# Contents

**Chapter 2. Likelihood in the Symbolic Context** . . . . . . . . . . . . .   31

Richard EMILION and Edwin DIDAY

**Chapter 3. Dimension Reduction and Visualization of Symbolic Interval-Valued Data Using Sliced Inverse Regression** . . . . . . . . .   49

Han-Ming WU, Chiun-How KAO and Chun-houh CHEN

## Chapter 6. Incremental Calculation Framework for Complex Data . <span>119</span>
Huiwen WANG, Yuan WEI and Siyang WANG

## Part 3. Network Data . . . . . . . . . . . . . . . . . . . . . . . . . . . . . <span>139</span>

## Chapter 7. Recommender Systems and Attributed Networks . . . . <span>141</span>
Françoise FOGELMAN-SOULIÉ, Lanxiang MEI, Jianyu ZHANG, Yiming LI,
Wen GE, Yinglan LI and Qiaofei YE

## Chapter 8. Attributed Networks Partitioning Based on Modularity Optimization . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . 169

David COMBE, Christine LARGERON, Baptiste JEUDY,
Françoise FOGELMAN-SOULIÉ and Jing WANG

## Part 4. Clustering . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . 187

## Chapter 9. A Novel Clustering Method with Automatic Weighting of Tables and Variables . . . . . . . . . . . . . . . . . . . . . . . . . . . 189

Rodrigo C. DE ARAÚJO, Francisco DE ASSIS TENORIO DE CARVALHO
and Yves LECHEVALLIER

## Chapter 10. Clustering and Generalized ANOVA for Symbolic Data Constructed from Open Data . . . . . . . . . . . . . . . . . . . .  209

Simona Korenjak-Černe, Nataša Kejžar and Vladimir Batagelj