# Statistical Data Mining and Knowledge Discovery
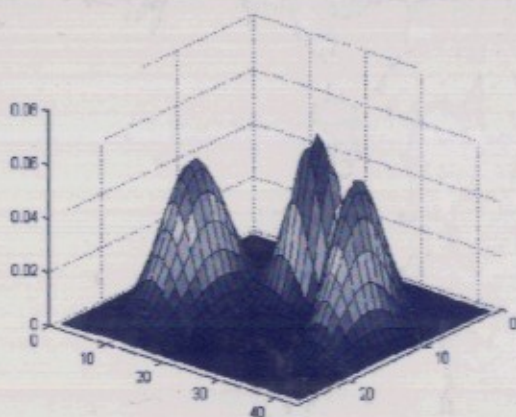


Edited by
## Hamparsum Bozdogan

# Contents

## 3  Econometric and Statistical Data Mining, Prediction and Policy-Making  57

*Arnold Zellner* University of Chicago, Chicago, USA

## 4  Data Mining Strategies for the Detection of Chemical Warfare Agents  79

*Jeffrey. L. Solka,[1,2] Edward J. Wegman,[1] and David J. Marchette[2]*
[2]Naval Surface Warfare Center (NSWCDD), Dahlgren, VA, [1]George Mason University, Fairfax, VA, USA

## 5  Disclosure Limitation Methods Based on Bounds for Large Contingency Tables With Applications to Disability  93

*Adrian Dobra, Elena A. Erosheva and Stephen E. Fienberg* Duke University, Durham, University of Washington, Seattle, and Carnegie-Mellon University, Pittsburgh, USA

**32 Data Mining in Federal Agencies**            **529**

*David L. Banks and Robert T. Olszewski* U.S. Food and Drug Administration,
Rockville, MD, and University of Pittsburgh, Pittsburgh, PA, USA

**33 STING: Evaluation of Scientific & Technological Innovation and**
**Progress**          **549**

*S. Sirmakessis[1], K. Markellos[2], P. Markellou[3], G. Mayritsakis[4],*
*K. Perdikouri[5], A. Tsakalidis[6], and Georgia Panagopoulou[7]* [1-6]Computer
Technology Institute, and[7]National Statistical Services of Greece,
IT Division, Greece