# Modern Statistical Methods for Astronomy

## With **R** Applications

ERIC D. FEIGELSON

Pennsylvania State University

G. JOGESH BABU

Pennsylvania State University

# Contents

*The color plates are to be found between pages 398 and 399.*

# Preface

## Motivation and goals

For many years, astronomers have struggled with the application of sophisticated statistical methodologies to analyze their rich datasets and address complex astrophysical problems. On one hand, at least in the United States, astronomers receive little or no formal training in statistics. The traditional method of education has been informal exposure to a few familiar methods during early research experiences. On the other hand, astronomers correctly perceive that a vast world of applied mathematical and statistical methodologies has emerged in recent decades. But systematic, broad training in modern statistical methods has not been available to most astronomers.

This volume seeks to address this problem at three levels. First, we present fundamental principles and results of broad fields of statistics applicable to astronomical research. The material is roughly at a level of advanced undergraduate courses in statistics. We also outline some recent advanced techniques that may be useful for astronomical research to give a flavor of the breadth of modern methodology. It is important to recognize that we give only incomplete introductions to the fields, and we guide the astronomer towards more complete and authoritative treatments.

Second, we present tutorials on the application of both simple and more advanced methods applied to contemporary astronomical research datasets using the **R** statistical software package. **R** has emerged in recent years as the most versatile public-domain statistical software environment for researchers in many fields. In addition to a coherent language for data analysis and common statistical tools, over 3000 packages have been added for advanced analyses in the **CRAN** archive. We have culled these packages for functionalities that may be useful to astronomers. **R** can also be linked to other analysis systems and languages such as C, FORTRAN and Python, so that legacy codes can be included in an **R**-based analysis and *vice versa*.

Third, we hope the book communicates to astronomers our enthusiasm for statistics as a substantial and fascinating intellectual enterprise. Just as astronomers use the latest engineering to build their telescopes and apply advanced physics to interpret cosmic phenomena, they can benefit from exploring the many roads of analyzing and interpreting data through modern statistical analysis.

Another important purpose of this volume is to give astronomers and other physical scientists a bridge to the vast library of specialized texts and monographs in statistics and allied fields. We strongly encourage researchers who are engaged in statistical data analysis to read more detailed treatments in the 'Recommended reading' at the end of each chapter; they are carefully chosen from many available volumes. Most of the material in the book which is not

specifically referenced in the text is presented in more detail in these recommended readings. To further this goal, the present book does not shy away from technical language that, though unfamiliar in the astronomical community, is critical for further learning from the statistical literature. For example, the astronomers' "upper limits" are "left-censored data points", a "power-law distribution" is a "Pareto distribution", and "$1/f$ noise" is a "long-memory process". The text make these connections between the languages of astronomy and statistics, and the comprehensive index can assist the reader in finding material in both languages.

The reader may find the appendices useful. An introduction to **R** is given in Appendix B. It includes an overview of the programming language and an outline of its statistical functionalities, including the many **CRAN** packages. **R** applications to astronomical datasets are given at the end of each chapter which implement methods discussed in the text. Appendix C presents 18 astronomical datasets illustrating the range of statistical challenges that arise in contemporary research. The full datasets and **R scripts** are available online at http://astrostatistics.psu.edu/MSMA. Readers can thus easily reproduce the **R** results in the book.

In this volume, we do not present mathematical proofs underlying statistical results, and we give only brief outlines of a few computational algorithms. We do not review research at the frontiers of astrostatistics, except for a few topics where astronomers have contributed critically important methodology (such as the treatment of truncated data and irregularly spaced time series). Only a small fraction of the many methodological studies in the recent astronomical literature are mentioned. Some fields of applied statistics useful for astronomy (such as wavelet analysis and image processing) are covered only briefly. Finally, only $\sim$2500 **CRAN** packages were examined for possible inclusion in the book; roughly one new package is added every day and many others are extended.

## Audience

The main audience envisioned for this volume is graduate students and researchers in observational astronomy. We hope it serves both as a textbook in a course on data analysis or astrostatistics, and as a reference book to be consulted as specific research problems are encountered. Researchers in allied fields of physical science, such as high-energy physics and Earth sciences, may also find portions of the volume helpful. Statisticians can see how existing methods relate to questions in astronomy, providing background for astrostatistical research initiatives.

Our presentation assumes that the reader has a background in basic linear algebra and calculus. Familiarity of elementary statistical methods commonly used in the physical sciences is also useful; this preparatory material is covered in volumes such as Bevington & Robinson (2002) and Cowan (1998).

## Outline and classroom use

The introduction (Chapter 1) reviews the long historical relationship between astronomy and statistics and philosophical discussions of the relationship between statistical and scientific inference. We then start with probability theory and proceed to lay foundations of statistical inference: hypothesis testing, estimation, modeling, resampling and Bayesian inference

(Chapters 2 and 3). Probability distributions are discussed in Chapter 4 and nonparametric statistics are covered in Chapter 5.

The volume proceeds to various fields of applied statistics that rest on these foundations. Data smoothing is covered in Chapters 5 and 6. Regression is discussed in Chapter 7, followed by analysis and classification of multivariate data (Chapters 8 and 9). Treatments of nondetections are covered in Chapter 10, followed by the analysis of time-variable astronomical phenomena in Chapter 11. Chapter 12 considers spatial point processes. The book ends with appendices introducing the **R** software environment and providing astronomical datasets illustrative of a variety of statistical problems.

We can make some recommendation regarding classroom use. The first part of a semester course in astrostatistics for astronomy students would be devoted to the principles of statistical inference in Chapters 1–4 and learning the basics of **R** in Appendix B. The second part of the semester would be topics of applied statistical methodology selected from Chapters 5–12. We do not provide predefined student exercises with definitive answers, but rather encourage both instructors and students to develop open-ended explorations of the contemporary astronomical datasets based on the **R** tutorials distributed throughout the volume. Suggestions for both simple and advanced problems are given in the dataset presentations (Appendix C).

## Astronomical datasets and R scripts

The datasets and **R** scripts in the book can be downloaded from Penn State's Center for Astrostatistics at http://astrostatistics.psu.edu/MSMA. The **R** scripts are self-contained; simple cut-and-paste will ingest the datasets, perform the statistical operations, and produce tabular and graphical results.

Extensive resources to pursue issues discussed in the book are available on-line. The **R** system can be downloaded from http://www.r-project.org and **CRAN** packages are installed on-the-fly within an **R** session. The primary astronomy research literature, including full-text articles, is available through the NASA–Smithsonian *Astrophysics Data System* (http://adswww.harvard.edu). Thousands of astronomical datasets are available from the Vizier service at the Centre des Données Stellaires (http://vizier.u-strasbg.fr) and the emerging *International Virtual Observatory Alliance* (http://ivo.net). The primary statistical literature can be accessed through MathSciNet (http://www.ams.org/mathscinet/) provided by the American Mathematical Society. Considerable statistical information is available on Wikipedia (http://en.wikipedia.org/wiki/Index_of_statistics_articles). Astronomers should note, however, that the best way to learn statistics is often through textbooks and monographs written by statisticians, such as those in the recommended reading.

# Acknowledgements

*in Statistics for Astronomers* since 2005 and taught at Penn State and Bangalore's Indian Institute of Astrophysics. We are grateful to our dozens of statistician colleagues who have taught at the *Summer Schools in Statistics for Astronomers* for generously sharing their knowledge and perspectives. David Hunter and Arnab Chakraborty developed **R** tutorials for astronomers. Donald Percival generously gave detailed comments on the time series analysis chapter. We are grateful to Nancy Butkovich and her colleagues for providing excellent library services. Finally, we acknowledge the National Science Foundation, National Aeronautics and Space Administration, and the Eberly College of Science for supporting astrostatistics at Penn State over many years.

<div align="right">

Eric D. Feigelson
G. Jogesh Babu
Center for Astrostatistics
Pennsylvania State University
University Park, PA, U.S.A.

</div>