

Intelligent Systems Reference Library

Volume 118

Series editors

Janusz Kacprzyk, Polish Academy of Sciences, Warsaw, Poland
e-mail: kacprzyk@ibspan.waw.pl

Lakhmi C. Jain, University of Canberra, Canberra, Australia;
Bournemouth University, UK;
KES International, UK
e-mails: jainlc2002@yahoo.co.uk; Lakhmi.Jain@canberra.edu.au
URL: <http://www.kesinternational.org/organisation.php>

About this Series

The aim of this series is to publish a Reference Library, including novel advances and developments in all aspects of Intelligent Systems in an easily accessible and well structured form. The series includes reference works, handbooks, compendia, textbooks, well-structured monographs, dictionaries, and encyclopedias. It contains well integrated knowledge and current information in the field of Intelligent Systems. The series covers the theory, applications, and design methods of Intelligent Systems. Virtually all disciplines such as engineering, computer science, avionics, business, e-commerce, environment, healthcare, physics and life science are included.

More information about this series at <http://www.springer.com/series/8578>

Dionisios N. Sotiropoulos · George A. Tsihrintzis

Machine Learning Paradigms

Artificial Immune Systems and their
Applications in Software Personalization



Springer

Dionisios N. Sotiropoulos
University of Piraeus
Piraeus
Greece

George A. Tsirhrintzis
University of Piraeus
Piraeus
Greece

ISSN 1868-4394

ISSN 1868-4408 (electronic)

Intelligent Systems Reference Library

ISBN 978-3-319-47192-1

ISBN 978-3-319-47194-5 (eBook)

DOI 10.1007/978-3-319-47194-5

Library of Congress Control Number: 2016953915

© Springer International Publishing AG 2017

This work is subject to copyright. All rights are reserved by the Publisher, whether the whole or part of the material is concerned, specifically the rights of translation, reprinting, reuse of illustrations, recitation, broadcasting, reproduction on microfilms or in any other physical way, and transmission or information storage and retrieval, electronic adaptation, computer software, or by similar or dissimilar methodology now known or hereafter developed.

The use of general descriptive names, registered names, trademarks, service marks, 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.

The publisher, the authors and the editors are safe to assume that the advice and information in this book are believed to be true and accurate at the date of publication. Neither the publisher nor the authors or the editors give a warranty, express or implied, with respect to the material contained herein or for any errors or omissions that may have been made.

Printed on acid-free paper

This Springer imprint is published by Springer Nature

The registered company is Springer International Publishing AG

The registered company address is: Gewerbestrasse 11, 6330 Cham, Switzerland

To my beloved family and friends

Dionisios N. Sotiropoulos

*To my wife and colleague, Prof.-Dr. Maria
Virvou, and our daughters, Evina,
Konstantina and Andreani*

George A. Tsirhrintzis

Foreword

There are many real-world problems of such high complexity that traditional scientific approaches, based on physical and statistical modeling of the data generation mechanism, do not succeed in addressing them efficiently. The cause of inefficiencies often lies in multi-dimensionality, nonlinearities, chaotic phenomena and the presence of a plethora of degrees of freedom and unknown parameters in the mechanism that generates data. As a result, loss of information that is crucial to solve a problem is inherent in the data generation process itself, making a traditional mathematical solution intractable.

At the same time, however, biological systems have evolved to address similar problems in efficient ways. In nature, we observe abundant examples of high level intelligence, such as

- Biological neural networks, i.e., networks of interconnected biological neurons in the nervous system of most multi-cellular animals, are capable of learning, memorizing and recognizing patterns in signals such as images, sounds or odors.
- Ants exhibit a collective, decentralized and self-organized intelligence that allows them to discover the shortest route to food in very efficient ways.

A third example of a biological system that exhibits high level intelligence is the vertebrate immune system. The immune system in vertebrates is a decentralized system of biological structures and processes within an organism that protects against pathogens that threaten the organism and may cause disease. Even though it is virtually impossible for a system to learn and memorize all possible forms of pathogens that may potentially threaten an organism, the immune system is capable of detecting a wide variety of agents, ranging from viruses to parasitic worms, and distinguishing them from the organism's own healthy tissue.

Dionisios N. Sotiropoulos and George A. Tsihrintzis have authored the book at hand on computational aspects of the vertebrate immune system with the intent to promote the use of artificial immune systems in addressing machine learning problems. Artificial immune systems are a class of intelligent systems, inspired by the vertebrate immune system and capable of learning and memorizing. In the

recent years, the discipline of computer science has shown intense research interest into applying artificial immune system techniques in various pattern recognition, clustering and classification problems.

The book at hand is a significant addition to this field. The authors present artificial immune systems from the practical signal processing point of view, emphasizing relevant algorithms for clustering and classification. Additionally, the authors illustrate the use of the proposed algorithms through a number of case studies and application on a variety of real data. Particularly interesting is the authors' proposal to use artificial immune system approaches to tackle classification problems which exhibit high, even extreme, class imbalance (so-called *one-class classification problems*).

A practical application of their approach may be found in the design of recommender systems that require and use only positive examples of the preferences of their users. The authors show, through application on real data, that artificial immune systems may be an efficient way to address such problems.

The book is addressed to any graduate student and researcher in computer science. As such, it is self-contained, with the necessary number of introductory chapters on learning and learning paradigms before specific chapters on artificial immune systems. I believe that the authors have done a good job on addressing the tackled issues. I consider the book a good addition to the areas of learning, bio-inspired computing and artificial immune systems. I am confident that it will help graduate students, researchers and practitioners to understand and expand artificial immune systems and apply them in real-world problems.

Dayton, OH, USA
June 2016

Prof.-Dr. Nikolaos G. Bourbakis
IEEE Fellow, President, Biological and
Artificial Intelligence Foundation (BAIF)
OBR Distinguished Professor of Informatics and
Technology and Director
Assistive Technologies Research Center and Director
Assistive Technologies Research Center

Preface

In the monograph at hand, we explore theoretical and experimental justification of the use of *Artificial Immune Systems* as a *Machine Learning Paradigm*. Our inspiration stems from the fact that vertebrates possess an immune system, consisting of highly complex biological structures and processes, that efficiently protect them against disease. A biological immune system is capable of detecting a wide variety of agents, including viruses, parasites and cancer cells, and distinguishing them from the organism's own healthy tissue. This is achieved in the *adaptive immune* subsystem of the immune system.

More specifically, the adaptive immune (sub)system continuously performs a self/non-self discrimination process. In machine learning terms, the adaptive immune system addresses *a pattern classification problem with extreme class imbalance*. Over the recent years, classification problems with class imbalance have attracted the interest of researchers worldwide. However, little attention has been paid so far to the use of artificial immune systems in addressing classification problems with a high or extreme degree of class imbalance.

We address the fundamental problems of pattern recognition, i.e. (*clustering*, *classification* and *one-class classification*), by developing artificial immune system-based machine learning algorithms. We measure the efficiency of these algorithms against state of the art pattern recognition paradigms such as *support vector machines*. Particular emphasis is placed on pattern classification in the context of the class imbalance problem. In machine learning terms, we address degenerated binary classification problems where the class of interest to be recognized is known through only a limited number of positive training instances. In other words, the target class occupies only a negligible volume of the entire pattern space, while the complementary space of negative patterns remains completely unknown during the training process. A practical application of this approach may be found in the design of recommender systems that require the use of only positive examples of the preferences of their users. We show, through application on real data, that artificial immune systems address such problems efficiently.

The general experimentation framework adopted throughout the current monograph is an open collection of one thousand (1000) pieces from ten (10) classes of

western music. This collection has been extensively used in applications concerning music information retrieval and music genre classification. The experimental results presented in this monograph demonstrate that the general framework of artificial immune system-based classification algorithms constitutes a valid machine learning paradigm for clustering, classification and one-class classification.

In order to make the book as self-contained as possible, we have divided it into two parts. Specifically, the first part of the book presents machine learning fundamentals and paradigms with an emphasis on one-class classification problems, while the second part is devoted to biological and artificial immune systems and their application to one-class classification problems. The reader, depending on his/her previous exposure to machine learning, may choose either to read the book from its beginning or to go directly to its second part. It is our hope that this monograph will help graduate students, researchers and practitioners to understand and expand artificial immune systems and apply them in real-world problems.

Piraeus, Greece
June 2016

Dionisios N. Sotiropoulos
George A. Tsirhrintzis

Acknowledgments

We would like to thank Prof.-Dr. Lakhmi C. Jain for agreeing to include this monograph in the Intelligent Systems Reference Library (ISRL) book series of Springer that he edits. We would also like to thank Prof. Nikolaos Bourbakis of Wright State University, USA, for writing a foreword to the monograph. Finally, we would like to thank the Springer staff for their excellent work in typesetting and publishing this monograph.

Contents

Part I Machine Learning Fundamentals

1	Introduction	3
	References	7
2	Machine Learning	9
2.1	Introduction	10
2.2	Machine Learning Categorization According to the Type of Inference	11
2.2.1	Model Identification	12
2.2.2	Shortcoming of the Model Identification Approach	16
2.2.3	Model Prediction	17
2.3	Machine Learning Categorization According to the Amount of Inference	21
2.3.1	Rote Learning	22
2.3.2	Learning from Instruction	22
2.3.3	Learning by Analogy	22
2.4	Learning from Examples	23
2.4.1	The Problem of Minimizing the Risk Functional from Empirical Data	24
2.4.2	Induction Principles for Minimizing the Risk Functional on Empirical Data	26
2.4.3	Supervised Learning	26
2.4.4	Unsupervised Learning	29
2.4.5	Reinforcement Learning	31
2.5	Theoretical Justifications of Statistical Learning Theory	32
2.5.1	Generalization and Consistency	34
2.5.2	Bias-Variance and Estimation-Approximation Trade-Off	36
2.5.3	Consistency of Empirical Minimization Process	38
2.5.4	Uniform Convergence	40

2.5.5 Capacity Concepts and Generalization Bounds	42
2.5.6 Generalization Bounds	46
References	50
3 The Class Imbalance Problem	51
3.1 Nature of the Class Imbalance Problem	51
3.2 The Effect of Class Imbalance on Standard Classifiers	55
3.2.1 Cost Insensitive Bayes Classifier	55
3.2.2 Bayes Classifier Versus Majority Classifier	59
3.2.3 Cost Sensitive Bayes Classifier	66
3.2.4 Nearest Neighbor Classifier	72
3.2.5 Decision Trees	73
3.2.6 Neural Networks	74
3.2.7 Support Vector Machines	76
References	76
4 Addressing the Class Imbalance Problem	79
4.1 Resampling Techniques	79
4.1.1 Natural Resampling	80
4.1.2 Random Over-Sampling and Random Under-Sampling	80
4.1.3 Under-Sampling Methods	80
4.1.4 Over-Sampling Methods	83
4.1.5 Combination Methods	87
4.2 Cost Sensitive Learning	88
4.2.1 The MetaCost Algorithm	89
4.3 One Class Learning	91
4.3.1 One Class Classifiers	91
4.3.2 Density Models	95
4.3.3 Boundary Methods	97
4.3.4 Reconstruction Methods	100
4.3.5 Principal Components Analysis	102
4.3.6 Auto-Encoders and Diabolo Networks	103
References	105
5 Machine Learning Paradigms	107
5.1 Support Vector Machines	107
5.1.1 Hard Margin Support Vector Machines	107
5.1.2 Soft Margin Support Vector Machines	113
5.2 One-Class Support Vector Machines	118
5.2.1 Spherical Data Description	118
5.2.2 Flexible Descriptors	123
5.2.3 v - SVC	124
References	128

Part II Artificial Immune Systems

6 Immune System Fundamentals	133
6.1 Introduction	133
6.2 Brief History and Perspectives on Immunology	134
6.3 Fundamentals and Main Components	137
6.4 Adaptive Immune System	139
6.5 Computational Aspects of Adaptive Immune System	140
6.5.1 Pattern Recognition	140
6.5.2 Immune Network Theory	141
6.5.3 The Clonal Selection Principle	144
6.5.4 Immune Learning and Memory	146
6.5.5 Immunological Memory as a Sparse Distributed Memory	149
6.5.6 Affinity Maturation	151
6.5.7 Self/Non-self Discrimination	154
References	156
7 Artificial Immune Systems	159
7.1 Definitions	159
7.2 Scope of AIS	162
7.3 A Framework for Engineering AIS	164
7.3.1 Shape-Spaces	166
7.3.2 Affinity Measures	168
7.3.3 Immune Algorithms	181
7.4 Theoretical Justification of the Machine Learning Ability of the Adaptive Immune System	183
7.5 AIS-Based Clustering	191
7.5.1 Background Immunological Concepts	191
7.5.2 The Artificial Immune Network (AIN) Learning Algorithm	192
7.5.3 AiNet Characterization and Complexity Analysis	196
7.6 AIS-Based Classification	199
7.6.1 Background Immunological Concepts	200
7.6.2 The Artificial Immune Recognition System (AIRS) Learning Algorithm	202
7.6.3 Source Power of AIRS Learning Algorithm and Complexity Analysis	212
7.7 AIS-Based Negative Selection	214
7.7.1 Background Immunological Concepts	217
7.7.2 Theoretical Justification of the Negative Selection Algorithm	217
7.7.3 Real-Valued Negative Selection with Variable-Sized Detectors	220

7.7.4	AIS-Based One-Class Classification	223
7.7.5	V-Detector Algorithm	225
References.	229
8	Experimental Evaluation of Artificial Immune System-Based Learning Algorithms	237
8.1	Experimentation	237
8.1.1	The Test Data Set	238
8.1.2	Artificial Immune System-Based Music Piece Clustering and Database Organization	242
8.1.3	Artificial Immune System-Based Customer Data Clustering in an e-Shopping Application	248
8.1.4	AIS-Based Music Genre Classification	254
8.1.5	Music Recommendation Based on Artificial Immune Systems	299
References.	321
9	Conclusions and Future Work	325