

21st Century Statistical and Computational Challenges in Astrophysics

Eric D. Feigelson^{1,2,3}, Rafael S. de Souza⁴,
Emille E. O. Ishida⁵, and Gutti Jogesh Babu^{2,1,3}

¹ Department of Astronomy & Astrophysics, Penn State University, University Park PA, USA, 16802, e5f@psu.edu

² Department of Statistics, Penn State University, University Park PA, USA, 16802

³ Center for Astrostatistics, Penn State University, University Park PA, USA, 16802

⁴ Key Laboratory for Research in Galaxies and Cosmology, Shanghai Astronomical Observatory, Chinese Academy of Sciences, 80 Nandan Road, Shanghai 200030, China

⁵ Université Clermont Auvergne, CNRS/IN2P3, LPC, F-63000 Clermont-Ferrand, France

XXXX. XXX. XXX. XXX. YYYY. AA:1–27

[https://doi.org/10.1146/\(\(please add article doi\)\)](https://doi.org/10.1146/((please add article doi)))

Copyright © YYYY by Annual Reviews.
All rights reserved

Keywords

astronomy, astrophysics, astrostatistics, cosmology, galaxies, stars, exoplanets, gravitational waves, Bayesian inference, likelihood-free modeling, signal detection, periodic time series, machine learning, measurement errors

Abstract

Modern astronomy has been rapidly increasing our ability to see deeper into the universe, acquiring enormous samples of cosmic populations. Gaining astrophysical insights from these datasets requires a wide range of sophisticated statistical and machine learning methods. Long-standing problems in cosmology include characterization of galaxy clustering and estimation of galaxy distances from photometric colors. Bayesian inference, central to linking astronomical data to nonlinear astrophysical models, addresses problems in solar physics, properties of star clusters, and exoplanet systems. Likelihood-free methods are growing in importance. Detection of faint signals in complicated noise is needed to find periodic behaviors in stars and detect explosive gravitational wave events. Open issues concern treatment of heteroscedastic measurement errors and understanding probability distributions characterizing astrophysical systems. The field of astrostatistics needs increased collaboration with statisticians in the design and analysis stages of research projects, and to jointly develop new statistical methodologies. Together, they will draw more astrophysical insights into astronomical populations and the cosmos itself.

Contents

1. ASTRONOMY, ASTROPHYSICS, AND ASTROSTATISTICS.....	2
2. TWO LONG-STANDING PROBLEMS IN COSMOLOGY.....	4
2.1. Galaxy Clustering	4
2.2. The Photo-z Conundrum	6
3. BAYESIAN MODELING OF THE SUN, STARS, AND PLANETARY SYSTEMS.....	8
3.1. Solar Magnetic Fields	8
3.2. Star Cluster Properties	8
3.3. Modeling Multiplanet Systems.....	10
4. LIKELIHOOD-FREE MODELING	10
5. CHALLENGES IN SIGNAL DETECTION	12
5.1. Periodicity Detection in Irregular Time Series.....	12
5.2. Gravitational Wave Detection	13
6. TWO EXAMPLES OF MACHINE LEARNING IN ASTRONOMY	16
6.1. Photometry of Blended Galaxies.....	16
6.2. The Accelerating Expansion of the Universe	17
7. TWO CLASSIC PROBLEMS IN ASTROSTATISTICS.....	18
7.1. Heteroscedastic Measurement Errors	18
7.2. Domain-Specific Probability Distributions	19
8. ASTROSTATISTICAL CHALLENGES FOR STATISTICIANS.....	20

1. ASTRONOMY, ASTROPHYSICS, AND ASTROSTATISTICS

Astronomy, the oldest science, has had profound relationships to statistical concepts since antiquity. Hipparchus chose the midrange estimator to reconcile inconsistent measurements of the duration of a year; Ptolemy and al-Biruni instead chose the mean value in their celestial calculations (Sheynin 1974). Galileo outlined a theory of observational errors in a discussion of the Comet of 1572. Newton had little interest in probabilistic arguments but his followers developing celestial mechanics laid the foundations of modern statistics. Laplace minimized both the sum of absolute and squared residuals between observations and model predictions, but the L_2 method was integrated to the Gaussian error distribution in least squares theory (Stigler 1986). Many leading astronomers contributed to least squares methodology during the 19th century.

But the links between astronomy and statistics sundered during the 20th century as the former turned to physics and the latter to human affairs. Statistician involvement in astronomy was rare and unsuccessful. Pearson & Bell (1908) presented a study quantifying a strong correlation between the range of variability and brightness of variable stars, but it was subject to confusion between observed brightness and intrinsic luminosity and the mixing of pulsating and eclipsing variables. Decades later, Neyman & Scott (1952, 1958) developed a model for galaxy clustering involving power law clusters placed at seed locations obtained by a Poisson point process. However, it proved to be much too simple to describe for the hierarchical anisotropic spatial distribution of galaxies that emerged later from redshift surveys (de Lapparent et al. 1986).

The divorce of astronomers from established methodology was at times severe. Schlesinger (1910) had to plead with his community that least squares solutions to spectroscopic binary star orbits had a stronger foundation than approximate subjective meth-

ods. Hubble (1930) obtained fits of elliptical galaxy light distributions to a spherical self-gravitating model by trial-and-error, and Zwicky (1937) made the seminal discovery of Dark Matter in the Coma cluster of galaxies with a curve fitted by eye. The first widely recognized application of maximum likelihood estimation did not appear until Lynden-Bell et al. (1988) applied it to a dynamical model of galaxy redshifts, discovering a previously unrecognized nearby concentration of galaxies known as the Great Attractor.

Connections between the fields began to reestablish themselves in the 1990s and have rapidly grown. In the astronomy research literature, terms like Bayesian and machine learning have increased exponentially in the past decade, with Deep Learning rising meteorically since 2017. While the number of professional statisticians working on astronomical problems is still small, the interest in advanced methodology – particularly machine learning – in the astronomical community is very strong.

Observational astronomy today constitute a considerable enterprise with billions of dollars supporting ~20,000 scientists producing ~15,000 refereed papers annually. The science is devoted to the characterization and understanding of phenomena outside of Earth: our Sun and Solar System, other stars and their planetary systems, the Milky Way Galaxy and other galaxies, diffuse material between the stars and galaxies, and the Universe as a whole. Progress is propelled by rapid development of technologies improving our observations at wavelengths from the longest radio waves to the shortest gamma-rays. Not all discoveries involve electromagnetic waves. Telescopes with strange designs detect energetic particles like cosmic rays and neutrinos, and most recently gravitational waves in space-time itself, give unique insights into explosive phenomena across the Universe.

Theoretical astrophysics seeks to interpret telescopic results using physical processes known from terrestrial studies. An astonishing range of physics is involved in cosmic phenomena: gravitational physics and fundamental theory; atomic and nuclear physics; thermodynamics, hydrodynamics and magnetohydrodynamics; molecular and solid state physics. Many observations cannot be closely linked to physics and are interpreted using heuristic statistical models like linear regressions; these are often power law relationships because variables are plotted with logarithmic transformation.

But in other cases, convincing physical explanations are available and data are used to find best-fit parameters of complicated nonlinear astrophysical models. One of the most important and successful astrophysical models in recent decades is the Λ CDM cosmological model: the expanding Universe with attractive cold Dark Matter and repulsive Dark Energy. It gives a fundamental understanding of the evolution of the Universe from the Big Bang 13.7 billion years ago to the present day and into the future. Ordinary matter of stars, planets and people constitute only a small fraction of the ‘stuff’ in the Universe that is dominated by enigmatic, invisible material and forces. An indicator of the powerful links between astronomy and physics are the eleven Nobel Prizes in Physics awarded for astrophysics in the past 50 years.

While still a smaller enterprise than astrophysics, astrostatistics is playing an increasing role in the analysis of astronomical observations and linking data to astrophysical theory. Consider the field of time domain astronomy. While the stars seen with the naked eye seem unchanging throughout our lives, in fact many objects have variable characteristics from exoplanets orbiting nearby stars to accreting black holes in distant quasars. Enormous investment in telescopes for repeated measurements over time are made, such as the Vera C. Rubin Observatory now under construction in the high Atacama Desert of Chile. Its main task will be an astronomical survey, the Legacy Survey of Space and Time (LSST),

essentially a decade-long ‘movie’ of variable objects in the sky.

In this review, we present a non-comprehensive selection of issues important to current understanding of cosmic phenomena where progress seems impossible without sophisticated statistical analysis. In some cases, astrostatisticians have had considerable success with established methods. In other cases, new developments are underway or the problems need creative ideas. Our approach complements the more integrative review of Schafer (2015). Although our treatment is very incomplete, we hope to communicate to statisticians the unusual intellectual culture of astronomy with its amazing instruments, rapid progress, fascinating science, and exciting methodological challenges.

2. TWO LONG-STANDING PROBLEMS IN COSMOLOGY

2.1. Galaxy Clustering

The distribution of galaxies in space proves to be surprising complex from the viewpoint of spatial point processes. Figure 1 shows the distribution of galaxies in a slice of the sky out to redshift 0.14, equivalent to distances \sim 600 megaparsecs (Mpc, a parsec is about 3.26 lightyears). The dataset pictured here is part of the Sloan Digital Sky Survey (SDSS), one of the most successful astronomical projects of modern times that produced thousands of important studies.

The galaxy distribution is very nonstationary and anisotropic on smaller scales (< 200 Mpc) but mostly stationary on large scales. The pattern, commonly called the ‘large-scale structure’ of the Universe, roughly resembles a collection of contiguous soap bubbles where galaxies are distributed along curves ‘filaments’ and sheets surrounding ‘voids’ (Zeldovich et al. 1982). Particularly prominent filaments are called ‘Great Walls’. Rich clusters appear at the intersections of filaments and are sometimes collected into ‘superclusters’ that can include tens of thousands of galaxies. The pattern is far more complex than statistically established stationary models like Neyman-Scott, Matérn or Cox processes (Baddeley et al. 2015).

Despite this complexity, the galaxy two-point (pair) correlation function is a simple power law (Pareto) function with a universal slope from 0.01 to 100 Mpc (Peebles 1973). This, however, does not prove to be a powerful discriminant between cosmological models. But the detection of a faint bump in the correlation function around 300 Mpc separation validated an important prediction of specific Big Bang theories. Known as the baryonic acoustic oscillation signal, its discovery was one of the major achievements of the SDSS (Eisenstein et al. 2005).

Many statistical studies of large-scale structure rely on isotropic two- and three-point correlation functions as well as Fourier power spectra. But other studies seek to locate particular clusters, filaments or voids. Several filament finding techniques have been investigated. The pruned Minimal Spanning Tree is a popular procedure (Barrow et al. 1985) but the resulting filaments are often noisy with spurs and chaining. Stoica et al. (2010) develop a model based on point processes marked by a filament identifier with probabilities favoring networks of multiply-connected aligned segments. Global Bayesian solutions are found using a simulated annealing algorithm based on Metropolis-Hastings dynamics. Other filament-finding algorithms are based on scalar measures of size and shape (Sahni et al. 1998), thresholded loci of density saddle points (Novikov et al. 2006), segmented watershed transform (Platen et al. 2007), measures based on the Hessian eigenvalues (Bond et al. 2010), and density ridges traced by a subspace constrained mean shift algorithm (Chen

et al. (2015) (Moews et al. 2020).

Some of these algorithms are applied to the observed point spatial distribution and others are applied to its smoothed density estimator. The latter situation also occurs when filamentary structures appear in continuous media traced in real-valued images rather than point processes. Filaments are characteristic in hot flare plasma on the surface of the Sun and in cold molecular clouds within the Milky Way Galaxy. Men'shchikov et al. (2010) defines molecular cloud filaments using the tint fill algorithm from image processing; they are astrophysically explained as outcomes of supersonic magnetohydrodynamic turbulence (Beattie & Federrath 2019).

In the case of galaxy clustering, powerful simulations of clustering patterns are available from astrophysical models involving local gravitational contraction within the expanding Universe following the Big Bang. These are computationally intensive calculations with names like the *Millenium* (Springel et al. 2005) and *Illustris* (Nelson et al. 2015) simulations. The models and observed galaxy clustering patterns are in considerable agreement, validating the dominance of Dark Matter in large-scale structure formation (Figure 2). Important issues regarding bias in visible galaxy formation and quasar feedback are still under

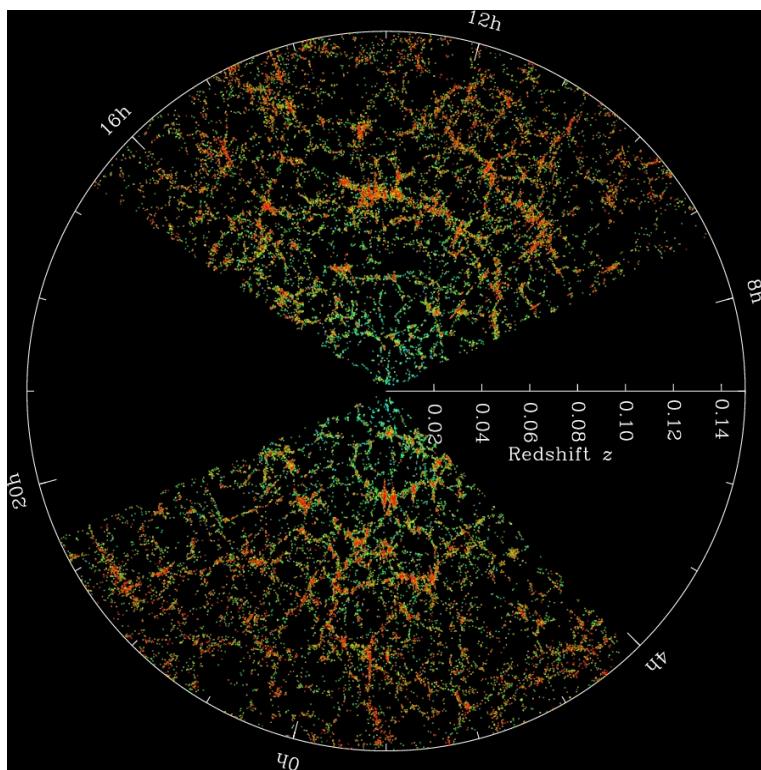


Figure 1

Observed large scale structure from the Sloan Digital Sky Survey main galaxy redshift sample. The slice is 2.5° thick, and extends to redshift 0.14. Galaxies are color-coded by g-r color. Image from <http://classic.sdss.org/legacy>

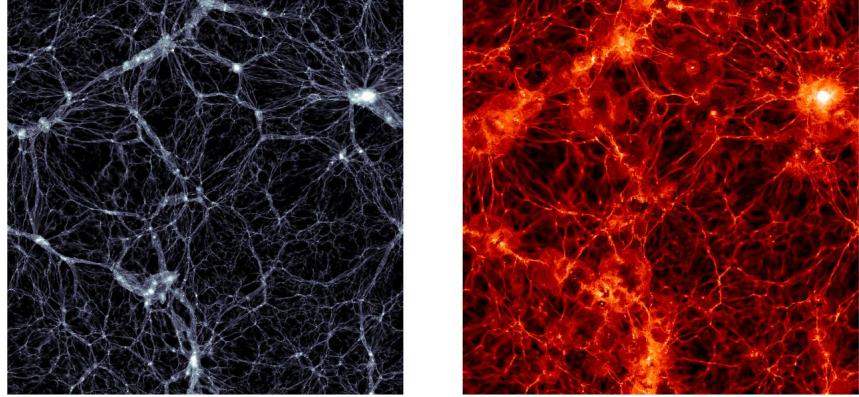


Figure 2

Astrophysical simulation of large-scale structure. A 2-dimensional slice today (13.7 billion years after the Big Bang) of a small portion of the 3-dimensional Illustris calculation of the growth of large-scale structure in the expanding Universe. The images show Dark Matter (left) and ordinary matter (right) density covering an area of $\sim 100 \times 100 \text{ Mpc}^2$ [Haider et al. 2016].

investigation.

But important methodological problems remain. First, there has been little comparison of the performance of the various filament-finding methods or other measures of clustering. Most are algorithmic with little foundation in mathematical or statistical theory. Second, astronomers have few tools to compare complicated observed and simulated two- and three-dimensional point distributions with complexities, such as Figures 1 and 2. These tests must treat millions of points for quantitative comparison of model predictions and observed data.

2.2. The Photo-z Conundrum

It is relatively easy in astronomy to measure locations on the celestial sphere but very difficult to measure distances. Distances to stars in our region of the Milky Way Galaxy can be obtained from annual parallax measurements, a method known to Copernicus though not achieved until the 19th century. But distances to other galaxies must rely in indirect procedures. For galaxies lying billions of parsecs away, estimating distances is paramount to understand their evolution across cosmic time. With the best telescopes, we can now see galaxies forming only ~ 1 Gyr after the Big Bang. But as the Universe is expanding, the spectrum is shifted towards longer (redder) observed wavelengths (λ_{obs}) compared to their rest-frame λ_{rest} wavelengths. The redshift is defined to be $z = (\lambda_{obs} - \lambda_{rest})/\lambda_{rest}$. The astronomer obtains a spectrum of a galaxy with a spectrograph on a telescope to measure the wavelength in which specific spectral features are found. This is compared with the expected wavelength for that feature in rest frame. For example, a galaxy with redshift $z = 3$ has its hydrogen Balmer break shifted from the blue around 400 nm to the infrared around 1200 nm.

Once the expansion rate of the Universe, known as Hubble's constant H_0 , is measured, then the galaxy's distance in parsecs can be inferred from the redshift. After great effort over several decades, H_0 is now known to be around 70 km/s/Mpc to within $\sim 2\%$. Once the distance is obtained, then the galaxy's size, luminosity, mass, cluster environment, and

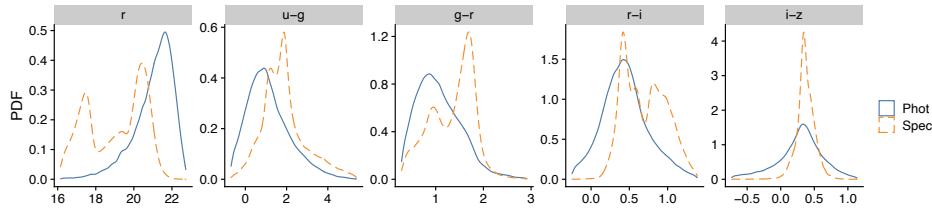


Figure 3

Distributions for magnitude in r-band and four colors (u-g, g-r, r-i, i-z) for the spectroscopic (training) sample compared with the photometric (target) sample in the data sets presented by Beck et al. (2017).

stellar mass function are inferred from observed properties, as well as its chemical, merger and star formation history. Obtaining accurate redshifts is thus essential for understanding the formation and properties of galaxies and to trace the evolution of large-scale structures (§2.1).

However, it is expensive in telescope time to obtain sufficiently high-resolution spectra in order to measure spectral lines and directly obtain redshifts. It is much cheaper to measure brightnesses of many galaxies simultaneously in wide spectral bands; this is called photometry and produces a low-resolution ‘spectral energy distribution’ or SED. Here each band blurs together many spectral lines with continuum starlight. A wide-field imaging camera may obtain photometry 100-fold more efficiently than a multi-object spectrograph for faint high-redshift galaxies (Hildebrandt et al. 2008). The redshift information is encoded in the photometry but not as directly as in a high-resolution spectrum.

Therefore, photometric redshift (photo- z) estimation has become a vital tool in the extragalactic astronomy and observational cosmology. The challenge of photo- z accuracy then depends on the statistical procedures used to calibrate photometric measurements to spectroscopic redshifts.

A plethora of methods have been proposed and tested for this task (Benítez 2000, Beck et al. 2016, Budavári 2009, Elliott et al. 2015, Cavaoti et al. 2016). They range from several variants of least squares weighted by photometric measurement errors to Bayesian inference (Leistedt et al. 2016), k-nearest neighbor procedures, kernel density estimation, Gaussian mixture models, generalized linear models, Self-Organizing Maps, Random Forest (Carliles et al. 2010), hyper-optimized gradient boosting regression, Gaussian processes regression, and hybrid schemes. Galaxy morphology from images can supplement photometric measurements when neural networks are used. More recently, some more advanced methods have been used, including the use of different flavors of deep convolutional networks to derive photometric redshift directly from multi-band images (D’Isanto & Polsterer 2018, Pasquet et al. 2019), and to build a morphology-aware photo- z estimator (Menou 2019). Comparisons of performance are presented by Dahlen et al. (2013), Rau et al. (2015), Salvato et al. (2019), and Schmidt et al. (2020).

However most methods fail to achieve better than $\sim 2\%$ accuracy. The relationship between photometry and spectroscopic redshifts are not only nonlinear, but degeneracies, heteroscedasticity, and ‘catastrophic outliers’ abound in the datasets. A further caveat encountered by machine learning based methods is a common mismatch between the pa-

rameter distribution of the training (spectroscopic) and target (photometric) datasets as portrayed in Figure 3. This is connected with the nature of spectroscopic measurements which, within the same survey, demands higher quality data and consequently delivers lower redshift and higher metallicities galaxies than the photometric counterpart.

In cases where spectroscopic and photometric samples have similar coverage in magnitude/color space, it is feasible to adapt the spectroscopic sample using Domain Adaptation (Beck et al. 2017). However, in most real-data scenarios this assumption will not hold. An alternative approach relies on active learning to improve the training set by making sensible decisions to query extra galaxies from which one could measure the spectra (Vilalta et al. 2017).

3. BAYESIAN MODELING OF THE SUN, STARS, AND PLANETARY SYSTEMS

3.1. Solar Magnetic Fields

Studies of the Sun are among the active fields of astrophysical modeling, particularly addressing the many manifestations of its magnetic field generated by gas motions in its interior. Asensio Ramos (2006) presents a Bayesian model for the spatial variations of linear and circular polarization which can be linked to the underlying magnetic field strength through physical processes like radiative transfer and the Zeeman effect. The method was applied to a polarization map of a small portion of the quiet Sun (Figure 4) to investigate the role of magnetic fields in granulation arising from convection in its upper layers. The calculation is made for each pixel in the map, and maps are constructed of the Kullback-Liebler divergence between the prior and posterior distributions for several physical quantities (such as the magnetic field strength, filling factor, and geometry) after marginalization over other parameters. The resulting maps are similar to those obtained by traditional weighted least squares calculations (designated ‘ χ^2 minimization’ in the astronomical community) but with improved values for magnetic field strength and treatment of degeneracies between the physical parameters.

3.2. Star Cluster Properties

Over a century ago, it was discovered that most stars lie along a curved locus in the Hertzsprung-Russell (H-R) diagram, a plot of log-luminosity against surface temperature or color. As nuclear physics developed, astrophysical models were constructed showing that this ‘main sequence’ represents the long-lived phase powered by the fusion of hydrogen to helium in the stellar cores. When hydrogen is depleted the cores, the stars move off the main sequence to higher luminosities and redder surface colors, the ‘red giant’ regime. The theory of stellar evolution has been well-developed for decades; the evolution of a star in the H-R diagram can be calculated based on nuclear, gravitational, atomic and fluid physics for a given mass and abundance of elements, summarized in a parameter called ‘metallicity’ (Z). Astronomers now use the evolutionary models to infer precise stellar properties that are otherwise difficult to ascertain, such as their ages (τ) and masses (M).

A long-standing astrostatistical collaboration has applied Bayesian inference to this classical problem in stellar astronomy (van Dyk et al. 2009). The data are easily obtained from photometric measurements of the brightness B_{ij} in several color bands j for $i = 1, \dots, N$ stars in a coeval star cluster, together with measurement errors $\sigma(B_{ij})$. The likelihood is a

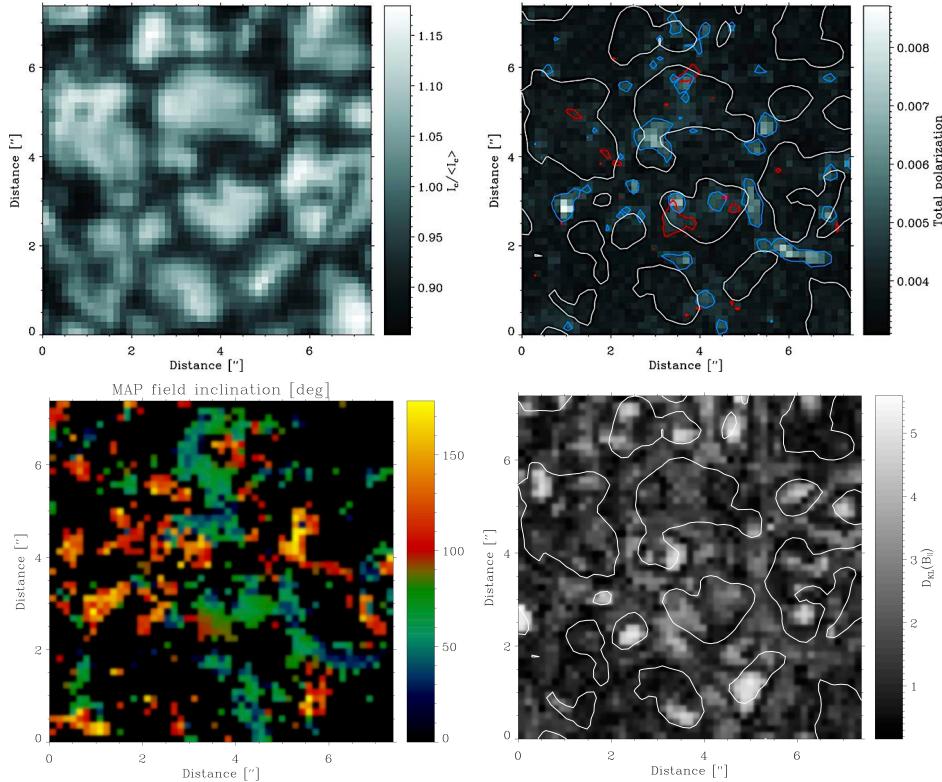


Figure 4

Bayesian analysis of magnetic fields in the quiet Sun's photosphere (Asensio Ramos 2009). *Top left:* Map of the continuum brightness for a small region of the Sun. *Top right:* Map of the polarization brightness where the white contours mark areas of bright continuum, red contours mark bright linear polarization, and blue contours mark bright circular polarization. *Bottom left:* Maximum-a-posterior map of the magnetic field inclination inferred from the Bayesian model. *Bottom right:* Map of the Kullback-Leibler divergence between prior and posterior for the magnetic field inclination.

product of Gaussians,

$$\mathcal{L}(M, \tau, Z | B, \sigma(B)) \propto \prod_{j=1}^k \prod_{i=1}^N \exp \left[\frac{-[B_{ij} - G_j(M, \tau, Z)]^2}{2\sigma^2(B_{ij})} \right], \quad 1.$$

where G_j are the predicted brightness from astrophysical stellar evolution models. Astronomers have an informative prior distribution for stellar masses; stellar Initial Mass Function (IMF) is known to follow a log-normal distribution for stars in the range $0.1 - 8$ solar masses (§7.2). However, there is less knowledge about the star formation history of the Galaxy, so an uninformative uniform prior in log-age over a wide range is assumed. Gaussian priors are given to other stellar parameters such as metallicity. Posterior values are calculated by Markov chain Monte Carlo using a one-at-a-time Gibbs sampler with a uniform Metropolis jumping rule. This is a high-dimensional model with $3N+5$ parameters;

a large look-up table of G values is used to reduce stellar model calculations. More complex likelihood models can be used to account for contamination by irrelevant field stars or for multiple age components in the star cluster.

The method is now widely used to estimate ages of star clusters in the Milky Way galaxy. [Bossini et al. \(2019\)](#) find ages ranging from 10 to 10,000 million years for hundreds of star clusters newly identified from the Gaia catalog of 1.3 billion Galactic stars. The estimated masses agree very closely to values obtained by a more accurate method (modeling of asteroseismological oscillations) available for a handful of stars. New star clusters have been recently discovered ([Cantat-Gaudin et al. 2019](#)), and there are many more to come.

3.3. Modeling Multiplanet Systems

One of the great excitements of modern astronomy is the clear evidence that most stars in the sky have their own planetary systems. The original discovery was based on radial velocity measurement of the host star by [Mayor & Queloz \(1995\)](#) who were consequently awarded the 2019 Nobel Prize in Physics. The star is pulled back and forth as it is orbited by planets which are usually too close to the star to be resolved in a telescope image. Spectra to obtain the star's radial velocity are repeatedly obtained to detect the resulting periodic Doppler redshift and blueshift. Extremely precise spectrographs are needed, as the velocity shifts may be less than 1 meter/sec, a tiny fraction of the velocity of light at 300 million meters/sec.

The resulting dataset is a time series of radial velocity measurements. Unfortunately, the cadence of observing times is sparse (as the spectra require considerable allocation of telescope time) and irregularly spaced ([§5.1](#)). Furthermore, the observed Doppler shifts are subject to both instrumental noise and a ‘jitter’ intrinsic to the star caused by magnetic activity similar to that seen in the Sun ([§3.1](#)).

Fortunately, we have a reliable astrophysical model for the orbits of planets around a star, essentially the same that Newton's Laws gave for celestial mechanics in our own Solar System two centuries ago. The parameters for the orbit of each planet include the period (in days), radial velocity amplitude (in meters/sec), eccentricity, four angles, a velocity calibration, and at least one parameter for the stellar jitter. With nine parameters for each planet, and a crucial additional parameter on the number of planets in the system, modeling a radial velocity time series involves optimizing a non-linear high dimensional problem.

This problem is typically formulated in a Bayesian inferential framework ([Ford 2005](#)). Choices of Bayesian priors and computational method have been carefully considered in light of astrophysical constraints ([Clyde et al. 2007](#), [Sharma 2017](#)). With clever choice of variable transformations, MCMC method, and implementation on GPU or cloud supercomputers, the calculation can be successful. An illustration of this approach, using a combination of radial velocity and transit measurements, is the discovery of a third planet orbiting Kepler-47 that gives new insights into planet formation processes and planet habitability around binary stars ([Orosz et al. 2019](#)).

4. LIKELIHOOD-FREE MODELING

Two main forms of statistical models can be distinguished: those describe by probability distributions for which an explicit likelihood can be written; and implicit models, as with cosmological simulations (Figure 2) from which one can simulate samples but without ex-

plicit likelihood formulations. The latter are often called generative models. In astronomy, an important example of a simple linear parametric models is the $M_{BH}-\sigma_{gal}$ relation that describes the empirical relation between the mass of the supermassive black hole present in the center of a galaxy and the velocity dispersion of the stars in its bulge (Gebhardt et al. 2000). This is a heuristic model that does not (yet) derive from astrophysical insights.

Approximate Bayesian Computation (ABC) enables parameter inference for complex physical systems in cases where the true likelihood function is unknown, unavailable, or computationally too expensive. It relies on the forward simulation of mock data and comparison between observed and synthetic results. The underlying idea is to validate the theoretical model by confronting its simulated outcomes with those observed in the real world. A model which produces sufficiently realistic synthetic results is more likely to be close to the truth than another one which does not, as with measured and simulated large-scale structure properties (§2.1). This concept is natural to astronomers; strategies very similar to – and sometimes inspired by – ABC were developed in the astronomical literature (Lin & Kilbinger 2015, Killelsey et al. 2017).

Shortly after Schafer & Freeman (2012) highlighted the potential of ABC in astronomy, Cameron & Pettitt (2012) used it to investigate the morphological evolution of distant galaxies. Since then the subject has attracted increasingly more attention including cosmological applications based on supernovae (Weyant et al. 2013), galaxy clusters number counts (Ishida et al. 2015), simulation images calibration (Akeret et al. 2015) and weak lensing peak counts (Lin & Kilbinger 2015). ABC methods have also tackled subjects like galaxy star formation histories (Hahn et al. 2017), inter-galactic medium (Davies et al. 2018), exoplanetary systems (Hsu et al. 2018), accretion disks around supermassive black holes (Witzel et al. 2018), luminosity functions (Riechers et al. 2019), and the stellar IMF (Cisewski-Kehe et al. 2019). This movement is accompanied by the development of specialized software, designed by astronomers, which have played an important role in the dissemination of the technique such as CosmoABC (Ishida et al. 2015), abcpmc (Akeret et al. 2015), and astroABC (Jennings & Madigan 2017).

As the astronomical community became more familiar with these statistical tools, it also became aware of its bottlenecks – and started to propose solutions for crucial aspects the paradigm. For example, the high number of simulations required by ABC can be prohibitive for many astronomical cases where calculations are computationally expensive. Kacprzak et al. (2018) proposes a quantile regression strategy which approximates the behavior of the distance function, thereby identifying regions of the parameter space with low probability to avoid generating samples dissimilar from the observed data. Alsing et al. (2018) propose a Density Estimation Likelihood-Free Inference algorithm where information from all available simulations are used to estimate a joint distribution of data and parameters, leading to a more efficient sampling than the traditional Population Monte Carlo approach.

They also address the choice of summary statistics, proposing a 2-step algorithm that reduces dimensionality by combining heuristic summary statistics and Fisher information (Alsing & Wandelt 2018). Charnock et al. (2018) suggest replacing the user defined summary by a neural network which is able to find nonlinear functions of the data that maximize Fisher information. Alsing et al. (2019) push the paradigm forward with neural density estimators to learn the likelihood function from a set of simulated data sets and coupling it with active learning (Settles 2012) to adaptively acquire simulations in the most relevant regions of parameter.

In less than 10 years, likelihood-free inference techniques like ABC and its derivatives

have flourished within the astronomical community. Given the high complexity of data expected from the next generation of large scale surveys, these efforts will prove valuable to the future of parametric inference in modern astronomical data.

5. CHALLENGES IN SIGNAL DETECTION

5.1. Periodicity Detection in Irregular Time Series

The study of variable objects in the sky – time domain astronomy – is burgeoning with more than 2000 studies annually (Griffin 2019). A common study involves the multi-epoch measurement of brightness of variable stars or quasars; astronomers call these time series ‘light curves’. Yet, except for Fourier analysis and occasionally wavelet analysis pioneered by Starck et al. (1997), astronomers rarely use statistical methods of time series analysis established for use in signal processing and econometrics (Box et al. 2015).

The principal reason is that most multi-epoch programs have irregularly spaced cadences. Some causes are unavoidable: for a single mountaintop telescope using visible light, a star or quasar is unobservable during daylight and is totally unobservable for half the year due to the annual solar motion. Other causes are sociological: telescope allocation committees must juggle projects by many scientists and can rarely give regularly spaced cadences. There are a few exceptions to these difficulties, such as NASA’s Kepler satellite that reported the brightness of $\sim 200,000$ stars every 30 minutes for four years, or the HAT South network of three small telescopes on different continents to reduce diurnal gaps. But irregular cadences are the norm in astronomical time series.

As outlined in §1 cosmic objects exhibit an incredible variety of temporal characteristics in all parts of the electromagnetic spectrum, some of which are strictly periodic. These might be a rapidly rotating neutron star with narrowly beamed radio emission (§5.2) or an exoplanet transiting a single star with periodic variations in velocity and brightness (§3.3). These situations are sufficiently common and important that astronomers have developed a suite of methods for characterizing both periodic and aperiodic variations in irregular time series. However, their statistical foundations are often inadequate, and the quality of scientific inferences consequently are often unreliable.

Astronomers first developed nonparametric periodograms that provide signals in flexible time domain situations: irregular cadences; non-sinusoidal shapes (e.g., for brief transits); and heteroscedastic measurement errors (§7.1). The most widely used for variable star characterization is Stellingwerf (1978) ‘phase dispersion minimization’ (PDM) periodogram. Here the observations are folded modulo a trial period P , the data are grouped into a small number k of evenly spaced bins, the weighted mean x_i and standard deviation σ_i for $i = 1, \dots, k$ bins are obtained, and the PDM(P) is the ratio of within-bin to between-bin dispersions weighted by measurement errors. The periodogram consists of PDM values for a large number of trial periods. If periodic behavior is present then, at the correct period P_0 , the high (low) brightnesses are collected in one or a few bins, and the PDM periodogram has a dip at P_0 . If there is no periodicity, then the high (low) brightnesses are distributed randomly among the bins and the PDM shows only noise values. Another formulation called the Analysis of Variance periodogram is mathematically related to the PDM (Schwarzenberg-Czerny 1996).

While PDM(P) has a definable distribution with large sample Gaussian white noise, the conditions are often unfavorable. Too few data points may be present in some bins for accurate variance estimates; indeed, some bins can be empty for some trial periods. Scatter

within a bin is often non-Gaussian with outliers from poorly quantified measurements. The choice of k is arbitrary, the probabilities for False Alarms in the multiple trials has not been addressed, and the periodogram is subject to strong aliasing effects when true periods are present. Perhaps the most insidious problem is that many stars show autocorrelated but aperiodic variations that interact with the irregular cadence to produce false peaks in the periodogram for short duration datasets. None of these statistical issues have been examined by mathematical statisticians despite its use in ~ 1500 astronomical papers over four decades.

Variants of epoch-folding periodograms have been developed. Dworetzky (1983) proposes an unbinned L_1 version of the PDM nicknamed the ‘minimum string length’ periodogram that alleviates some problems, but in practice seems less sensitive than PDM. The detection of exoplanet transits, where the shape of the dip in brightness is box-shaped as the planet passes in front of the star, is commonly based on Box Least Squares regression procedure (Kovács et al. 2002). Caceres et al. (2019) recently developed an efficient algorithm, nicknamed Transit Comb Filter, designed for a differenced transit light curve that calculates a matched filter to a periodic double-spike pattern rather than a periodic box pattern.

The most commonly used tool for period searching in irregular cadence astronomical light curves is the Lomb-Scargle periodogram (LSP, Scargle 1982 with 4000 citations) that assumes sinusoidal periodic behavior. It is a generalization of the Schuster periodogram in Fourier analysis for irregularly cadences. It often shows higher signal-to-noise and weaker aliases than PDM in tests with simulated data.

But debates have waged on evaluating the statistical significance, or False Alarm Probability (FAP), of LSP peaks for realistic data. The analytic exponential distribution of peak power applies only to idealized data: an infinite stream of regularly spaced Gaussian white noise with known variance and a single sinusoid superposed (Percival & Walden 1993). Koen (1990) argues that the substitute of the sample variance for the population variance in this exponential formula can lead to badly underestimated errors in FAPs. Several groups develop links between the LSP and Bayesian models, including an odds ratio for LSP peak significance (Brethorst 2003, Mortier et al. 2015). LSP FAPs based on generalized extreme value (GEV) distributions have a strong performance in the analyses of Baluev (2008) and Süveges et al. (2015). Sulis et al. (2017) combine bootstrap resampling with GEV distributions to estimate FAPs, and Delisle et al. (2020) develop it further with computationally efficient approximations that work in the presence of correlated noise. GEV-based approaches have garnered some support in the community, but many practitioners are still using naive, unreliable FAPs. In his review on understanding LSPs, VanderPlas (2018) concludes:

Unfortunately, there is no silver bullet for answering these broader, more relevant questions of uncertainty of Lomb-Scargle results. Perhaps the most fruitful path toward understanding of such effects for a particular set of observations with particular noise characteristics and a particular observing window is via simulated data injected into the detection pipeline.

5.2. Gravitational Wave Detection

The dawn of a new era of astronomy came in 2015 with the detection of Gravitational Waves (GWs). A century earlier, Albert Einstein predicted that changes in gravitational

fields would cause tiny ripples in space-time to propagate through the Universe. The effect is so weak that only the most sophisticated instruments can detect the strongest gravitational events, such as the inspiralling and merger of two black holes or neutron stars (Figure 5). These remarkable events are rare, but a growing number have now been detected in other galaxies with the Laser Interferometer Gravitational-Wave Observatory (LIGO). Other observatories are seeking GWs, such as ground-based radio pulsar timing arrays and planned space-based multi-satellite interferometers. This new field of astronomy, which garnered the 2017 Nobel Prize in Physics, will open new opportunities for confronting both astrophysical theories (binary star evolution, black hole formation, merging galaxies, quasar evolution) and fundamental physics (high density matter, particle physics, theory of gravitation) in ways that are otherwise inaccessible.

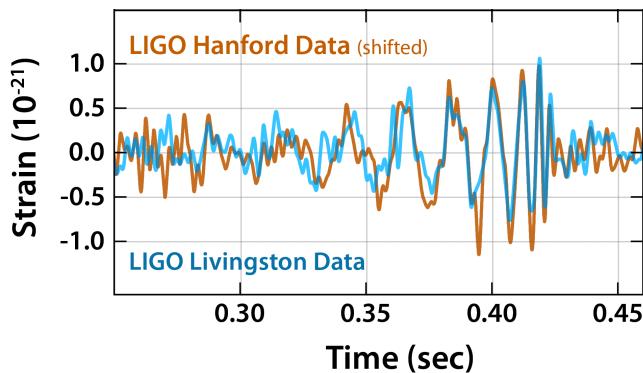


Figure 5

Discovery of the first gravitational wave event, GW15091, from the U.S. Laser Interferometric Gravitational Observatory. The plot shows a chirp signal from two widely separated interferometers after shifting by a short time delay. This event arises from the inspiraling and coalescence of two ~ 30 solar mass black holes in an unknown distant galaxy. Credit: Caltech/MIT/LIGO Lab.

The LIGO observatories (now in the U.S. and Italy, soon to expand in Japan and India), have laser beams of infrared light traveling along two perpendicular kilometer-scale vacuum chambers where mirrors and beam splitters allow tiny motions in the length of the device to be measured using interference patterns. LIGO is designed to detect GW events on timescales of $10^{-4} - 10^{-1}$ seconds corresponding to violent changes in stellar mass bodies such as white dwarfs, neutron stars and black holes arising from supernova explosions. In some cases, the GW event is accompanied by a burst of gamma-rays or other electromagnetic radiation; this was first seen in GW170817 (Abbott et al. 2017). Time delays between detections at the different observatories allow triangulation of the GW location in the celestial sphere, so other observatories can quickly search for electromagnetic counterparts. The scramble to detect signals from different observatories associated with GW and similar explosive events is nicknamed ‘multi-messenger astronomy’.

The statistical challenge with LIGO is to detect sudden short-lived chirp-like events in a continuous time series (Figure 5) where noise is dominated by instrumental effects that can be continuous (perhaps caused by vibrations in the mirror structures) or transient (perhaps caused by minor Earth tremors). We seek brief, weak signals in non-stationary, non-Gaussian noise. An elaborate procedure for removing or flagging these non-astrophysical

variations has been developed (McIver [2012] Cabero et al. [2019]). The resulting residual time series is nearly Gaussian white noise, and GW chirps are most often sought by matched filtering with astrophysical models. Unsupervised classification of GW candidates with convolutional neural networks have also been applied (George & Huerta [2018]).

Gravitational waves are predicted to span a wide range of frequencies requiring different technologies to detect. Low-frequency gravitational waves are emitted from pairs of supermassive black holes in distant galaxies in tight orbits that power quasars and other active galactic nuclei, each of which is millions times more massive than those detected by LIGO. It is hoped these can be detected with a clever method using millisecond pulsars. Radio telescopes have been observing Galactic millisecond pulsars, produced by spun-up rapidly rotating neutron stars, for decades but only recently have been harnessed for low-frequency GW detection (Taylor et al. [2016] Mingarelli et al. [2017]). Time series of pulsar pulse arrival times from pulsars distributed through the Milky Way Galaxy provide a Galaxy-scale detector for nanohertz GWs. The astrophysical signals here are not short chirps but rather continuous periodic GWs emitted by pairs of orbiting supermassive black holes in distant galaxies. A stochastic autocorrelated GW noise from many orbiting black holes is expected.

Careful tracking and modeling of neutron star pulse arrival time give a very precise measurement of distance from Earth to the pulsar. When a set of pulsars around the Galaxy are monitored in a pulsar timing array, it provides the ability to detect minute variations in space-time due to the passage of long wavelength GWs. Correlating the residuals from model predictions across pairs of pulsars leverages the common influence of a gravitational-wave background against unwanted, uncorrelated noise. So far, no astrophysical signal has been reported (NANOGrav Collaboration et al. [2018]).

The NAGOGgrav time series data present considerable scientific, statistical and computational challenges. This analysis starts with multicomponent, nonlinear model of a non-Gaussian and non-stationary time series (NANOGrav Collaboration et al. [2015] [2018]). A model of pulse arrival times must take into account: the astronomical characteristics of the pulsar (spin period, proper motion); gravitational effects from any stellar companions; time-varying dispersion measure variations from the propagation of the pulse through the inhomogeneous Galactic interstellar medium; autoregressive ‘red’ noise of unknown origin with timescales of weeks to years; white noise (radio receiver noise and some additional component of unknown origin). And finally, the model seeks a non-zero amplitude for a GW signal which gives a unique correlation pattern across multiple pulsars. The GW component of the model typically rests on assumptions of an isotropic distribution of black holes and circular orbits of black hole binaries. Best-fit models are typically obtained using Bayesian inference with uninformative priors.

The models here, and in more broadly in astrostatistics, traditionally assume the non-GW behaviors are stationary and Gaussian. But newer precise data are beginning to show traces of non-Gaussian and non-stationary features. Common analysis has not yet benefited from statistical methods such as multi-level hierarchical modeling, Gibbs sampling, or nonlinear time series modeling. Additionally, testing for GW-induced distortions in this complicated model is often addressed with computationally expensive techniques involving Gaussian Processes regression and Bayesian inference with MCMC sampling (van Haasteren & Vallisneri [2014]). More efficient approaches are needed.

6. TWO EXAMPLES OF MACHINE LEARNING IN ASTRONOMY

Machine learning techniques are growing exponentially in the astronomical literature covering a wide range of questions in planetary, stellar, extragalactic, time domain and cosmological fields. A new generation of young astronomers are quickly becoming proficient in their use.

6.1. Photometry of Blended Galaxies

The upcoming generation of deep and wide galaxy surveys from ground (LSST) and space (Euclid,WFIRST) observatories will present image processing challenges for which new methodologies are imperative (Schmitz et al. 2018). Of particular interest is the capability to extract photometric information overlapping, or blended, objects. As the instruments become more sensitive, the blending probability increases; half of the sources observed by LSST are predicted to have some level of overlap. This blending can involve two galaxies that are physically interacting, two unrelated galaxies that overlap by chance, or chance superposition of a Galactic star with a distant galaxy. While astronomically motivated methods image segmentation algorithms are widely used (Bertin & Arnouts 1996) (Mancone et al. 2013), interest has increased in data-driven approaches to address this problem (Melchior et al. 2018).

Most recently, deep learning has been used to process galaxy images, particularly for the classification of galaxy morphologies (Dieleman et al. 2015) (Barchi et al. 2017) (Domínguez Sánchez et al. 2018) (Huertas-Company et al. 2018) (Khalifa et al. 2018). Techniques seek the recovery of galaxy features in noisy images with generative adversarial networks (Schawinski et al. 2017), the search for strong lensing effects with deep learning networks (Lanusse et al. 2018), and for deblending galaxies (Figure 6) (Reiman & Göhre 2019). Boucaud et al. (2020) implement a modular version of U-net architecture (Ronneberger et al. 2015) to recover the fractional segmentation map; this is an image with pixel values between 0 and 1 estimating the fraction of flux belonging to a given galaxy. Figure 7 displays the input blended galaxies on the top left and the recovered masks on the top right. The method

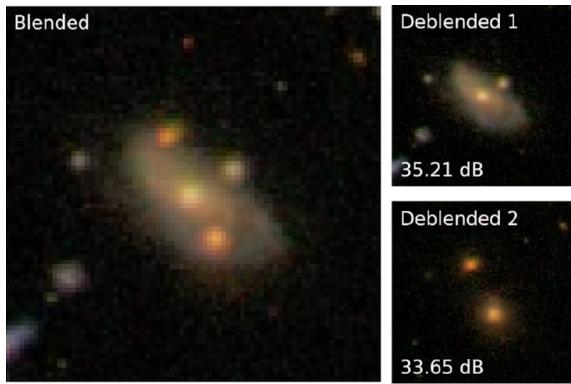


Figure 6

An artificial image constructed of two blended galaxy groups from the SDSS survey (left) and a successful decomposition using a branched generative adversarial network (Reiman & Göhre 2019).

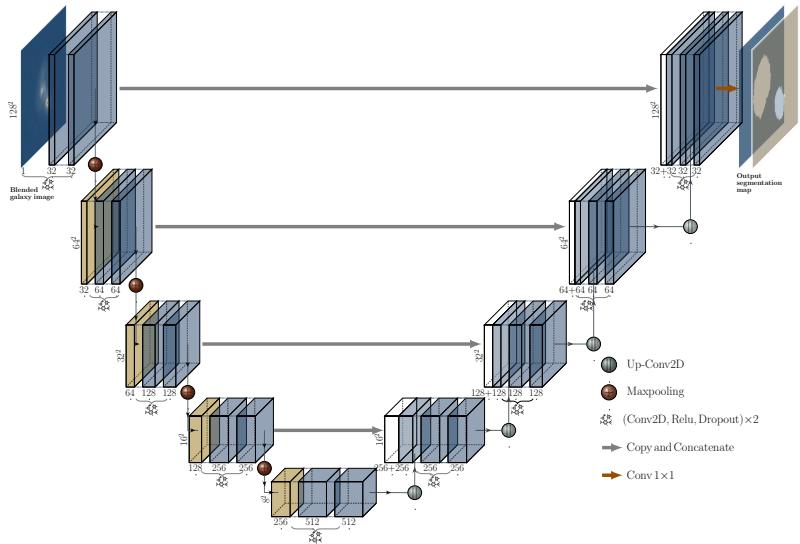


Figure 7

A U-net network that takes as input an image of blended system and outputs two fractional segmentation maps, each pixel corresponding to the fraction of flux belonging to a given galaxy (Boucaud et al. 2020).

outperformed more traditional approaches, stressing the power of deep learning architecture tailored for astronomical image processing.

6.2. The Accelerating Expansion of the Universe

The expansion of the Universe since the Big Bang 13.7 billion years ago is naturally decelerated by the gravitational attraction of Dark Matter and ordinary matter. So it was a major surprise two decades ago when distance measurements of Type Ia supernovae (SNe Ia) provided the first empirical evidence that our Universe is experiencing an accelerated expansion. SNe Ia are particularly useful because they behave as standard candles; that is, their intrinsic luminosities can be inferred enabling accurate estimates of their distances for a given cosmological model. Whilst the universal acceleration first detected by SNe Ia is corroborated by independent observations of the cosmic background radiation and other lines of evidence, SNe Ia continue to play a major role in our quest to understand Dark Energy, the unknown phenomena that causes this accelerated expansion. Dedicated surveys are underway to improve the statistics of Dark Energy measurements; LSST is expected to measure $\sim 300,000$ SNe Ia light curves over the period of 10 years. However, other types of SNe and transient objects confuse the situation. Only a few percent of LSST SNe are expected to have spectroscopic information which provide accurate labels. So the community must rely on sparse photometric light curves to classify transients and thus obtain SNe Ia samples.

Ideally, the SNe photometric classification problem could be formulated as a case of straight-forward supervised learning. A subset of SNe with both spectroscopic and photo-

metric information can give training sets for different SNe types. This formulation of the problem was recently submitted to the scrutiny of the methodology community through the Kaggle PLAsTiCC competition (The PLAsTiCC team et al. 2018), which attracted more than 1000 participating teams. Nevertheless, a realistic treatment is more complex for several reasons.

A particularly difficult characteristic of the SN classification problem is the discrepancy between spectroscopic and photometric samples, violating a central assumption underlying most learning algorithms. Due to the heavy observational burden, spectroscopy will not be available for the faintest LSST SNe, leading to biased training sets (Childress et al. 2017). Several efforts are underway to treat this issue in machine learning-based analyses (Revsbech et al. 2018), including the use of deep learning.

This leads to the problem of optimizing the distribution of spectroscopic resources to construct a training sample that maximizes accurate classification with a minimum number of labels. Within the framework of active learning, Ishida et al. (2019) provide observers with a spectroscopic follow-up protocol on a night-by-night basis. The method iteratively identifies which objects in the target (photometric) sample would most likely improve the classifier if included in the training data - allowing sequential updates of the learning model with a minimum number of labeled instances. While the approach should outperform standard supervised learning models, caveats still remains. Better modeling of uncertainty quantification of the photometric measurement errors, awareness of anomaly objects, and a seamless inclusion of astronomical information into the loss function are needed.

7. TWO CLASSIC PROBLEMS IN ASTROSTATISTICS

While major projects are underway to advance machine learning methods for Big Data problems in astronomy, many research efforts can benefit from innovations in more classical statistical inference methodology.

7.1. Heteroscedastic Measurement Errors

Heteroscedasticity occurs when errors are not uniform across a sample but exhibit dependencies on the measurement process or some (often unknown) properties of the objects under study. In regression, the heteroscedasticity may be included within hierarchical statistical models as nuisance hyperparameters (Carroll et al. 2006). Heteroscedastic measurement errors are ubiquitous in astronomy but with a crucial difference: astronomers measure the errors directly, so the variance of the error enters into the input data along with the property under study. This is possible because astronomers build carefully calibrated instruments with known noise characteristics, and additional noise from the celestial environment or observational conditions can be measured simultaneously with the property of interest. In imagery, it is the noise of the dark sky next to a star or galaxy; in spectroscopy, the noise of the spectrum next to an emission line of interest; and in time series, the temporal behavior before and after an event of interest. It is thus standard procedure throughout astronomy for each property measured in each object to be accompanied by its own heteroscedastic measurement error with known variance.

The only common use of these measurement errors is in regression with weighted least squares which minimizes the sum $\sum_{i=1}^n (x_i - M(\theta)_i)^2 / \sigma_i^2$. Here x_i is the observed property of interest for the i -th measurement, M_i is the model prediction with parameters θ , and σ_i^2

is the observed variance of the measurement. Astronomers call this ‘minimum- χ^2 ’ regression hoping that the sum is chi-squared distributed. The problem is that other components to the total error are often present which may be added in quadrature to the denominator of the sum. The model is often misspecified in realistic situations so that the total error does not account for the sample variance and the sum is no longer chi-squared distributed.

There is great need for statistical methodology treating known heteroscedastic measurement errors in astronomy beyond this regression problem. First, in many cases the value of x_i is not much greater than σ_i , and the object is considered to be undetected and the measurement x_i is replaced by a left-censored value like $< 3 \times \sigma_i$. The problem is that standard survival analysis treating censoring does not also treat heteroscedastic weighting of the clearly detected data points. Second, astronomers have other statistical needs such as cluster and classification which need to incorporate measurement errors. Although this may be possible for parametric treatment like Gaussian mixture models, heteroscedastic errors are not easily adapted to nonparametric algorithmic clustering and classification procedures.

One approach to the problem is a partially Bayesian solution called the iterative conditional modes which uses componentwise maximization routine to find the mode of the posterior (Berry et al. 2002). Another promising direction in the machine learning community is the development of a framework for Bayesian deep learning (Gal & Ghahramani 2016). Similarly, the Probably Approximately Correct Bayesian (PAC-Bayesian) Learning formalism has been used to give theoretical guarantees to cutting edge topics in machine learning as deep learning and domain adaptation (Guedj 2019).

7.2. Domain-Specific Probability Distributions

Astronomers often seek analytic formulae to model the distributions of properties of scientific interest. One of the most famous is the ‘Schechter function that accurately describes the distribution of galaxy luminosities with a truncated power law (Pareto) at low luminosities and an exponential distribution at high luminosities. But astronomers did not realize this is just Eulers gamma distribution. A similar issue is the stellar IMF which resembles a power law at high masses (known as the 1955 ‘Salpeter function) and a lognormal at low masses. This is usually modeled as a piece-wise composite function such as Kroupas or Chabriers IMF. But Maschberger (2013) proposes a continuous generalized log-logistic distribution with three parameters,

$$P_{L3}(m) \propto \frac{m^{-\alpha}}{\mu} \left(1 + \left(\frac{m}{\mu} \right)^{1-\alpha} \right)^{-\beta}. \quad 2.$$

In a different area of astrophysical study, the distribution of particle energies in plasmas emerging from the Sun and measured throughout the Solar System is found to follow a ‘kappa distribution that is Maxwellian (Gaussian) at low energies and power law at high energies

$$f^\kappa(v) = \frac{1}{(\pi \kappa v_\kappa^2)^{3/2}} \frac{\Gamma(\kappa + 1)}{\Gamma(\kappa - 1/2)} \left(1 + \frac{v^2}{\kappa v_\kappa^2} \right)^{-(\kappa+1)}. \quad 3.$$

Astrophysical insights can emerge from linking such distribution functions to generative stochastic processes. For example, statisticians Reed & Jorgensen (2004) show how a broad class of double-Pareto functions can be produced by killed multiplicative Brownian motion processes. Astrophysicist Collier (1993) similarly shows that kappa distributions can arise from random walks in velocity governed by Lvy flight probability distributions.

8. ASTROSTATISTICAL CHALLENGES FOR STATISTICIANS

Contemporary astronomical data analysis often elude the capabilities of classical statistical techniques, and inevitably requires the use and development of sophisticated, and sometimes novel, statistical tools.

Astronomy requires expertise in vast fields of statistics and information science: non-parametric and parametric inference (especially Bayesian), high-dimensional nonlinear regression, censoring and truncation, measurement error theory, spatial point processes, image analysis, time series analysis, multivariate analysis, clustering and classification, and many other forms of machine learning. Statistical models range from simple heuristic power law regressions to high-dimensional non-linear models from astrophysical theory. Samples sizes range from a dozen to billions of objects. A hierarchy of problems with several layers of uncertainty is often involved; for example, a survey subject to flux limits is then filtered by morphology, subject to classification, and multivariate relationships are sought involving properties that may or may not be subject to the original truncation. When errors are non-Gaussian, analysis can benefit from generalized linear modeling (de Souza et al. 2015a[b]).

Clearly, cross-disciplinary collaboration and research is not only desired but imperative. Yet most of the astrostatistical innovations reviewed here have been developed by astronomers who have little formal training in the mathematical and computational sciences. Complex astrostatistical procedures are often developed in isolation from the mainstream of methodological studies. Education of astronomers in methodology is generally informal with emphasis on Python-based software and hack weeks rather than thorough university-based courses. Research in astronomy can be well-funded, but resources for methodological development are scarce; astrostatistical efforts are typically informally embedded within formal science and software projects. An institute or observatory will employ dozens of researchers with university degrees in astronomy and physics but very few with degrees in statistics, applied mathematics or computer science.

Despite these structural constraints, a vibrant field of astrostatistics has emerged since the 1990s with international conferences, training workshops, and several cross-disciplinary scholarly organizations. An informal online Facebook group on astrostatistics started in 2013 has grown to nearly 5000 members. The Cosmostatistics Initiative founded in 2014 created a successful interdisciplinary science development environment where innovative astrostatistics projects are developed and disseminated. The Statistical and Applied Mathematical Sciences Institute has run several astrostatistical programs.

In some respects, statisticians can readily enjoy the fruits of astronomical observations for statistical study: a vast range of data are freely available. The astronomical community has a long-standing tradition, often legally binding, of making both raw and analyzed data accessible on the Web. The U.S. National Aeronautical and Space Administration and European Space Agency operate large science archive centers for satellite observatories. Ground-based data are available from the European Southern Observatory, U.S. National Optical Astronomy Observatory and National Radio Astronomy Observatory, and other major institutions. The International Virtual Observatory Alliance provides an integrated interface to these and hundreds of other distributed databases. The Vizier and SIMBAD services from Frances Centre de Donnes Astronomique, and the NASA/IPAC Extragalactic Database, provide convenient Web-based access to published tabular data. The Smithsonian/NASA Astrophysics Data System provides superb bibliographic services with full-text, references and citations for nearly the entire astronomical research literature.

But in other respects, conducting astrostatistical research requires considerable care. The datasets are often subject to selection biases inaccessible to non-experts. The scientific questions are often complex, so analyses must be designed within the context of established research programs. Statistical studies in cosmology can be particularly challenging. Analyses of cosmic microwave background maps need to be interpreted within the Λ CDM cosmological model, and three-dimensional cosmography of Dark Matter combines complexities of both statistical weak lensing signatures and the imprecision of photo-z distances.

Interested statisticians can meet astronomers through newly formed scholarly societies and at specialized conferences. The methodology community has the International Statistical Institutes Astrostatistics Special Interest Group, American Statistical Associations Astrostatistics Interest Group, IEEE Task Force on AstroData Mining, and the independent International Astrostatistics Association and International AstroInformatics Association. The astronomical community has the International Astronomical Unions Commission on Astroinformatics and Astrostatistics, and the American Astronomical Society's Working Group on Astroinformatics and Astrostatistics. Ongoing conference series include *Statistical Challenges in Modern Astronomy*, *Astroinformatics*, and *Astronomical Data Analysis* as well as sessions at Joint Statistical Meetings, World Statistics Congresses, and American Astronomical Society meetings. The large LSST project has an Information and Statistics Science Collaboration.

The purpose of astronomy is to discover, characterize, and gain physical insight into cosmic phenomena. Tremendous successes have emerged in the past century, propelled by amazingly sophisticated and sensitive instrumentation. Astrostatistics plays an important integrative role, bridging raw data to intelligible information, and information to insightful astrophysical models. Astrostatistics is a young, fertile field in which interdisciplinary communities can grow, providing essential expertise to further our understanding of the Universe we inhabit.

Acknowledgements: Astrostatistics at Penn State is supported by NSF grant AST-1614690, NASA grant 80NSSC17K0122, and the Eberly College of Science through the Center for Astrostatistics. EEOI is supported by a 2018-20 CNRS MOMENTUM fellowship.

LITERATURE CITED

- Abbott BP, Abbott R, 3676 colleagues. 2017. Multi-messenger observations of a binary neutron star merger. *Astrophys. J.l* 848:L12
- Akeret J, Refregier A, Amara A, Seehars S, Hasner C. 2015. Approximate Bayesian computation for forward modeling in cosmology. *J. Cosmology & Astroparticle Physics* 2015:043
- Alsing J, Charnock T, Feeney S, Wandelt B. 2019. Fast likelihood-free cosmology with neural density estimators and active learning. *Mon. Not. Royal Astro. Soc.* 488:4440–4458
- Alsing J, Wandelt B. 2018. Generalized massive optimal data compression. *Mon. Not. Royal Astro. Soc.* 476:L60–L64
- Alsing J, Wandelt B, Feeney S. 2018. Massive optimal data compression and density estimation for scalable, likelihood-free inference in cosmology. *Mon. Not. Royal Astro. Soc.* 477:2874–2885
- Asensio Ramos A. 2006. The Minimum Description Length Principle and Model Selection in Spectropolarimetry. *Astrophys. J.* 646:1445–1451
- Asensio Ramos A. 2009. Evidence for Quasi-Isotropic Magnetic Fields from Hinode Quiet-Sun Observations. *Astrophys. J.* 701:1032–1043
- Baddeley A, Rubak E, Turner R. 2015. Spatial point patterns: Methodology and applications with r. Chapman & Hall/CRC Interdisciplinary Statistics. CRC Press

- Baluev RV. 2008. Assessing the statistical significance of periodogram peaks. *Mon. Not. Royal Astro. Soc.* 385:1279–1285
- Barchi PH, da Costa FG, Sautter R, Moura TC, Stalder DH, et al. 2017. Improving galaxy morphology with machine learning. *J. Comput. Interdiscip. Sci.* 7:114
- Barrow JD, Bhavsar SP, Sonoda DH. 1985. Minimal spanning trees, filaments and galaxy clustering. *Monthly Notices of the Royal Astronomical Society* 216:17–35
- Beattie JR, Federrath C. 2019. Filaments and striations: anisotropies in observed, supersonic, highly magnetized turbulent clouds. *Monthly Notices of the Royal Astronomical Society* 492:668–685
- Beck R, Dobos L, Budavári T, Szalay AS, Csabai I. 2016. Photometric redshifts for the SDSS Data Release 12. *Mon. Not. Royal Astro. Soc.* 460:1371–1381
- Beck R, Lin CA, Ishida EEO, Gieseke F, de Souza RS, et al. 2017. On the realistic validation of photometric redshifts. *Mon. Not. Royal Astro. Soc.* 468:4323–4339
- Benítez N. 2000. Bayesian Photometric Redshift Estimation. *Astrophys. J.* 536:571–583
- Berry SM, Carroll RJ, Ruppert D. 2002. Bayesian smoothing and regression splines for measurement error problems. *Journal of the American Statistical Association* 97:160–169
- Bertin E, Arnouts S. 1996. SExtractor: Software for source extraction. *Astron. & Astrophys.* 117:393–404
- Bond NA, Strauss MA, Cen R. 2010. Crawling the cosmic network: identifying and quantifying filamentary structure. *Mon. Not. Royal Astro. Soc.* 409:156–168
- Bossini D, Vallenari A, Bragaglia A, Cantat-Gaudin T, Sordo R, et al. 2019. Age determination for 269 Gaia DR2 open clusters. *Astron. & Astrophys.* 623:A108
- Boucaud A, Huertas-Company M, Heneka C, Ishida EEO, Sedaghat N, et al. 2020. Photometry of high-redshift blended galaxies using deep learning. *Mon. Not. Royal Astro. Soc.* 491:2481–2495
- Brethorst GL. 2003. Frequency estimation and generalized Lomb-Scargle periodograms, In *Statistical Challenges in Astronomy*, eds. ED Feigelson, GJ Babu. New York, NY: Springer New York
- Budavári T. 2009. A Unified Framework for Photometric Redshifts. *Astrophys. J.* 695:747–754
- Cabero M, Lundgren A, Nitz AH, Dent T, Barker D, et al. 2019. Blip glitches in Advanced LIGO data. *Classical and Quantum Gravity* 36:155010
- Caceres GA, Feigelson ED, Jogesh Babu G, Bahamonde N, Christen A, et al. 2019. Autoregressive Planet Search: Methodology. *Astron. J.* 158:57
- Cameron E, Pettitt AN. 2012. Approximate Bayesian Computation for astronomical model analysis: a case study in galaxy demographics and morphological transformation at high redshift. *Mon. Not. Royal Astro. Soc.* 425:44–65
- Cantat-Gaudin T, Krone-Martins A, Sedaghat N, Farahi A, de Souza RS, et al. 2019. Gaia DR2 unravels incompleteness of nearby cluster population: new open clusters in the direction of Perseus. *Astron. & Astrophys.* 624:A126
- Carliles S, Budavári T, Heinis S, Priebe C, Szalay AS. 2010. Random Forests for Photometric Redshifts. *Astrophys. J.* 712:511–515
- Carroll R, Ruppert D, Stefanski L, Crainiceanu C. 2006. Measurement error in nonlinear models: A modern perspective, second edition. CRC Press
- Carvuto S, Brescia M, Vellucci C, Longo G, Amaro V, Tortora C. 2016. Probability density estimation of photometric redshifts based on machine learning, In *2016 IEEE Symposium Series on Computational Intelligence (SSCI)*, pp. 1–6
- Charnock T, Lavaux G, Wandelt BD. 2018. Automatic physical inference with information maximizing neural networks. *Physical Review D* 97:083004
- Chen YC, Ho S, Freeman PE, Genovese CR, Wasserman L. 2015. Cosmic web reconstruction through density ridges: method and algorithm. *Mon. Not. Royal Astro. Soc.* 454:1140–1156
- Childress MJ, et al. 2017. OzDES multifibre spectroscopy for the Dark Energy Survey: Three year results and first data release. *MNRAS* 472:273–288
- Cisewski-Kehe J, Weller G, Schafer C. 2019. A preferential attachment model for the stellar initial

- mass function. *Electron. J. Statist.* 13:1580–1607
- Clyde MA, Berger JO, Bullard F, Ford EB, Jefferys WH, et al. 2007. Current Challenges in Bayesian Model Choice, In *Statistical Challenges in Modern Astronomy IV*, eds. GJ Babu, ED Feigelson, vol. 371 of *Astronomical Society of the Pacific Conference Series*, p. 224
- Collier MR. 1993. On generating Kappa-like distribution functions using velocity space Lévy flights. *Geophysical Research Letters* 20:1531–1534
- Dahlen T, Mobasher B, Faber SM, Ferguson HC, Barro G, et al. 2013. A Critical Assessment of Photometric Redshift Methods: A CANDELS Investigation. *Astrophys. J.* 775:93
- Davies FB, Hennawi JF, Eilers AC, Lukić Z. 2018. A New Method to Measure the Post-reionization Ionizing Background from the Joint Distribution of Ly α and Ly β Forest Transmission. *Astrophys. J.* 855:106
- de Lapparent V, Geller MJ, Huchra JP. 1986. A Slice of the Universe. *Astrophys. J.l* 302:L1
- de Souza RS, Cameron E, Killedar M, Hilbe J, Vilalta R, et al. 2015a. The overlooked potential of Generalized Linear Models in astronomy, I: Binomial regression. *Astronomy and Computing* 12:21–32
- de Souza RS, Hilbe JM, Buelens B, Riggs JD, Cameron E, et al. 2015b. The overlooked potential of generalized linear models in astronomy - III. Bayesian negative binomial regression and globular cluster populations. *MNRAS* 453:1928–1940
- Delisle JB, Hara N, Ségransan D. 2020. Efficient modeling of correlated noise. I. Statistical significance of periodogram peaks. *Astron. & Astrophys.* 635:A83
- Dieleman S, Willett KW, Dambre J. 2015. Rotation-invariant convolutional neural networks for galaxy morphology prediction. *Mon. Not. Royal Astro. Soc.* 450:1441–1459
- D’Isanto A, Polsterer KL. 2018. Photometric redshift estimation via deep learning. Generalized and pre-classification-less, image based, fully probabilistic redshifts. *Astron. & Astrophys.* 609:A111
- Domínguez Sánchez H, Huertas-Company M, Bernardi M, Tuccillo D, Fischer JL. 2018. Improving galaxy morphologies for SDSS with Deep Learning. *Mon. Not. Royal Astro. Soc.* 476:3661–3676
- Dworetzky MM. 1983. A period-finding method for sparse randomly spaced observations or “How long is a piece of string ?”. *Mon. Not. Royal Astro. Soc.* 203:917–924
- Eisenstein DJ, Zehavi I, Hogg DW, Scoccimarro R, Blanton MR, et al. 2005. Detection of the Baryon Acoustic Peak in the Large-Scale Correlation Function of SDSS Luminous Red Galaxies. *Astrophys. J.* 633:560–574
- Elliott J, de Souza RS, Krone-Martins A, Cameron E, Ishida EEO, et al. 2015. The overlooked potential of Generalized Linear Models in astronomy-II: Gamma regression and photometric redshifts. *Astronomy and Computing* 10:61–72
- Ford EB. 2005. Quantifying the Uncertainty in the Orbits of Extrasolar Planets. *Astron. J.* 129:1706–1717
- Gal Y, Ghahramani Z. 2016. Dropout as a bayesian approximation: Representing model uncertainty in deep learning, In *Proceedings of The 33rd International Conference on Machine Learning*, eds. MF Balcan, KQ Weinberger, vol. 48, pp. 1050–1059, PMLR
- Gebhardt K, Bender R, Bower G, Dressler A, Faber SM, et al. 2000. A Relationship between Nuclear Black Hole Mass and Galaxy Velocity Dispersion. *ApJ* 539:L13–L16
- George D, Huerta EA. 2018. Deep Learning for real-time gravitational wave detection and parameter estimation: Results with Advanced LIGO data. *Physics Letters B* 778:64–70
- Griffin RE, ed. 2019. Southern horizons in time-domain astronomy, vol. 339 of *IAU Symposium*
- Guedj B. 2019. A Primer on PAC-Bayesian Learning, In *Proceedings of the 2nd congress of the Société Mathématique de France*, pp. 391–414
- Hahn C, Tinker JL, Wetzel A. 2017. Star Formation Quenching Timescale of Central Galaxies in a Hierarchical Universe. *Astrophys. J.* 841:6
- Haider M, Steinhauser D, Vogelsberger M, Genel S, Springel V, et al. 2016. Large-scale mass distribution in the Illustris simulation. *Mon. Not. Royal Astro. Soc.* 457:3024–3035
- Hildebrandt H, Wolf C, Benítez N. 2008. A blind test of photometric redshifts on ground-based

- data. *Astron. & Astrophys.* 480:703–714
- Hsu DC, Ford EB, Ragozzine D, Morehead RC. 2018. Improving the Accuracy of Planet Occurrence Rates from Kepler Using Approximate Bayesian Computation. *Astron. J.* 155:205
- Hubble EP. 1930. Distribution of luminosity in elliptical nebulae. *Astrophys. J.* 71:231–276
- Huertas-Company M, Primack JR, Dekel A, Koo DC, Lapiner S, et al. 2018. Deep learning identifies high- z galaxies in a central blue nugget phase in a characteristic mass range. *Astrophys. J.* 858:114
- Ishida EEO, Beck R, González-Gaitán S, de Souza RS, Krone-Martins A, et al. 2019. Optimizing spectroscopic follow-up strategies for supernova photometric classification with active learning. *MNRAS* 483:2–18
- Ishida EEO, Vitenti SDP, Penna-Lima M, Cisewski J, de Souza RS, et al. 2015. COSMOABC: Likelihood-free inference via Population Monte Carlo Approximate Bayesian Computation. *Astronomy and Computing* 13:1–11
- Jennings E, Madigan M. 2017. astroABC : An Approximate Bayesian Computation Sequential Monte Carlo sampler for cosmological parameter estimation. *Astronomy and Computing* 19:16–22
- Kacprzak T, Herbel J, Amara A, Réfrégier A. 2018. Accelerating Approximate Bayesian Computation with Quantile Regression: application to cosmological redshift distributions. *J. Cosmology & Astroparticle Physics* 2018:042
- Khalifa NE, Hamed Tah M, , Hassani AE, Selim I. 2018. Deep Galaxy V2: Robust Deep Convolutional Neural Networks for Galaxy Morphology Classifications, In *2018 International Conference on Computing Sciences and Engineering (ICCSE)*, pp. 1–6
- Killedar M, Borgani S, Fabjan D, Dolag K, Granato G, et al. 2017. Simulation-based marginal likelihood for cluster strong lensing cosmology. *Mon. Not. Royal Astro. Soc.* 473:1736–1750
- Koen C. 1990. Significance Testing of Periodogram Ordinates. *Astrophys. J.* 348:700
- Kovács G, Zucker S, Mazeh T. 2002. A box-fitting algorithm in the search for periodic transits. *Astron. & Astrophys.* 391:369–377
- Lanusse F, Ma Q, Li N, Collett TE, Li CL, et al. 2018. CMU DeepLens: deep learning for automatic image-based galaxy-galaxy strong lens finding. *Mon. Not. Royal Astro. Soc.* 473:3895–3906
- Leistedt B, Mortlock DJ, Peiris HV. 2016. Hierarchical Bayesian inference of galaxy redshift distributions from photometric surveys. *Mon. Not. Royal Astro. Soc.* 460:4258–4267
- Lin CA, Kilbinger M. 2015. A new model to predict weak-lensing peak counts. II. Parameter constraint strategies. *Astron. & Astrophys.* 583:A70
- Lynden-Bell D, Faber SM, Burstein D, Davies RL, Dressler A, et al. 1988. Photometry and Spectroscopy of Elliptical Galaxies. V. Galaxy Streaming toward the New Supergalactic Center. *Astrophys. J.* 326:19
- Mancone CL, Gonzalez AH, Moustakas LA, Price A. 2013. PyGFit: A Tool for Extracting PSF Matched Photometry. *Publ. Astro. Soc. Pacific* 125:1514
- Maschberger T. 2013. On the function describing the stellar initial mass function. *Mon. Not. Royal Astro. Soc.* 429:1725–1733
- Mayor M, Queloz D. 1995. A Jupiter-mass companion to a solar-type star. *Nature* 378:355–359
- McIver J. 2012. Data quality studies of enhanced interferometric gravitational wave detectors. *Classical and Quantum Gravity* 29:124010
- Melchior P, Moolekamp F, Jerdee M, Armstrong R, Sun AL, et al. 2018. scarlet: Source separation in multi-band images by Constrained Matrix Factorization. *Astronomy and Computing* 24:129 – 142
- Menou K. 2019. Morpho-photometric redshifts. *Mon. Not. Royal Astro. Soc.* 489:4802–4808
- Men'shchikov A, André P, Didelon P, Könyves V, Schneider N, et al. 2010. Filamentary structures and compact objects in the Aquila and Polaris clouds observed by Herschel. *Astron. & Astrophys.* 518:L103
- Mingarelli CMF, Lazio TJW, Sesana A, Greene JE, Ellis JA, et al. 2017. The local nanohertz gravitational-wave landscape from supermassive black hole binaries. *Nature Astronomy* 1:886–

- Moews B, Schmitz MA, Lawler AJ, Zuntz J, Malz AI, et al. 2020. Ridges in the Dark Energy Survey for cosmic trough identification. *arXiv e-prints* :arXiv:2005.08583
- Mortier A, Faria JP, Correia CM, Santerne A, Santos NC. 2015. BGLS: A Bayesian formalism for the generalised Lomb-Scargle periodogram. *Astron. & Astrophys.* 573:A101
- NANOGrav Collaboration, Arzoumanian Z, Brazier A, Burke-Spolaor S, Chamberlin S, et al. 2015. The NANOGrav Nine-year Data Set: Observations, Arrival Time Measurements, and Analysis of 37 Millisecond Pulsars. *Astrophys. J.* 813:65
- NANOGrav Collaboration, Arzoumanian Zea, NANOGrav Collaboration. 2018. The NANOGrav 11-year Data Set: High-precision Timing of 45 Millisecond Pulsars. *Astrophys. J.s* 235:37
- Nelson D, Pillepich A, Genel S, Vogelsberger M, Springel V, et al. 2015. The *Illustris* simulation: Public data release. *Astronomy and Computing* 13:12–37
- Neyman J, Scott EL. 1952. A Theory of the Spatial Distribution of Galaxies. *Astrophys. J.* 116:144
- Neyman J, Scott EL. 1958. Statistical approach to problems of cosmology. *Journal of the Royal Statistical Society: Series B (Methodological)* 20:1–29
- Novikov D, Colombi S, Dor O. 2006. Skeleton as a probe of the cosmic web: the two-dimensional case. *Monthly Notices of the Royal Astronomical Society* 366:1201–1216
- Orosz JA, Welsh WF, Haghighipour N, Quarles B, Short DR, et al. 2019. Discovery of a third transiting planet in the kepler-47 circumbinary system. *The Astronomical Journal* 157:174
- Pasquet J, Bertin E, Treyer M, Arnouts S, Fouchez D. 2019. Photometric redshifts from SDSS images using a convolutional neural network. *Astron. & Astrophys.* 621:A26
- Pearson K, Bell J. 1908. On some points with regard to the light-fluctuation of variable stars. *Mon. Not. Royal Astro. Soc.* 69:128
- Peebles PJE. 1973. Statistical Analysis of Catalogs of Extragalactic Objects. I. Theory. *Astrophys. J.* 185:413–440
- Percival DB, Walden AT. 1993. Spectral Analysis for Physical Applications. Cambridge University Press
- Platen E, Van De Weygaert R, Jones BJT. 2007. A cosmic watershed: the WVF void detection technique. *Monthly Notices of the Royal Astronomical Society* 380:551–570
- Rau MM, Seitz S, Brimioulle F, Frank E, Friedrich O, et al. 2015. Accurate photometric redshift probability density estimation - method comparison and application. *Mon. Not. Royal Astro. Soc.* 452:3710–3725
- Reed W, Jorgensen M. 2004. The double pareto-lognormal distributiona new parametric model for size distributions. *Communications in Statistics. Theory and Methods* 8:1733–1753
- Reiman DM, Göhre BE. 2019. Deblending galaxy superpositions with branched generative adversarial networks. *Mon. Not. Royal Astro. Soc.* 485:2617–2627
- Revsbech EA, Trotta R, van Dyk DA. 2018. STACCATO: a novel solution to supernova photometric classification with biased training sets. *MNRAS* 473:3969–3986
- Riechers DA, Pavesi R, Sharon CE, Hodge JA, Decarli R, et al. 2019. COLDz: Shape of the CO Luminosity Function at High Redshift and the Cold Gas History of the Universe. *Astrophys. J.* 872:7
- Ronneberger O, Fischer P, Brox T. 2015. U-net: Convolutional networks for biomedical image segmentation, In *Medical Image Computing and Computer-Assisted Intervention – MICCAI 2015*, eds. N Navab, J Hornegger, WM Wells, AF Frangi, pp. 234–241, Cham: Springer International Publishing
- Sahni V, Sathyaprakash BS, Shandarin SF. 1998. Shapefinders: A New Shape Diagnostic for Large-Scale Structure. *Astrophys. J.l* 495:L5–L8
- Salvato M, Ilbert O, Hoyle B. 2019. The many flavours of photometric redshifts. *Nature Astronomy* 3:212–222
- Scargle JD. 1982. Studies in astronomical time series analysis. II. Statistical aspects of spectral analysis of unevenly spaced data. *Astrophys. J.* 263:835–853

- Schafer CM. 2015. A Framework for Statistical Inference in Astrophysics. *Annual Review of Statistics and Its Application* 2:141–162
- Schafer CM, Freeman PE. 2012. Likelihood-free inference in cosmology: Potential for the estimation of luminosity functions. In *Statistical Challenges in Modern Astronomy V*. Springer, 3–19
- Schawinski K, Zhang C, Zhang H, Fowler L, Santhanam GK. 2017. Generative adversarial networks recover features in astrophysical images of galaxies beyond the deconvolution limit. *Mon. Not. Royal Astro. Soc.* 467:L110–L114
- Schlesinger F. 1910. The determination of the orbit of a spectroscopic binary by the method of least-squares. *Publications of the Allegheny Observatory of the University of Pittsburgh* 1:33–44
- Schmidt SJ, Malz AI, Soo JYH, Almosallam IA, Brescia M, et al. 2020. Evaluation of probabilistic photometric redshift estimation approaches for LSST. *arXiv e-prints* :arXiv:2001.03621
- Schmitz MA, Heitz M, Bonneel N, Ngol F, Coeurjolly D, et al. 2018. Wasserstein dictionary learning: Optimal transport-based unsupervised nonlinear dictionary learning. *SIAM Journal on Imaging Sciences* 11:643–678
- Schwarzenberg-Czerny A. 1996. Fast and Statistically Optimal Period Search in Uneven Sampled Observations. *Astrophys. J.l* 460:L107
- Settles B. 2012. Active learning. Morgan & Claypool Publishers
- Sharma S. 2017. Markov Chain Monte Carlo Methods for Bayesian Data Analysis in Astronomy. *Ann. Rev. Astron. Astrophys.* 55:213–259
- Sheynin OB. 1974. On the prehistory of the theory of probability. *Archive for History of Exact Sciences* 12:97–141
- Springel V, White SDM, Jenkins A, Frenk CS, Yoshida N, et al. 2005. Simulations of the formation, evolution and clustering of galaxies and quasars. *Nature* 435:629–636
- Starck JL, Siebenmorgen R, Gredel R. 1997. Spectral Analysis Using the Wavelet Transform. *ApJ* 482:1011–1020
- Stellingwerf RF. 1978. Period determination using phase dispersion minimization. *Astrophys. J.* 224:953–960
- Stigler S. 1986. The history of statistics: The measurement of uncertainty before 1900. Belknap Series of Harvard University. Belknap Press of Harvard University Press
- Stoica RS, Martínez, V. J., Saar, E. 2010. Filaments in observed and mock galaxy catalogues. *A&A* 510:A38
- Sulis S, Mary D, Bigot L. 2017. A bootstrap method for sinusoid detection in colored noise and uneven sampling. application to exoplanet detection, In *25th European Signal Processing Conference (EUSIPCO)*, pp. 1095–1099
- Süveges M, Guy LP, Eyer L, Cuypers J, Holl B, et al. 2015. A comparative study of four significance measures for periodicity detection in astronomical surveys. *Mon. Not. Royal Astro. Soc.* 450:2052–2066
- Taylor SR, Vallisneri M, Ellis JA, Mingarelli CMF, Lazio TJW, van Haasteren R. 2016. Are We There Yet? Time to Detection of Nanohertz Gravitational Waves Based on Pulsar-timing Array Limits. *Astrophys. J.l* 819:L6
- The PLAsTiCC team, Allam Tarek J, Bahmanyar A, Biswas R, Dai M, et al. 2018. The Photometric LSST Astronomical Time-series Classification Challenge (PLAsTiCC): Data set. *arXiv e-prints* :arXiv:1810.00001
- van Dyk DA, DeGennaro S, Stein N, Jefferys WH, von Hippel T. 2009. Statistical analysis of stellar evolution. *Ann. Appl. Stat.* 3:117–143
- van Haasteren R, Vallisneri M. 2014. New advances in the Gaussian-process approach to pulsar-timing data analysis. *Physical Review D* 90:104012
- VanderPlas JT. 2018. Understanding the Lomb-Scargle Periodogram. *Astrophys. J.s* 236:16
- Vilalta R, Ishida EEO, Beck R, Sutrisno R, de Souza RS, Mahabal A. 2017. Photometric redshift estimation: An active learning approach, In *2017 IEEE Symposium Series on Computational Intelligence (SSCI)*, pp. 1–8

- Weyant A, Schafer C, Wood-Vasey WM. 2013. Likelihood-free Cosmological Inference with Type Ia Supernovae: Approximate Bayesian Computation for a Complete Treatment of Uncertainty. *Astrophys. J.* 764:116
- Witzel G, Martinez G, Hora J, Willner SP, Morris MR, et al. 2018. Variability Timescale and Spectral Index of Sgr A* in the Near Infrared: Approximate Bayesian Computation Analysis of the Variability of the Closest Supermassive Black Hole. *Astrophys. J.* 863:15
- Zeldovich IB, Einasto J, Shandarin SF. 1982. Giant voids in the Universe. *Nature* 300:407–413
- Zwicky F. 1937. On the Masses of Nebulae and of Clusters of Nebulae. *Astrophys. J.* 86:217