# Statistics for Astrophysics
# Clustering and Classification

Didier Fraix-Burnet & Stéphane Girard (Editors)

**Detailed Table of Contents**

**INTRODUCTION TO R**

Didier Fraix-Burnet

# ELEMENTS OF STATISTICS

Gérard Grégoire

## SOME BASIC ELEMENTS IN CLUSTERING AND CLASSIFICATION

Gérard Grégoire

SUPERVISED AND UNSUPERVISED CLASSIFICATION USING MIXTURE MODELS
Stéphane Girard

**MODEL-BASED CLUSTERING OF HIGH-DIMENSIONAL DATA IN ASTROPHYSICS**
Charles Bouveyron

CLUSTERING OF VARIABLES FOR MIXED DATA

Jérôme Saracco

## INTRODUCTION TO KERNEL METHODS: CLASSIFICATION OF MULTIVARIATE DATA

Mathieu Fauvel

# MODELLING STRUCTURED DATA WITH PROBABILISTIC GRAPHICAL MODELS

Florence Forbes

## CONCEPTS OF CLASSIFICATION AND TAXONOMY – PHYLOGENETIC CLASSIFICATION

Didier Fraix-Burnet