# Springer Series in Statistics

# Springer Series in Statistics

Eric D. Kolaczyk

# Statistical Analysis of Network Data

Methods and Models

Eric D. Kolaczyk
Department of Mathematics & Statistics
Boston University
111 Cummington St.
Boston MA 02215
USA

Printed on acid-free paper

springer.com

To the 'network' that is my family . . .

. . . those both near and far.

# Preface

In recent years there has been an explosion of *network data* – that is, measurements that are either of or from a system conceptualized as a network – from seemingly all corners of science. The combination of an increasingly pervasive interest in scientific analysis at a systems level and the ever-growing capabilities for high-throughput data collection in various fields has fueled this trend. Researchers from biology and bioinformatics to physics, from computer science to the information sciences, and from economics to sociology are more and more engaged in the collection and statistical analysis of data from a network-centric perspective.

Accordingly, the contributions to statistical methods and modeling in this area have come from a similarly broad spectrum of areas, often independently of each other. Many books already have been written addressing network data and network problems in specific individual disciplines. However, there is at present no single book that provides a modern treatment of a core body of knowledge for statistical analysis of network data that cuts across the various disciplines and is organized rather according to a statistical taxonomy of tasks and techniques. This book seeks to fill that gap and, as such, it aims to contribute to a growing trend in recent years to facilitate the exchange of knowledge across the pre-existing boundaries between those disciplines that play a role in what is coming to be called 'network science.'

The book is written for students and researchers with a 'mature' knowledge of statistics and hence is intended not only for statisticians but also for people involved with network data in various other areas, like those mentioned above. Background in calculus and linear algebra and some reasonable foundation in statistics and probability are expected. Beyond that, I have attempted to build all necessary material as needed.

In an effort to reach this admittedly diverse audience successfully, I have aimed in each chapter to communicate the material in a manner that strikes an appropriate balance between concepts, on the one hand, and technical depth and rigor, on the other. It is expected that the interested reader will want – and, indeed, is encouraged – to pursue the relevant primary sources for details I may have chosen to omit. The book is in this sense intended to serve as an *entrée* to the larger literature. Copious use of references has been made throughout the book for this very purpose. In addi-

tion, the exercises at the end of each chapter provide further opportunities to explore some of the topics in greater depth. There are both analytical and computational exercises to be found, with the latter frequently designed to be fairly open-ended in nature, so as to encourage exploration. Finally, the methods and models presented herein are illustrated throughout the book with examples from a wide range of disciplines. I have found this overall approach to the pedagogy of the material to work well when I taught classes of precisely the diversity that I envision for the readership of this book.

*Eric D. Kolaczyk*
Boston, Massachusetts
March, 2009

# Contents