

Why Astrostatistics?

Eric Feigelson
Penn State University

NARIT-EACOA Summer Workshop on Astrostatistics & Astroinformatics
August 2019

What is astronomy?

Astronomy is the observational study of matter beyond Earth: planets in the Solar System, stars in the Milky Way Galaxy, galaxies in the Universe, and diffuse matter between these concentrations.

Astrophysics is the study of the intrinsic nature of astronomical bodies and the processes by which they interact and evolve. This is an indirect, inferential intellectual effort based on the assumption that physics – gravity, electromagnetism, quantum mechanics, etc – apply universally to distant cosmic phenomena.

What is statistics? *(No consensus !!)*

- “... briefly, and in its most concrete form, the object of statistical methods is the reduction of data”
(R. A. Fisher, 1922)
- “Statistics is the mathematical body of science that pertains to the collection, analysis, interpretation or explanation, and presentation of data.”
(Wikipedia, 2014)
- “Statistics is the study of the collection, analysis, interpretation, presentation and organization of data.”
(Wikipedia, 2015)
- “A statistical inference carries us from observations to conclusions about the populations sampled”
(D. R. Cox, 1958)

Does statistics relate to scientific models?

The pessimists ...

“Essentially, all models are wrong, but some are useful.”

(Box & Draper 1987)

“There is no need for these hypotheses to be true, or even to be at all like the truth; rather ... they should yield calculations which agree with observations” (Osiander’s Preface to Copernicus’ *De Revolutionibus*, quoted by C. R. Rao in *Statistics and Truth*)

"The object [of *statistical* inference] is to provide ideas and methods for the critical analysis and, as far as feasible, the interpretation of empirical data ... The extremely challenging issues of *scientific* inference may be regarded as those of synthesising very different kinds of conclusions if possible into a coherent whole or theory ... The use, if any, in the process of simple *quantitative* notions of probability and their numerical assessment is unclear."

(D. R. Cox, 2006)

The positivists ...

“The goal of science is to unlock nature’s secrets. ... Our understanding comes through the development of theoretical models which are capable of explaining the existing observations as well as making testable predictions. ...

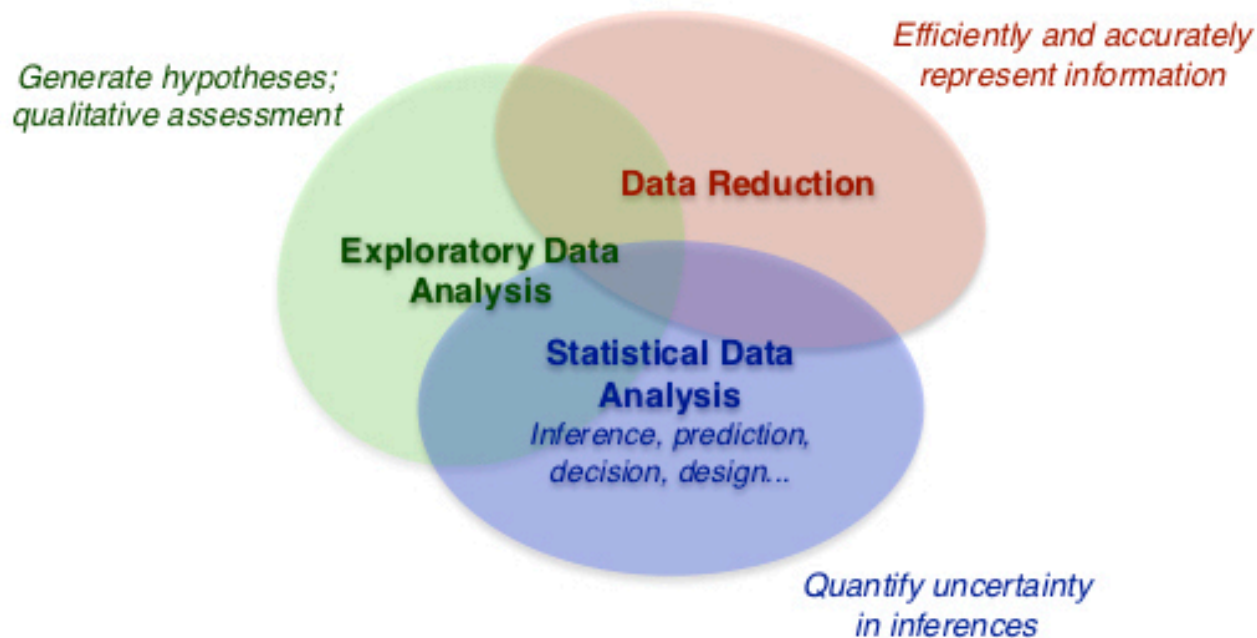
“Fortunately, a variety of sophisticated mathematical and computational approaches have been developed to help us through this interface, these go under the general heading of statistical inference.”

(P. C. Gregory, *Bayesian Logical Data Analysis for the Physical Sciences*, 2005)

Data analysis

Building & Appraising Arguments Using Data

Modes of Data Analysis



Inference: Learning about the data generating process (population, signals...) from observed data—just one of several interacting modes of analyzing data

Recommended steps in the statistical analysis of scientific data

The application of statistics can reliably quantify information embedded in scientific data and help adjudicate the relevance of theoretical models. But this is not a straightforward, mechanical enterprise. It requires:

- exploration of the data
- careful statement of the scientific problem
- model formulation in mathematical form
- choice of statistical method(s)
- calculation of statistical quantities ← *easiest step with R*
- judicious scientific evaluation of the results

Astronomers often do not adequately pursue each step

- Modern statistics is vast in its scope and methodology. It is difficult to find what may be useful (jargon problem!), and there are usually several ways to proceed. Very confusing.
- Some statistical procedures are based on mathematical proofs which determine the applicability of established results. It is perilous to violate mathematical truths! Some issues are debated among statisticians, or have no known solution.
- Scientific inferences should not depend on arbitrary choices in methodology & variable scale. Prefer nonparametric & scale-invariant methods. Try multiple methods.
- It can be difficult to interpret the meaning of a statistical result with respect to the scientific goal. Statistics is only a tool towards understanding nature from incomplete information.

***We should be knowledgeable in our use of statistics
and judicious in its interpretation***

Astronomy & Statistics: A glorious past

*For most of western history,
the astronomers were the statisticians!*

Ancient Greeks to 18th century

Best estimate of the length of a year from discrepant data?

- Middle of range: Hipparchos (4th century B.C.)
- Observe only once! (medieval)
- Mean: Brahe (16th c), Galileo (17th c), Simpson (18th c)
- Median with bootstrap (21st c)

19th century

Discrepant observations of planets/moons/comets used to estimate orbital parameters using Newtonian celestial mechanics

- Legendre, Laplace & Gauss develop least-squares regression and normal error theory (~1800-1820)
- Prominent astronomers contribute to least-squares theory (~1850-1900)

The lost century of astrostatistics....

In the late-19th and 20th centuries, statistics moved towards human sciences (demography, economics, psychology, medicine, politics) and industrial applications (agriculture, mining, manufacturing).

During this time, astronomy recognized the power of modern physics: electromagnetism, thermodynamics, quantum mechanics, relativity. Astronomy & physics were wedded into astrophysics.

Thus, astronomers and statisticians substantially broke contact; e.g. the curriculum of astronomers heavily involved physics but little statistics. Statisticians today know little modern astronomy.

The state of astrostatistics today

(not good ... but getting better!)

Many astronomical studies are confined to a narrow suite of familiar statistical methods:

- Fourier transform for temporal analysis (Fourier 1807)
- Least squares regression (Legendre 1805, Pearson 1901)
- Kolmogorov-Smirnov goodness-of-fit test (Kolmogorov, 1933)
- Principal components analysis for tables (Hotelling 1936)

Even traditional methods are sometimes misused!

see Beware the Kolmogorov-Smirnov test! page on ASAIP

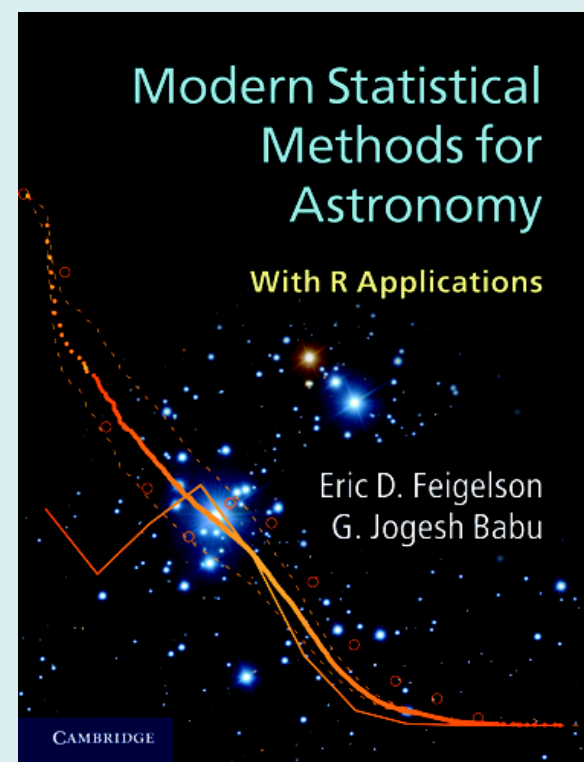
Under-utilized methodology:

- modeling (MLE, EM Algorithm, BIC, bootstrap)
- multivariate classification (LDA, SVM, CART, RFs)
- time series (autoregressive models, state space models)
- spatial point processes (Ripley's K, kriging)
- nondetections (survival analysis)
- image analysis (computer vision methods, False Detection Rate)
- statistical computing (R)

Advertisement ...

Modern Statistical Methods for Astronomy with R Applications

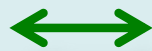
E. D. Feigelson & G. J. Babu,
Cambridge Univ Press, 2012



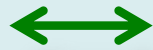
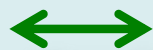


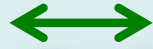







*Winner 2012 PROSE Award for
Best Astronomy & Cosmology Book*

An astrostatistics lexicon ...

Cosmology



Statistics

| | | |
|-------------------------------|--|--------------------------------------|
| Galaxy clustering |  | Spatial point processes, clustering |
| Galaxy morphology |  | Regression, mixture models |
| Galaxy luminosity fn |  | Gamma distribution |
| Power law relationships |  | Pareto distribution |
| Weak lensing morphology |  | Geostatistics, density estimation |
| Strong lensing morphology |  | Shape statistics |
| Strong lensing timing |  | Time series with lag |
| Faint source detection |  | False Discovery Rate |
| Multiepoch survey lightcurves |  | Multivariate classification |
| CMB spatial analysis |  | Markov fields, ICA, etc |
| Λ CDM parameters |  | Bayesian inference & model selection |
| Comparing data & simulation |  | <i>Uncertainty Quantification</i> |

Recent resurgence in astrostatistics

- Improved access to statistical software. R/CRAN, Matlab & Python
- Papers in astronomical literature doubled to ~500/yr in past decade
- Short training courses (Penn State, India, Brazil, Greece, China, Italy, France, Germany, Spain, Sweden, IAU/AAS/CASCA/... meetings)
- Cross-disciplinary research collaborations (Harvard, Carnegie-Mellon, Penn State, CEA-Saclay/, Cornell, Imperial College London, Swinburne, Yale, ...)
- Cross-disciplinary conferences (*Statistical Challenges in Modern Astronomy*, *Astronomical Data Analysis 1991--*, SAMSI 2006/2012, *Astronomical Data Analysis (1995--)*, *Astroinformatics (2012--)*)
- Scholarly societies:
 - International Stat Institute SIGAstro
 - International Astrostatistical Association
 - International Astro Union Commission B3
 - American Astro Soc Working Group
 - American Stat Assn Interest Group
 - LSST Science Collaboration
 - IEEE Astro Data Miner Task Force
 - International Astroinformatics Association

To treat massive data streams and databases ...

Rapid rise of astroinformatics

Statistics guides the scientist on what to compute

Informatics helps the scientist perform the computation

Methodology: Computationally intensive astronomy, data mining, multivariate regression & classification, machine learning, Monte Carlo methods, NlogN algorithms, etc.

Software & hardware: Parallel processing on multi-processors machines, cloud computing, CUDA & GPU computing, database management & promulgation, data streams, etc.

Workshops & training schools emerging. IAU Symposium 2016, IEEE Symposium 2018. Growing perception that more community training is needed.

New resources in astrostatistics

General textbooks

Practical Statistics for Astronomers

Wall & Jenkins, 2nd ed, 2012

Modern Statistical Methods for Astronomy with R Application,

Feigelson & Babu, 2012

Statistics, Data Mining, and Machine Learning in Astronomy: A Practical Python Guide for the Analysis of Survey Data,

Ivecic, Connolly, VanderPlas & Gray, 2014

Bayesian textbooks

Bayesian Logical Data Analysis for the Physical Sciences: A Comparative Approach with Mathematica Support

Gregory, 2005

Bayesian Models for Astrophysical Data: Using R, JAGS, Python, and Stan


Hilbe, de Souza & Ishida 2017

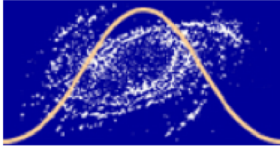
Practical Bayesian Inference: A Primer for Physical Scientists

Bailer-Jones 2017

Astrostatistics and Astroinformatics Portal

<http://asaip.psu.edu>

**Astrostatistics and
Astroinformatics Portal (ASAIP)**



Log in

Search This Site

Go

Search ASAIP-related sites

[Home](#) [Members](#) [Recent Papers](#) [Resources](#) [Organizations](#) [Articles](#) [Forums](#) [Meetings](#) [ASAIP updates](#)

You are here: [Home](#)

Welcome to ASAIP

The Astrostatistics and Astroinformatics Portal (<http://asaip.psu.edu>) is a new Web site serving the cross-disciplinary communities of astronomers, statisticians and computer scientists. It is intended to foster research into advanced methodologies for astronomical research, and to promulgate such methods into the broader astronomy community. The WWW public is welcome to read materials in ASAIP. Use the navigation bar above, or the search box at the upper right, to find material throughout the ASAIP Web site.

Meetings, jobs, societies, ... other resources

A vision of astrostatistics by 2025 ...

- Astronomy graduate curriculum has 1 year of statistical and computational methodology
- Some astronomers have M.S. in statistics and computer science
- Astrostatistics and astrophysics is a well-funded, cross-disciplinary research field involving a few percent of astronomers (cf. astrochemists) pushing the frontiers of methodology.
- Astronomers regularly use many methods coded in R.
- *Statistical Challenges in Modern Astronomy* meetings are held annually with ~300 participants

Prelude to R

A brief history of statistical computing

1960s – c2000: Statistical analysis developed by academic statisticians, but implementation relegated to commercial companies (SAS, BMDP, Statistica, Stata, Minitab, etc).

1980s: John Chambers (ATT, USA) develops S system, C-like command line interface.

1990s: Ross Ihaka & Robert Gentleman (Univ Auckland NZ) mimic S in an open source system, R. R Core Development Team expands, GNU GPL release.

Early–2000s: Comprehensive R Analysis Network (CRAN) for user–provided specialized packages grows exponentially. Important packages incorporated into base–R.