

# **Statistics and Computing**

*Series Editors:*

J. Chambers

D. Hand

W. Härdle

For other titles published in this series, go to  
<http://www.springer.com/series/3022>



Roberto Baragona · Francesco Battaglia · Irene Poli

# Evolutionary Statistical Procedures

An Evolutionary Computation Approach  
to Statistical Procedures Designs  
and Applications

 Springer

Prof. Roberto Baragona  
Sapienza University of Rome  
Department of Communication  
and Social Research  
Via Salaria 113  
00198 Rome  
Italy  
roberto.baragona@uniroma1.it

Prof. Francesco Battaglia  
Sapienza University of Rome  
Department of Statistical Sciences  
Piazzale Aldo Moro 5  
00100 Roma  
Italy  
francesco.battaglia@uniroma1.it

Prof. Irene Poli  
Ca' Foscari University of Venice  
Department of Statistics  
Cannaregio 873  
30121 Venice  
Italy  
irenpoli@unive.it

**Series Editors:**

J. Chambers  
Department of Statistics  
Sequoia Hall  
390 Serra Mall  
Stanford University  
Stanford, CA 94305-4065

D. Hand  
Department of Mathematics  
Imperial College London,  
South Kensington Campus  
London SW7 2AZ  
United Kingdom

W. Härdle  
C.A.S.E. Centre for Applied  
Statistics and Economics  
School of Business and  
Economics  
Humboldt-Universität zu Berlin  
Unter den Linden 6  
10099 Berlin  
Germany

ISSN 1431-8784  
ISBN 978-3-642-16217-6      e-ISBN 978-3-642-16218-3  
DOI 10.1007/978-3-642-16218-3  
Springer Heidelberg Dordrecht London New York

© Springer-Verlag Berlin Heidelberg 2011

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. Violations are liable to prosecution under the German Copyright Law.

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:* SPI Publisher Services

Printed on acid-free paper

Springer is part of Springer Science+Business Media (www.springer.com)

# Preface

In many application fields, artificial intelligence, data mining, pattern recognition operations research, to name but a few, often problems arise that may be reduced at their very essence to optimization problems.

Unfortunately, neither the objective function nor the solution search space display that nice properties that may be conveniently exploited by widespread familiar numerical analysis tools. Though these latter offer powerful devices to cope with a great deal of both theoretical and practical problems in so many disciplines, the hypotheses on which they rely are far from being fulfilled within the frameworks that so often constitute the background of the application fields we mentioned so far. Here well behaved analytic functions and compact domains are not commonplace and only raw hazardous simplifying assumptions may constrain so many classes of problems in such a way they may be treated by means of comfortable numerical procedures.

It may happen that too much simplification does not allow the actual problem to be solved satisfactorily, that is we obtain useless though well grounded solutions. Heuristic methods have been developed to yield reliable solutions to many particular problems, but only the development of general heuristics offered a theoretical framework for dealing with a large class of problems.

One outstanding approach in this kind of methods proved to be evolutionary computing. This rapidly growing research field came nowadays to a well established discipline on its own enjoying solid theoretical foundations and large evidence of effectiveness as far as complex non conventional optimization problems are concerned. History dates back to the fifties and since then an enormous body of theory has been developed which makes evolutionary computing a suitable framework for building applied methodology while it is an active and thriving research field. Its influence spread out through so many disciplines, from biology to informatics and engineering and to economics and social sciences.

We shall try to make an account of the influence of evolutionary computing in Statistics and close related fields. Evolutionary computation is particularly useful in Statistics, in all cases when the statistician has to select, inside a very large discrete set just one element, be it a method, a model, a parameter value, or such.

Therefore a common application of evolutionary computation is to the selection of variables, both in regression problems and in time series linear models. In time series analysis it has been proposed also for building non linear models. For the same reason, evolutionary computation may be employed in the problem of outlier detection, and several papers were published both for the independent observations case and for time series.

A recent, very promising application is in the design of experiments, where the optimal choice of a combination of factors, and their levels, is needed, and cost constraints are very strong.

Finally, a typical field of application of evolutionary computation to Statistics is cluster analysis. Here, the use of an evolutionary algorithm provides a valid alternative when the number of units and variables is large, and the standard cluster techniques allow only approximate solutions.

This book brings together most literature on the use of evolutionary computation methods in statistical analysis, and contains also some unpublished material. It is based on the over 15 years experience and research work of the authors in this field.

This book requires a basic knowledge of mathematics and statistics, and may be useful to research students and professionals to appreciate the suitability for solving complex statistical problems of evolutionary computation. We believe that these algorithms will become a common standard in Statistics in a few years.

Much of the research work included in this book has been possible due to generous funding from Institutions that we are happy to acknowledge with thanks. F. Battaglia and R. Baragona gratefully acknowledge funding from MIUR, Italy (PRIN 2007) and European Commission through Marie Curie Research and Training Network COMISEF – Computational Methods in Statistics, Econometrics and Finance (contract MRTN-CT-2006-034270). I. Poli also would like to acknowledge the MIUR for the PRIN 2007 project, the European Commission for the PACE integrated project (IST-FET, [www.istpace.org](http://www.istpace.org)), and the Fondazione di Venezia for the DICE project. She also would like to thank the researchers at the European Centre for Living Technology for the very fruitful collaboration at the development of the evolutionary perspective in Statistics, in particular Davide de March, Debora Slanzi, Laura Villanova, Matteo Borrotti, Michele Forlin.

We are willing to express our gratitude to many colleagues who commented on large parts of the manuscript letting some obscurities and subtleties be disclosed. In particular we need to mention Antonietta di Salvatore and Domenico Cucina who patiently and carefully read all of the versions of the manuscript, and Debora Slanzi who provided us with lot of computer programs and graphical displays, specially as regards design of experiments. We are indebted for useful and animated discussions to Matt Protopapas, specially concerned with issues related to non linear time series, and to Sanghamitra Bandyopadhyay and Ujjwal Maulik who pointed out in particular some important aspects related to multiobjective optimization and cluster analysis. We have to thank the Springer Editors who successively took care of the manuscript for their assiduous assistance and encouragement, Lilith Braun and Peter

Niels Thomas. Also our thanks are deservedly due to the Desk Editor Alice Blanck for her accurate preparation of the final version of the manuscript, and to Samuel Roobesh, Project Manager at Integra Software Services Pvt. Ltd., for his attentive care in handling the production of this book.

Rome, Italy  
Rome, Italy  
Venice, Italy  
September 2010

Roberto Baragona  
Francesco Battaglia  
Irene Poli





# Contents

<b>1</b>	<b>Introduction</b>	1
1.1	Bio-inspired Optimization Methods	1
1.2	Topics Organization	4
<b>2</b>	<b>Evolutionary Computation</b>	5
2.1	Introduction	5
2.1.1	Evolutionary Computation Between Artificial Intelligence and Natural Evolution	5
2.1.2	The Contribution of Genetics	8
2.2	Evolutionary Computation Methods	10
2.2.1	Essential Properties	10
2.2.2	Evolutionary Programming	14
2.2.3	Evolution Strategies	16
2.2.4	Genetic Algorithms	18
2.2.5	Estimation of Distribution Algorithms	20
2.2.6	Differential Evolution	23
2.2.7	Evolutionary Behavior Algorithms	25
2.2.8	A Simple Example of Evolutionary Computation	27
2.3	Properties of Genetic Algorithms	36
2.3.1	Genetic Algorithms as a Paradigm of Evolutionary Computation	36
2.3.2	Evolution of Genetic Algorithms	41
2.3.3	Convergence of Genetic Algorithms	44
2.3.4	Issues in the Implementation of Genetic Algorithms	51
2.3.5	Genetic Algorithms and Random Sampling from a Probability Distribution	56
<b>3</b>	<b>Evolving Regression Models</b>	63
3.1	Introduction	63
3.2	Identification	64
3.2.1	Linear Regression	64
3.2.2	Generalized Linear Models	67
3.2.3	Principal Component Analysis	68

3.3	Parameter Estimation . . . . .	69
3.3.1	Regression Models . . . . .	69
3.3.2	The Logistic Regression Model . . . . .	70
3.4	Independent Component Analysis . . . . .	74
3.4.1	ICA algorithms . . . . .	76
3.4.2	Simple GAs for ICA . . . . .	77
3.4.3	GAs for Nonlinear ICA . . . . .	83
<b>4</b>	<b>Time Series Linear and Nonlinear Models . . . . .</b>	<b>85</b>
4.1	Models of Time Series . . . . .	86
4.2	Autoregressive Moving Average Models . . . . .	88
4.2.1	Identification of ARMA Models by Genetic Algorithms . . . . .	91
4.2.2	More General Models . . . . .	95
4.3	Nonlinear Models . . . . .	97
4.3.1	Threshold AR and Double Threshold GARCH Models . . . . .	97
4.3.2	Exponential Models . . . . .	100
4.3.3	Piecewise Linear Models . . . . .	103
4.3.4	Bilinear Models . . . . .	114
4.3.5	Real Data Applications . . . . .	116
4.3.6	Artificial Neural Networks . . . . .	118
<b>5</b>	<b>Design of Experiments . . . . .</b>	<b>125</b>
5.1	Introduction . . . . .	125
5.2	Experiments and Design of Experiments . . . . .	126
5.2.1	Randomization, Replication and Blocking . . . . .	128
5.2.2	Factorial Designs and Response Surface Methodology . . . . .	129
5.3	The Evolutionary Design of Experiments . . . . .	132
5.3.1	High-Dimensionality Search Space . . . . .	132
5.3.2	The Evolutionary Approach to Design Experiments . . . . .	133
5.3.3	The Genetic Algorithm Design (GA-Design) . . . . .	135
5.4	The Evolutionary Model-Based Experimental Design: The Statistical Models in the Evolution . . . . .	144
5.4.1	The Evolutionary Neural Network Design (ENN-Design) . . . . .	144
5.4.2	The Model Based Genetic Algorithm Design (MGA-Design) . . . . .	147
5.4.3	The Evolutionary Bayesian Network Design (EBN-Design) . . . . .	152
<b>6</b>	<b>Outliers . . . . .</b>	<b>159</b>
6.1	Outliers in Independent Data . . . . .	159
6.1.1	Exploratory Data Analysis for Multiple Outliers Detection . . . . .	160
6.1.2	Genetic Algorithms for Detecting Outliers in an i.i.d. Data Set . . . . .	162
6.2	Outliers in Time Series . . . . .	167
6.2.1	Univariate ARIMA Models . . . . .	169
6.2.2	Multivariate Time Series Outlier Models . . . . .	181

- 6.3 Genetic Algorithms for Multiple Outlier Detection . . . . . 184
  - 6.3.1 Detecting Multiple Outliers in Univariate Time Series . . . . . 186
  - 6.3.2 Genetic Algorithms for Detecting Multiple Outliers  
in Multivariate Time Series . . . . . 187
  - 6.3.3 An Example of Application to Real Data . . . . . 191
- 7 Cluster Analysis . . . . . 199**
  - 7.1 The Partitioning Problem . . . . . 199
    - 7.1.1 Classification . . . . . 200
    - 7.1.2 Algorithms for Clustering Data . . . . . 204
    - 7.1.3 Indexes of Cluster Validity . . . . . 212
  - 7.2 Genetic Clustering Algorithms . . . . . 219
    - 7.2.1 A Genetic Divisive Algorithm . . . . . 219
    - 7.2.2 Quick Partition Genetic Algorithms . . . . . 221
    - 7.2.3 Centroid Evolution Algorithms . . . . . 227
    - 7.2.4 The Grouping Genetic Algorithm . . . . . 230
    - 7.2.5 Genetic Clustering of Large Data Sets . . . . . 233
  - 7.3 Fuzzy Partition . . . . . 234
    - 7.3.1 The Fuzzy c-Means Algorithm . . . . . 234
    - 7.3.2 Genetic Fuzzy Partition Algorithms . . . . . 236
  - 7.4 Multivariate Mixture Models Estimation by Evolutionary  
Computing . . . . . 239
    - 7.4.1 Genetic Multivariate Mixture Model Estimates . . . . . 240
    - 7.4.2 Hybrid Genetic Algorithms and the EM Algorithm . . . . . 244
    - 7.4.3 Multivariate Mixture Model Estimates with Unknown  
Number of Mixtures . . . . . 246
  - 7.5 Genetic Algorithms in Classification and Regression Trees Models . . 248
  - 7.6 Clusters of Time Series and Directional Data . . . . . 248
    - 7.6.1 GAs-Based Methods for Clustering Time Series Data . . . . . 249
    - 7.6.2 GAs-Based Methods for Clustering Directional Data . . . . . 254
  - 7.7 Multiobjective Genetic Clustering . . . . . 258
    - 7.7.1 Pareto Optimality . . . . . 258
    - 7.7.2 Multiobjective Genetic Clustering Outline . . . . . 259
- References . . . . . 261**
- Index . . . . . 273**