# DATA MINING IN AGRICULTURE

# Springer Optimization and Its Applications

## VOLUME 34

*Aims and Scope*
Optimization has been expanding in all directions at an astonishing rate during the last few decades. New algorithmic and theoretical techniques have been developed, the diffusion into other disciplines has proceeded at a rapid pace, and our knowledge of all aspects of the field has grown even more profound. At the same time, one of the most striking trends in optimization is the constantly increasing emphasis on the interdisciplinary nature of the field. Optimization has been a basic tool in all areas of applied mathematics, engineering, medicine, economics and other sciences.

The *Springer Optimization and Its Applications* series publishes undergraduate and graduate textbooks, monographs and state-of-the-art expository works that focus on algorithms for solving optimization problems and also study applications involving such problems. Some of the topics covered include nonlinear optimization (convex and nonconvex), network flow problems, stochastic optimization, optimal control, discrete optimization, multi-objective programming, description of software packages, approximation techniques and heuristic approaches.

# DATA MINING IN AGRICULTURE

By

ANTONIO MUCHERINO
University of Florida, Gainesville, FL, USA

PETRAQ J. PAPAJORGJI
University of Florida, Gainesville, FL, USA

PANOS M. PARDALOS
University of Florida, Gainesville, FL, USA

Antonio Mucherino
Institute of Food & Agricultural
Information Technology Office
University of Florida
P.O. Box 110350
Gainesville, FL 32611
USA
amucherino@ufl.edu

Petraq J. Papajorgji
Institute of Food & Agricultural
Information Technology Office
University of Florida
P.O. Box 110350
Gainesville, FL 32611
USA
petraq@ifas.ufl.edu

Panos M. Pardalos
Department of Industrial & Systems Engineering
University of Florida
303 Weil Hall
Gainesville, FL 32611-6595
USA
pardalos@ise.ufl.edu

Printed on acid-free paper

*Dedicated to Sonia
who supported me morally
during the preparation of this book.*

*To the memory of my parents
Eleni and Jorgji Papajorgji
who taught me not to betray my principles
even in tough times.*

*Dedicated to my father and mother
Miltiades and Kalypso Pardalos
for teaching me to love nature
and to grow my own garden.*

# Preface

Data mining is the process of finding useful patterns or correlations among data. These patterns, associations, or relationships between data can provide information about a specific problem being studied, and information can then be used for improving the knowledge on the problem. Data mining techniques are widely used in various sectors of the economy. Initially they were used by large companies to analyze consumer data from different perspectives. Data was then analyzed and useful information was extracted with the goal of increasing profitability.

The idea of using information hidden in relationships among data inspired researchers in agricultural fields to apply these techniques for predicting future trends of agricultural processes. For example, data collected during wine fermentation can be used to predict the outcome of the fermentation while still in the early days of this process. In the same way, soil water parameters for a certain soil type can be estimated knowing the behavior of similar soil types.

The principles used by some data mining techniques are not new. In ancient Rome, the famous orator Cicero used to say *pares cum paribus facillime congregantur* (*birds of a feather flock together* or literally *equals with equals easily associate*). This old principle is successfully applied to classify unknown samples based on known classification of their neighbors. Before writing this book, we thoroughly researched applications of data mining techniques in the fields of agriculture and environmental studies. We found papers describing systems developed to classify apples, separating good apples from bad ones on a conveyor belt. We found literature describing a system that classifies chicken breast quality, and others describing systems able to predict climate forecasting and soil classification, and so forth. All these systems use various data mining techniques.

Therefore, given the scientific interest and the positive results obtained using the data mining techniques, we thought that it was time to provide future specialists in agriculture and environment-related fields with a textbook that will explain basic techniques and recent developments in data mining. Our goal is to provide students and researchers with a book that is easy to read and understand. The task was challenging. Some of the data mining techniques can be transformed into optimization problems, and their solutions can be obtained using appropriate optimization meth-

ods. Although this transformation helps finding a solution to the problem, it makes the presentation difficult to understand by students that do not have a strong mathematical background.

The clarity of the presentation was the major obstacle that we worked hard to overcome. Thus, whenever possible, examples in Euclidean space are provided and corresponding figures are shown to help understand the topic. We make abundant use of MATLAB® to create examples and the corresponding figures that visualize the solution. Besides, each technique presented is ranked using a well-known publication on the relevance of data mining techniques. For each technique, the reader will find published examples of its use by researchers around the world and simple examples that will help in its understanding. We made serious efforts to shed light on when to use the method and the quality of the expected results. An entire chapter is dedicated to the validation of the techniques presented in the book, and examples in MATLAB are used again to help the presentation. Another chapter discusses the potential implementation of data mining techniques in a parallel computing environment; practical applications often require high-speed computing environments. Finally, one appendix is devoted to the MATLAB environment and another one is dedicated to the implementation of one of the presented data mining techniques in C programming language.

It is our hope that readers will find this book to be of use. We are very thankful to our students that helped us shape this course. As always, their comments were useful and appropriate and helped us create a consistent course. We thank Vianney Houles, Guillermo Baigorria, Erhun Kundakcioglu, Sepehr M. Nasseri, Neng Fan, and Sonia Cafieri for reading all the material and for finding subtle inconsistencies. Last but certainly not least, we thank Vera Tomaino for reading the entire book very carefully and for working all exercises. Her input was very useful to us.

Finally, we thank Springer for trusting and giving us another opportunity to work with them.

Gainesville, Florida                                                                      *Antonio Mucherino*
January 2009                                                                              *Petraq J. Papajorgji*
                                                                                          *Panos M. Pardalos*

# Contents

# List of Figures