination	388
10.2.4 Direct Mutation of Weights and Thresholds	392
10.3 Solving Problems with GEP Neural Networks	394
10.3.1 Neural Network for the Exclusive-or Problem	394
10.3.2 Neural Network for the 6-Multiplexer	397
10.4 Evolutionary Dynamics of GEP-nets	402
11 Combinatorial Optimization	405
11.1 Multigene Families and Scheduling Problems	406
11.2 Combinatorial-specific Operators: Performance and Mechanisms ..	407
11.2.1 Inversion	408
11.2.2 Gene Deletion/Insertion	409
11.2.3 Restricted Permutation	410
11.2.4 Other Search Operators	411
11.2.4.1 Sequence Deletion/Insertion	411
11.2.4.2 Generalized Permutation	411
11.3 Two Scheduling Problems	413
11.3.1 The Traveling Salesperson Problem	413
11.3.2 The Task Assignment Problem	416
11.4 Evolutionary Dynamics of Simple GEP Systems	418
12 Evolutionary Studies	421
12.1 Genetic Operators and their Power	421
12.1.1 Their Performances	422
12.1.2 Evolutionary Dynamics of Different Types of Populations ..	425
12.1.2.1 Mutation	425
12.1.2.2 Transposition	427
12.1.2.3 Recombination	429
12.2 The Founder Effect	431

12.2.1 Choosing the Population Types	432
12.2.2 The Founder Effect in Simulated Evolutionary Processes ...	436
12.3 Testing the Building Block Hypothesis	438
12.4 The Role of Neutrality in Evolution	440
12.4.1 Genetic Neutrality in Unigenic Systems	442
12.4.2 Genetic Neutrality in Multigenic Systems	445
12.5 The Higher Organization of Multigenic Systems	448
12.6 The Open-ended Evolution of GEP Populations	450
12.7 Analysis of Different Selection Schemes	453
Bibliography	457
Index	463