

Statistics, Data Mining, and Machine Learning in Astronomy

A Practical Python Guide for the Analysis of Survey Data

Željko Ivezić, Andrew J. Connolly,

Statistics, Data Mining, and
Machine Learning in Astronomy



PRINCETON SERIES IN MODERN OBSERVATIONAL ASTRONOMY

David N. Spergel, SERIES EDITOR

Written by some of the world's leading astronomers, the Princeton Series in Modern Observational Astronomy addresses the needs and interests of current and future professional astronomers. International in scope, the series includes cutting-edge monographs and textbooks on topics generally falling under the categories of wavelength, observational techniques and instrumentation, and observational objects from a multiwavelength perspective.

Statistics, Data Mining, and Machine Learning in Astronomy

A PRACTICAL PYTHON GUIDE FOR THE ANALYSIS OF SURVEY DATA



Željko Ivezić, Andrew J. Connolly,
Jacob T. VanderPlas, and Alexander Gray

Copyright © 2014 by Princeton University Press
Published by Princeton University Press, 41 William Street,
Princeton, New Jersey 08540

In the United Kingdom: Princeton University Press, 6 Oxford Street,
Woodstock, Oxfordshire OX20 1TW

press.princeton.edu

All Rights Reserved

ISBN 978-0-691-15168-7

Library of Congress Control Number: 2013951369

British Library Cataloging-in-Publication Data is available

This book has been composed in Minion Pro w/ Universe light condensed for display
Printed on acid-free paper ∞

Typeset by S R Nova Pvt Ltd, Bangalore, India

Printed in the United States of America

10 9 8 7 6 5 4 3 2 1

Contents

<i>Preface</i>	vii
----------------	-----

I Introduction	1
1 About the Book and Supporting Material	3
1.1 What Do Data Mining, Machine Learning, and Knowledge Discovery Mean?	3
1.2 What is This Book About?	5
1.3 An Incomplete Survey of the Relevant Literature	8
1.4 Introduction to the Python Language and the Git Code Management Tool	12
1.5 Description of Surveys and Data Sets Used in Examples	14
1.6 Plotting and Visualizing the Data in This Book	31
1.7 How to Efficiently Use This Book	37
References	39
2 Fast Computation on Massive Data Sets	43
2.1 Data Types and Data Management Systems	43
2.2 Analysis of Algorithmic Efficiency	44
2.3 Seven Types of Computational Problem	46
2.4 Seven Strategies for Speeding Things Up	47
2.5 Case Studies: Speedup Strategies in Practice	50
References	63
II Statistical Frameworks and Exploratory Data Analysis	67
3 Probability and Statistical Distributions	69
3.1 Brief Overview of Probability and Random Variables	70
3.2 Descriptive Statistics	78
3.3 Common Univariate Distribution Functions	85
3.4 The Central Limit Theorem	105
3.5 Bivariate and Multivariate Distribution Functions	108
3.6 Correlation Coefficients	115
3.7 Random Number Generation for Arbitrary Distributions	118
References	122

4 Classical Statistical Inference	123
4.1 Classical vs. Bayesian Statistical Inference	123
4.2 Maximum Likelihood Estimation (MLE)	124
4.3 The goodness of Fit and Model Selection	131
4.4 ML Applied to Gaussian Mixtures: The Expectation Maximization Algorithm	134
4.5 Confidence Estimates: the Bootstrap and the Jackknife	140
4.6 Hypothesis Testing	144
4.7 Comparison of Distributions	149
4.8 Nonparametric Modeling and Histograms	163
4.9 Selection Effects and Luminosity Function Estimation	166
4.10 Summary	172
References	172
5 Bayesian Statistical Inference	175
5.1 Introduction to the Bayesian Method	176
5.2 Bayesian Priors	180
5.3 Bayesian Parameter Uncertainty Quantification	185
5.4 Bayesian Model Selection	186
5.5 Nonuniform Priors: Eddington, Malmquist, and Lutz–Kelker Biases	191
5.6 Simple Examples of Bayesian Analysis: Parameter Estimation	196
5.7 Simple Examples of Bayesian Analysis: Model Selection	223
5.8 Numerical Methods for Complex Problems (MCMC)	229
5.9 Summary of Pros and Cons for Classical and Bayesian methods	239
References	243
III Data Mining and Machine Learning	247
6 Searching for Structure in Point Data	249
6.1 Nonparametric Density Estimation	250
6.2 Nearest-Neighbor Density Estimation	257
6.3 Parametric Density Estimation	259
6.4 Finding Clusters in Data	270
6.5 Correlation Functions	277
6.6 Which Density Estimation and Clustering Algorithms Should I Use?	281
References	285
7 Dimensionality and Its Reduction	289
7.1 The Curse of Dimensionality	289
7.2 The Data Sets Used in This Chapter	291
7.3 Principal Component Analysis	292
7.4 Nonnegative Matrix Factorization	305
7.5 Manifold Learning	306
7.6 Independent Component Analysis and Projection Pursuit	313
7.7 Which Dimensionality Reduction Technique Should I Use?	316
References	318

8 Regression and Model Fitting	321
8.1 Formulation of the Regression Problem	321
8.2 Regression for Linear Models	325
8.3 Regularization and Penalizing the Likelihood	332
8.4 Principal Component Regression	337
8.5 Kernel Regression	338
8.6 Locally Linear Regression	339
8.7 Nonlinear Regression	340
8.8 Uncertainties in the Data	342
8.9 Regression that is Robust to Outliers	344
8.10 Gaussian Process Regression	349
8.11 Overfitting, Underfitting, and Cross-Validation	352
8.12 Which Regression Method Should I Use?	361
References	363
9 Classification	365
9.1 Data Sets Used in This Chapter	365
9.2 Assigning Categories: Classification	366
9.3 Generative Classification	368
9.4 K -Nearest-Neighbor Classifier	378
9.5 Discriminative Classification	380
9.6 Support Vector Machines	382
9.7 Decision Trees	386
9.8 Evaluating Classifiers: ROC Curves	394
9.9 Which Classifier Should I Use?	397
References	400
10 Time Series Analysis	403
10.1 Main Concepts for Time Series Analysis	404
10.2 Modeling Toolkit for Time Series Analysis	405
10.3 Analysis of Periodic Time Series	426
10.4 Temporally Localized Signals	452
10.5 Analysis of Stochastic Processes	455
10.6 Which Method Should I Use for Time Series Analysis?	465
References	465
IV Appendices	469
A An Introduction to Scientific Computing with Python	471
A.1 A Brief History of Python	471
A.2 The SciPy Universe	472
A.3 Getting Started with Python	474
A.4 IPython: The Basics of Interactive Computing	486
A.5 Introduction to NumPy	488
A.6 Visualization with Matplotlib	494
A.7 Overview of Useful NumPy/SciPy Modules	498

A.8 Efficient Coding with Python and NumPy	503
A.9 Wrapping Existing Code in Python	506
A.10 Other Resources	508
B AstroML: Machine Learning for Astronomy	511
B.1 Introduction	511
B.2 Dependencies	511
B.3 Tools Included in AstroML v0.1	512
C Astronomical Flux Measurements and Magnitudes	515
C.1 The Definition of the Specific Flux	515
C.2 Wavelength Window Function for Astronomical Measurements	515
C.3 The Astronomical Magnitude Systems	516
D SQL Query for Downloading SDSS Data	519
E Approximating the Fourier Transform with the FFT	521
References	525
<i>Visual Figure Index</i>	527
<i>Index</i>	533

Preface

Astronomy and astrophysics are witnessing dramatic increases in data volume as detectors, telescopes, and computers become ever more powerful. During the last decade, sky surveys across the electromagnetic spectrum have collected hundreds of terabytes of astronomical data for hundreds of millions of sources. Over the next decade, data volumes will enter the petabyte domain, and provide accurate measurements for billions of sources. Astronomy and physics students are not traditionally trained to handle such voluminous and complex data sets. Furthermore, standard analysis methods employed in astronomy often lag far behind the rapid progress in statistics and computer science. The main purpose of this book is to help minimize the time it takes a student to become an effective researcher.

This book provides the interface between astronomical data analysis problems and modern statistical methods. It is aimed at physical and data-centric scientists who have an understanding of the science drivers for analyzing large data sets but may not be aware of developments in statistical techniques over the last decade. The book targets researchers who want to use existing methods for the analysis of large data sets, rather than those interested in the development of new methods. Theoretical discussions are limited to the minimum required to understand the algorithms. Nevertheless, extensive and detailed references to relevant specialist literature are interspersed throughout the book.

We present an example-driven compendium of modern statistical and data mining methods, together with carefully chosen examples based on real modern data sets, and of current astronomical applications that will illustrate each method introduced in the book. The book is loosely organized by practical analysis problems, and offers a comparative analysis of different techniques, including discussions of the advantages and shortcomings of each method, and their scaling with the sample size. The exposition of the material is supported by appropriate publicly available Python code (available from the book website, rather than fully printed here) and data to enable a reader to reproduce all the figures and examples, evaluate the techniques, and adapt them to their own field of interest. To some extent, this book is an analog of the well-known *Numerical Recipes* book, but aimed at the analysis of massive astronomical data sets, with more emphasis on modern tools for data mining and machine learning, and with freely available code.

From the start, we desired to create a book which, in the spirit of reproducible research, would allow readers to easily replicate the analysis behind every example and figure. We believe this feature will make the book uniquely valuable as a practical guide. We chose to implement this using Python, a powerful and flexible programming language that is quickly becoming a standard in astronomy (a number of next-generation large astronomical surveys and projects use Python, e.g., JVLA, ALMA, LSST). The Python code base associated with this book, called AstroML, is maintained as a live web repository (GitHub), and is intended to be a growing collection of well-documented and well-tested tools for astronomical research. Any astronomical researcher who is currently developing

software for the analysis of massive survey data is invited and encouraged to contribute their own tools to the code.

The target audience for this text includes senior undergraduate and graduate students in physics and astronomy, as well as researchers using large data sets in a scientific context. Familiarity with calculus and other basic mathematical techniques is assumed, but no extensive prior knowledge in statistics is required (e.g., we assume that readers have heard of the Gaussian distribution, but not necessarily of the Lorentzian distribution). Though the examples in this book are aimed at researchers in the fields of astronomy and astrophysics, the organization of the book allows for easy mapping of relevant algorithms to problems from other fields. After the first introductory chapter, data organization and some aspects of fast computation are discussed in chapter 2, statistical foundations are reviewed in chapters 3–5 (statistical distributions, maximum likelihood and other classical statistics, and Bayesian methodology), exploratory data analysis is described in chapters 6 and 7 (“Searching for Structure in Point Data”; “Dimensionality and its Reduction”), and data-based prediction methods are described in chapters 8–10 (“Regression and Model Fitting”; “Classification”; “Time Series Analysis”).

Finally, we are indebted to a number of colleagues whose careful reading and resulting comments significantly improved this book. A summer study group consisting of Bob Abel, Yusra AlSayyad, Lauren Anderson, Vaishali Bhardwaj, James Davenport, Alexander Fry, Bryce Kalmbach, and David Westman identified many rough edges in the manuscript and tested the AstroML code. We thank Alan Weinstein for help and advice with LIGO data, and Carl Carter-Schwendler for motivational and expert discussions about Bayesian statistics. In addition, Tim Axelrod, Andy Becker, Joshua Bloom, Tamás Budavári, David Hogg, Robert Lupton, Chelsea MacLeod, Lovro Palaversa, Fernando Perez, Maria Süveges, Przemek Woźniak, and two anonymous reviewers provided extensive expert comments. Any remaining errors are entirely our own.

We dedicate this book to Cristin, Hieu, Ian, Nancy, Pamela, Tom, and Vedrana for their support, encouragement, and understanding during the periods of intensive work and absent-mindedness along the way to this finished text.

Authors, Seattle and Atlanta, 2012

PART I
Introduction



1 About the Book and Supporting Material

“Even the longest journey starts with the first step.” (Lao-tzu paraphrased)

This chapter introduces terminology and nomenclature, reviews a few relevant contemporary books, briefly describes the Python programming language and the Git code management tool, and provides details about the data sets used in examples throughout the book.

1.1. What Do Data Mining, Machine Learning, and Knowledge Discovery Mean?

Data mining, machine learning, and knowledge discovery refer to research areas which can all be thought of as outgrowths of multivariate statistics. Their common themes are analysis and interpretation of data, often involving large quantities of data, and even more often resorting to numerical methods. The rapid development of these fields over the last few decades was led by computer scientists, often in collaboration with statisticians. To an outsider, data mining, machine learning, and knowledge discovery compared to statistics are akin to engineering compared to fundamental physics and chemistry: applied fields that “make things work.” The techniques in all of these areas are well studied, and rest upon the same firm statistical foundations. In this book we will consider those techniques which are most often applied in the analysis of astronomical data.

While there are many varying definitions in the literature and on the web, we adopt and are happy with the following:

- **Data mining** is a set of techniques for analyzing and describing structured data, for example, finding patterns in large data sets. Common methods include density estimation, unsupervised classification, clustering, principal component analysis, locally linear embedding, and projection pursuit. Often, the term “knowledge discovery” is used interchangeably with data mining. Although there are many books written with “knowledge discovery” in their title, we shall uniformly adopt “data mining” in this book. The data mining techniques result in the understanding of data set properties, such as “My measurements of the size and temperature of stars form a well-defined sequence

in the size–temperature diagram, though I find some stars in three clusters far away from this sequence.” From the data mining point of view, it is not important to immediately contrast these data with a model (of stellar structure in this case), but rather to quantitatively describe the “sequence,” as well as the behavior of measurements falling “far away” from it. In short, data mining is about what the data themselves are telling us. Chapters 6 and 7 in this book primarily discuss data mining techniques.

- **Machine learning** is an umbrella term for a set of techniques for interpreting data by comparing them to models for data behavior (including the so-called nonparametric models), such as various regression methods, supervised classification methods, maximum likelihood estimators, and the Bayesian method. They are often called inference techniques, data-based statistical inferences, or just plain old “fitting.” Following the above example, a physical stellar structure model can predict the position and shape of the so-called main sequence in the size–temperature diagram for stars, and when combined with galaxy formation and evolution models, the model can even predict the distribution of stars away from that sequence. Then, there could be more than one competing model and the data might tell us whether (at least) one of them can be rejected. Chapters 8–10 in this book primarily discuss machine learning techniques.

Historically, the emphasis in data mining and knowledge discovery has been on what statisticians call *exploratory data analysis*: that is, learning qualitative features of the data that were not previously known. Much of this is captured under the heading of “unsupervised learning” techniques. The emphasis in machine learning has been on *prediction* of one variable based on the other variables—much of this is captured under the heading of “supervised learning.” For further discussion of data mining and machine learning in astronomy, see recent informative reviews [3, 7, 8, 10].

Here are a few concrete examples of astronomical problems that can be solved with data mining and machine learning techniques, and which provide an illustration of the scope and aim of this book. For each example, we list the most relevant chapter(s) in this book:

- Given a set of luminosity measurements for a sample of sources, quantify their luminosity distribution (the number of sources per unit volume and luminosity interval). *Chapter 3*
- Determine the luminosity distribution if the sample selection function is controlled by another measured variable (e.g., sources are detected only if brighter than some flux sensitivity limit). *Chapters 3 and 4*
- Determine whether a luminosity distribution determined from data is statistically consistent with a model-based luminosity distribution. *Chapters 3–5*
- Given a signal in the presence of background, determine its strength. *Chapter 5*
- For a set of brightness measurements with suspected outliers, estimate the best value of the intrinsic brightness. *Chapter 5*
- Given measurements of sky coordinates and redshifts for a sample of galaxies, find clusters of galaxies. *Chapter 6*
- Given several brightness measurements per object for a large number of objects, identify and quantitatively describe clusters of sources in the multi-dimensional color space. Given color measurements for an additional set of

sources, assign to each source the probabilities that it belongs to each of the clusters, making use of both measurements and errors. *Chapters 6 and 9*

- Given a large number of spectra, find self-similar classes. *Chapter 7*
- Given several color and other measurements (e.g., brightness) for a galaxy, determine its most probable redshift using (i) a set of galaxies with both these measurements and their redshift known, or (ii) a set of models predicting color distribution (and the distribution of other relevant parameters). *Chapters 6–8*
- Given a training sample of stars with both photometric (color) measurements and spectroscopic temperature measurements, develop a method for estimating temperature using only photometric measurements (including their errors). *Chapter 8*
- Given a set of redshift and brightness measurements for a cosmological supernova sample, estimate the cosmological parameters and their uncertainties. *Chapter 8*
- Given a set of position (astrometric) measurements as a function of time, determine the best-fit parameters for a model including proper motion and parallax motion. *Chapters 8 and 10*
- Given colors for a sample of spectroscopically confirmed quasars, use analogous color measurements to separate quasars from stars in a larger sample. *Chapter 9*
- Given light curves for a large number of sources, find variable objects, identify periodic light curves, and classify sources into self-similar classes. *Chapter 10*
- Given unevenly sampled low signal-to-noise time series, estimate the underlying power spectrum. *Chapter 10*
- Given detection times for individual photons, estimate model parameters for a suspected exponentially decaying burst. *Chapter 10*

1.2. What Is This Book About?

This book is about extracting knowledge from data, where “knowledge” means a quantitative summary of data behavior, and “data” essentially means the results of measurements. Let us start with the simple case of a scalar quantity, x , that is measured N times, and use the notation x_i for a single measurement, with $i = 1, \dots, N$. We will use $\{x_i\}$ to refer to the set of all N measurements. In statistics, the data x are viewed as realizations of a random variable X (random variables are functions on the sample space, or the set of all outcomes of an experiment). In most cases, x is a real number (e.g., stellar brightness measurement) but it can also take discrete values (e.g., stellar spectral type); missing data (often indicated by the special IEEE floating-point value NaN—Not a Number) can sometimes be found in real-life data sets.

Possibly the most important single problem in data mining is how to estimate the distribution $h(x)$ from which values of x are drawn (or which “generates” x). The function $h(x)$ quantifies the probability that a value lies between x and $x + dx$, equal to $h(x) dx$, and is called a *probability density function* (pdf). Astronomers sometimes use the terms “differential distribution function” or simply “probability distribution.” When x is discrete, statisticians use the term “probability mass function” (note that “density” and “mass” are already reserved words in physical sciences, but the

confusion should be minimal due to contextual information). The integral of the pdf,

$$H(x) = \int_{-\infty}^x h(x') dx', \quad (1.1)$$

is called the “cumulative distribution function” (cdf). The inverse of the cumulative distribution function is called the “quantile function.”

To distinguish the true pdf $h(x)$ (called the *population pdf*) from a data-derived estimate (called the *empirical pdf*), we shall call the latter $f(x)$ (and its cumulative counterpart $F(x)$).¹ Hereafter, we will assume for convenience that both $h(x)$ and $f(x)$ are properly normalized probability density functions (though this is not a necessary assumption), that is,

$$H(\infty) = \int_{-\infty}^{+\infty} h(x') dx' = 1 \quad (1.2)$$

and analogously for $F(\infty)$. Given that data sets are never infinitely large, $f(x)$ can never be exactly equal to $h(x)$. Furthermore, we shall also consider cases when measurement errors for x are not negligible and thus $f(x)$ will not tend to $h(x)$ even for an infinitely large sample (in this case $f(x)$ will be a “broadened” or “blurred” version of $h(x)$).

$f(x)$ is a *model* of the true distribution $h(x)$. Only samples from $h(x)$ are observed (i.e., data points); the functional form of $h(x)$, used to constrain the model $f(x)$, must be guessed. Such forms can range from relatively simple *parametric* models, such as a single Gaussian, to much more complicated and flexible *nonparametric* models, such as the superposition of many small Gaussians. Once the functional form of the model is chosen, the best-fitting member of that model family, corresponding to the best setting of the model’s parameters (such as the Gaussian’s mean and standard deviation) must be chosen.

A model can be as simple as an analytic function (e.g., a straight line), or it can be the result of complex simulations and other computations. Irrespective of the model’s origin, it is important to remember that we can never prove that a model is correct; we can only test it against the data, and sometimes reject it. Furthermore, within the Bayesian logical framework, we cannot even reject a model if it is the only one we have at our disposal—we can only compare models against each other and rank them by their success.

These analysis steps are often not trivial and can be quite complex. The simplest nonparametric method to determine $f(x)$ is to use a histogram; bin the x data and count how many measurements fall in each bin. Very quickly several complications arise: First, what is the optimal choice of bin size? Does it depend on the sample size, or other measurement properties? How does one determine the count error in each bin, and can we treat them as Gaussian errors?

¹Note that in this book we depart from a common notation in the statistical literature in which the true distribution is called $f(x)$ (here we use $h(x)$), and the data-derived estimate of the distribution is called $\hat{f}(x)$ (here we use $f(x)$).

An additional frequent complication is that the quantity x is measured with some uncertainty or error distribution, $e(x)$, defined as the probability of measuring value x if the true value is μ ,

$$e(x) = p(x|\mu, I), \quad (1.3)$$

where I stands for all other information that specifies the details of the error distribution, and “|” is read as “given.” Eq. 1.3 should be interpreted as giving a probability $e(x) dx$ that the measurement will be between x and $x + dx$.

For the commonly used Gaussian (or normal) error distribution, the probability is given by

$$p(x|\mu, \sigma) = \frac{1}{\sigma\sqrt{2\pi}} \exp\left(\frac{-(x-\mu)^2}{2\sigma^2}\right), \quad (1.4)$$

where in this case I is simply σ , the standard deviation (it is related to the uncertainty estimate popularly known as the “error bar”; for further discussion of distribution functions, see §3.3). The error distribution function could also include a bias b , and $(x - \mu)$ in the above expression would become $(x - b - \mu)$. That is, the bias b is a *systematic* offset of all measurements from the true value μ , and σ controls their “scatter” (bias is introduced formally in §3.2.2). How exactly the measurements are “scattered around” is described by the *shape* of $e(x)$. In astronomy, error distributions are often non-Gaussian or, even when they are Gaussian, σ might not be the same for all measurements, and often depends on the signal strength (i.e., on x ; each measured x_i is accompanied by a different σ_i). These types of errors are called *heteroscedastic*, as opposed to *homoscedastic* errors in which the error distribution is the same for each point.

Quantities described by $f(x)$ (e.g., astronomical measurements) can have different meanings in practice. A special case often encountered in practice is when the “intrinsic” or “true” (population pdf) $h(x)$ is a delta function, $\delta(x)$; that is, we are measuring some specific single-valued quantity (e.g., the length of a rod; let us ignore quantum effects here and postulate that there is no uncertainty associated with its true value) and the “observed” (empirical pdf) $f(x)$, sampled by our measurements x_i , simply reflects their error distribution $e(x)$. Another special case involves measurements with negligible measurement errors, but the underlying intrinsic or true pdf $h(x)$ has a finite width (as opposed to a delta function). Hence, in addition to the obvious effects of finite sample size, the difference between $f(x)$ and $h(x)$ can have two *very different origins* and this distinction is often not sufficiently emphasized in the literature: at one extreme it can reflect our measurement error distribution (we measure the same rod over and over again to improve our knowledge of its length), and at the other extreme it can represent measurements of a number of different rods (or the same rod at different times, if we suspect its length may vary with time) with measurement errors *much smaller* than the expected and/or observed length variation. Despite being extremes, these two limiting cases are often found in practice, and may sometimes be treated with the same techniques because of their mathematical similarity (e.g., when fitting a Gaussian to $f(x)$, we do not distinguish the case where its width is due to measurement errors from the case when we measure a population property using a finite sample).

The next level of complication when analyzing $f(x)$ comes from the sample size and dimensionality. There can be a large number of different scalar quantities, such as x , that we measure for each object, and each of these quantities can have a different error distribution (and sometimes even different selection function). In addition, some of these quantities may not be statistically independent. When there is more than one dimension, analysis can get complicated and is prone to pitfalls; when there are many dimensions, analysis is always complicated. If the sample size is measured in hundreds of millions, even the most battle-tested algorithms and tools can choke and become too slow.

Classification of a set of measurements is another important data analysis task. We can often “tag” each x measurement by some “class descriptor” (such quantities are called “categorical” in the statistics literature). For example, we could be comparing the velocity of stars, x , around the Galaxy center with subsamples of stars classified by other means as “halo” and “disk” stars (the latter information could be assigned codes H and D , or 0/1, or any other discrete attribute). In such cases, we would determine two independent distributions $f(x)$ —one for each of these two subsamples. Any new measurement of x could then be classified as a “halo” or “disk” star. This simple example can become nontrivial when x is heteroscedastic or multidimensional, and also raises the question of completeness vs. purity trade-offs (e.g., do we care more about never ever misclassifying a halo star, or do we want to minimize the total number of misclassifications for both disk and halo stars?). Even in the case of discrete variables, such as “halo” or “disk” stars, or “star” vs. “galaxy” in astronomical images (which should be called more precisely “unresolved” and “resolved” objects when referring to morphological separation), we can assign them a continuous variable, which often is interpreted as the probability of belonging to a class. At first it may be confusing to talk about the probability that an object is a star vs. being a galaxy because it cannot be both at the same time. However, in this context we are talking about *our current state of knowledge about a given object* and its classification, which can be elegantly expressed using the framework of probability.

In summary, this book is mostly about how to estimate the empirical pdf $f(x)$ from data (including multidimensional cases), how to statistically describe the resulting estimate and its uncertainty, how to compare it to models specified via $h(x)$ (including estimates of model parameters that describe $h(x)$), and how to use this knowledge to interpret additional and/or new measurements (including best-fit model reassessment and classification).

1.3. An Incomplete Survey of the Relevant Literature

The applications of data mining and machine learning techniques are not limited to the sciences. A large number of books discuss applications such as data mining for marketing, music data mining, and machine learning for the purposes of counter-terrorism and law enforcement. We shall limit our survey to books that cover topics similar to those from this book but from a different point of view, and can thus be used as supplemental literature. In many cases, we reference specific sections in the following books.

Numerical Recipes: The Art of Scientific Computing by Press, Teukolsky, Vetterling, and Flannery [27] is famous for its engaging text and concise mathematical

and algorithmic explanations (its Fortran version has been cited over 8000 times at the time of writing this book, according to the SAO/NASA Astrophysics Data System). While the whole book is of great value for the topics covered here, several of its 22 chapters are particularly relevant (“Random Numbers,” “Sorting and Selection,” “Fourier and Spectral Applications,” “Statistical Description of Data,” “Modeling of Data,” “Classification and Inference”). The book includes commented full listings of more than 400 numerical routines in several computer languages that can be purchased in machine-readable form. The supplemental code support for the material covered in the book served as a model for our book. We refer to this book as NumRec.

The Elements of Statistical Learning: Data Mining, Inference, and Prediction by Hastie, Tibshirani, and Friedman [16] is a classic book on these topics, and highly recommended for further reading. With 18 chapters and about 700 pages, it is more comprehensive than this book, and many methods are discussed in greater detail. The writing style is not heavy on theorems and the book should be easily comprehensible to astronomers and other physical scientists. It comes without computer code. We refer to this book as HTF09.

Two books by Wasserman, *All of Nonparametric Statistics* [39] and *All of Statistics: A Concise Course in Statistical Inference* [40] are closer to the statistician’s heart, and do not shy away from theorems and advanced statistics. Although “All” may imply very long books, together they are under 700 pages. They are good books to look into for deeper and more formal expositions of statistical foundations for data mining and machine learning techniques. We refer to these books as Wass10.

Statistics in Theory and Practice by Lupton [23] is a concise (under 200 pages) summary of the most important concepts in statistics written for practicing scientists, and with close to 100 excellent exercises (with answers). For those who took statistics in college, but need to refresh and extend their knowledge, this book is a great choice. We refer to this book as Lup93.

Practical Statistics for Astronomers by Wall and Jenkins [38] is a fairly concise (under 300 pages) summary of the most relevant contemporary statistical and probabilistic technology in observational astronomy. This excellent book covers classical parametric and nonparametric methods with a strong emphasis on Bayesian solutions. We refer to this book as WJ03.

Statistics: A Guide to the Use of Statistical Methods in the Physical Sciences by Barlow [4] is an excellent introductory text written by a physicist (200 pages). We highly recommend it as a starting point if you feel that the books by Lupton, and Wall and Jenkins are too advanced. We refer to this book as Bar89.

Data Analysis: A Bayesian Tutorial by Sivia [31] is an excellent short book (under 200 pages) to quickly learn about basic Bayesian ideas and methods. Its examples are illuminating and the style is easy to read and does not presume any prior knowledge of statistics. We highly recommend it! We refer to this book as Siv06.

Bayesian Logical Data Analysis for the Physical Sciences by Gregory [13] is more comprehensive (over 400 pages) than Sivia’s book, and covers many topics discussed here. It is a good book to look into for deeper understanding and implementation details for most frequently used Bayesian methods. It also provides code support (for Mathematica). We refer to this book as Greg05.

Probability Theory: The Logic of Science by Jaynes [20], an early and strong proponent of Bayesian methods, describes probability theory as extended logic. This

monumental treatise compares Bayesian analysis with other techniques, including a large number of examples from the physical sciences. The book is aimed at readers with a knowledge of mathematics at a graduate or an advanced undergraduate level. We refer to this book as Jay03.

Bayesian Methods in Cosmology provides an introduction to the use of Bayesian methods in cosmological studies [17]. Contributions from 24 cosmologists and statisticians (edited by M. P. Hobson, A. H. Jaffe, A. R. Liddle, P. Mukherjee, and D. Parkinson) range from the basic foundations to detailed descriptions of state-of-the-art techniques. The book is aimed at graduate students and researchers in cosmology, astrophysics, and applied statistics. We refer to this book as BayesCosmo.

Advances in Machine Learning and Data Mining for Astronomy is a recent book by over 20 coauthors from mostly astronomical backgrounds (edited by M. J. Way, J. D. Scargle, K. Ali, and A. N. Srivastava) [41]. This book provides a comprehensive overview (700 pages) of various data mining tools and techniques that are increasingly being used by astronomers, and discusses how current problems could lead to the development of entirely new algorithms. We refer to this book as WSAS.

Modern Statistical Methods for Astronomy With R Applications by Feigelson and Babu [9] is very akin in spirit to this book. It provides a comprehensive (just under 500 pages) coverage of similar topics, and provides examples written in the R statistical software environment. Its first chapter includes a very informative summary of the history of statistics in astronomy, and the number of references to statistics literature is larger than here. We refer to this book as FB2012.

Although not referenced further in this book, we highly recommend the following books as supplemental resources.

Pattern Recognition and Machine Learning by Bishop [6] provides a comprehensive introduction to the fields of pattern recognition and machine learning, and is aimed at advanced undergraduates and graduate students, as well as researchers and practitioners. The book is supported by a great deal of additional material, including lecture slides as well as the complete set of figures used in the book. It is of particular interest to those interested in Bayesian versions of standard machine learning methods.

Information Theory, Inference, and Learning Algorithms by MacKay [25] is an excellent and comprehensive book (over 600 pages) that unites information theory and statistical inference. In addition to including a large fraction of the material covered in this book, it also discusses other relevant topics, such as arithmetic coding for data compression and sparse-graph codes for error correction. Throughout, it addresses a wide range of topics—from evolution to sex to crossword puzzles—from the viewpoint of information theory. The book level and style should be easily comprehensible to astronomers and other physical scientists.

In addition to books, several other excellent resources are readily available.

The R language is familiar to statisticians and is widely used for statistical software development and data analysis. R is available as a free software environment² for statistical computing and graphics, and compiles and runs on a wide variety of UNIX

²<http://www.R-project.org/>

platforms, Windows and Mac OS. The capabilities of R are extended through user-created packages, which allow specialized statistical techniques, graphical devices, import/export capabilities, reporting tools, etc.

The Auton Lab, part of Carnegie Mellon University's School of Computer Science, researches new approaches to statistical data mining. The Lab is "interested in the underlying computer science, mathematics, statistics and AI of detection and exploitation of patterns in data." A large collection of software, papers, and other resources are available from the Lab's homepage.³

The IVOA (International Virtual Observatory Alliance) Knowledge Discovery in Databases group⁴ provides support to the IVOA by developing and testing scalable data mining algorithms and the accompanying new standards for VO interfaces and protocols. Their web pages contain tutorials and other materials to support the VO users (e.g., "A user guide for Data Mining in Astronomy").

The Center for Astrostatistics at Penn State University organizes annual summer schools in statistics designed for graduate students and researchers in astronomy. The school is an intensive week covering basic statistical inference, applied statistics, and the R computing environment. The courses are taught by a team of statistics and astronomy professors with opportunity for discussion of methodological issues. For more details, please see their website.⁵

The burgeoning of work in what has been called "astrostatistics" or "astroinformatics," along with the slow but steady recognition of its importance within astronomy, has given rise to recent activity to define and organize more cohesive communities around these topics, as reflected in manifestos by Loredo et al. [22] and Borne et al. [8]. Recent community organizations include the American Astronomical Society Working Group in Astroinformatics and Astrostatistics, the International Astronomical Union Working Group in Astrostatistics and Astroinformatics, and the International Astrostatistics Association (affiliated with the International Statistical Institute). These organizations promote the use of known advanced statistical and computational methods for astronomical research, encourage the development of new procedures and algorithms, organize multidisciplinary meetings, and provide educational and professional resources to the wider community. Information about these organizations can be found at the Astrostatistics and Astroinformatics Portal.⁶

With all these excellent references and resources already available, it is fair to ask why we should add yet another book to the mix. There are two main reasons that motivate this book: First, it is convenient to have the basic statistical, data mining, and machine learning techniques collected and described in a single book, and at a level of mathematical detail aimed at researchers entering into astronomy and physical sciences. This book grew out of materials developed for several graduate classes. These classes had to rely on a large number of textbooks with strongly varying styles and difficulty levels, which often caused practical problems. Second, when bringing a new student up to speed, one difficulty with the current array of texts on data mining and machine learning is that the implementation of the discussed methods is typically

³<http://www.autonlab.org/>

⁴<http://www.ivoa.net/cgi-bin/twiki/bin/view/IVOA/IvoaKDD>

⁵<http://astrostatistics.psu.edu/>

⁶<http://asaip.psu.edu>

left up to the reader (with some exceptions noted above). The lack of ready-to-use tools has led to a situation where many groups have independently implemented desired methods and techniques in a variety of languages. This reinventing of the wheel not only takes up valuable time, but the diverse approaches make it difficult to share data and compare results between different groups. With this book and the associated online resources, we hope to encourage, and to contribute to a common implementation of the basic statistical tools.

1.4. Introduction to the Python Language and the Git Code Management Tool

The material in this book is supported by publicly available code, available from <http://www.astroml.org>. The site includes the Python code to repeat all the examples in this text, as well as to reproduce all the figures from the book. We do not refer by name to code used to produce the figures because the code listing on the website is enumerated by the figure number in the book and thus is easy to locate. We believe and hope that these code examples, with minor modifications, will provide useful templates for your own projects. In this section, we first introduce the Python programming language and then briefly describe the code management tool Git.

1.4.1. Python

Python is an open-source, object-oriented interpreted language with a well-developed set of libraries, packages, and tools for scientific computation. In appendix A, we offer a short introduction to the key features of the language and its use in scientific computing. In this section, we will briefly list some of the scientific computing packages and tools available in the language, as well as the requirements for running the examples in this text.

The examples and figures in this text were created with the Python package AstroML, which was designed as a community resource for fast, well-tested statistical, data mining, and machine learning tools implemented in Python (see appendix B). Rather than reimplementing common algorithms, AstroML draws from the wide range of efficient open-source computational tools available in Python. We briefly list these here; for more detailed discussion see appendix A.

The core packages for scientific computing in Python are NumPy,⁷ SciPy,⁸ and Matplotlib.⁹ Together, these three packages allow users to efficiently read, store, manipulate, and visualize scientific data. Many of the examples and figures in this text require only these three dependencies, and they are discussed at length in appendix A.

There are a large number of other packages built upon this foundation, and AstroML makes use of several of them. An important one is Scikit-learn,¹⁰ a large and very well-documented collection of machine learning algorithms in Python.

⁷Numerical Python; <http://www.numpy.org>

⁸Scientific Python; <http://www.scipy.org>

⁹<http://matplotlib.org>

¹⁰<http://scikit-learn.org>

Scikit-learn is used extensively in the examples and figures within this text, especially those in the second half of the book. We also make use of PyMC¹¹ for Markov chain Monte Carlo methods, and HealPy¹² for spherical coordinates and spherical harmonic transformations.

There are a number of other useful Python packages that are not used in this book. For example, Erin Sheldon’s esutil package¹³ includes a wide variety of handy utilities, focused primarily on Numerical Python, statistics, and file input/output. The CosmoloPy package¹⁴ includes a set of functions for performing cosmological computations. Kapteyn¹⁵ contains many useful routines for astronomical data manipulation, most notably a very complete set of tools for translating between various sky coordinate systems. The AstroPython site¹⁶ acts as a community knowledge base for performing astronomy research with Python. Useful resources include a discussion forum, a collection of code snippets and various tutorials. AstroPy¹⁷ is another effort in this direction, with a focus on community development of a single core Python package for astronomers.

1.4.2. Code Management with Git

Complex analyses of large data sets typically produce substantial amounts of special-purpose code. It is often easy to end up with an unmanageable collection of different software versions, or lose code due to computer failures. Additional management difficulties are present when multiple developers are working on the same code. Professional programmers address these and similar problems using code management tools. There are various freely available tools such as CVS, SVN, Bazaar, Mercurial, and Git. While they all differ a bit, their basic functionality is similar: they support collaborative development of software and the tracking of changes to software source code over time.

This book and the associated code are managed using Git. Installing¹⁸ Git, using it for code management, and for distributing code, are all very user friendly and easy to learn.¹⁹ Unlike CVS, Git can manage not only changes to files, but new files, deleted files, merged files, and entire file structures.²⁰ One of the most useful features of Git is its ability to set up a remote repository, so that code can be checked in and out from multiple computers. Even when a computer is not connected to a repository (e.g., in the event of a server outage, or when no internet connection is available), the local copy can still be modified and changes reported to the repository later. In the event of a disk failure, the remote repository can even be rebuilt from the local copy.

¹¹<http://pymc-devs.github.com/pymc/>

¹²<http://healpy.readthedocs.org>

¹³<http://code.google.com/p/esutil/>

¹⁴<http://roban.github.com/CosmoloPy/>

¹⁵<http://www.astro.rug.nl/software/kapteyn/>

¹⁶<http://www.astropython.org/>

¹⁷see <http://www.astropy.org/>

¹⁸<http://git-scm.com/>

¹⁹For example, see <http://www.github.com/>

²⁰For a Git manual, see <http://progit.org/book/>

Because of these features, Git has become the de facto standard code management tool in the Python community: most of the core Python packages listed above are managed with Git, using the website <http://github.com> to aid in collaboration. We strongly encourage you to consider using Git in your projects. You will not regret the time spent learning how to use it.

1.5. Description of Surveys and Data Sets Used in Examples

Many of the examples and applications in this book require realistic data sets in order to test their performance. There is an increasing amount of high-quality astronomical data freely available online. However, unless a person knows exactly where to look, and is familiar with database tools such as SQL (Structured Query Language,²¹ for searching databases), finding suitable data sets can be very hard. For this reason, we have created a suite of data set loaders within the package AstroML. These loaders use an intuitive interface to download and manage large sets of astronomical data, which are used for the examples and plots throughout this text. In this section, we describe these data loading tools, list the data sets available through this interface, and show some examples of how to work with these data in Python.

1.5.1. AstroML Data Set Tools

Because of the size of these data sets, bundling them with the source code distribution would not be very practical. Instead, the data sets are maintained on a web page with http access via the data-set scripts in `astroML.datasets`. Each data set will be downloaded to your machine only when you first call the associated function. Once it is downloaded, the cached version will be used in all subsequent function calls.

For example, to work with the SDSS imaging photometry (see below), use the function `fetch_imaging_sample`. The function takes an optional string argument, `data_home`. When the function is called, it first checks the `data_home` directory to see if the data file has already been saved to disk (if `data_home` is not specified, then the default directory is `$HOME/astroML_data/`; alternatively, the `$ASTROML_DATA` environment variable can be set to specify the default location). If the data file is not present in the specified directory, it is automatically downloaded from the web and cached in this location.

The nice part about this interface is that the user does not need to remember whether the data has been downloaded and where it has been stored. Once the function is called, the data is returned whether it is already on disk or yet to be downloaded.

For a complete list of data set fetching functions, make sure AstroML is properly installed in your Python path, and open an IPython terminal and type

```
In [1]: from astroML.datasets import<TAB>
```

The tab-completion feature of IPython will display the available data downloaders (see appendix A for more details on IPython).

²¹ See, for example, <http://en.wikipedia.org/wiki/SQL>

1.5.2. Overview of Available Data Sets

Most of the astronomical data that we make available were obtained by the Sloan Digital Sky Survey²² (SDSS), which operated in three phases starting in 1998. The SDSS used a dedicated 2.5 m telescope at the Apache Point Observatory, New Mexico, equipped with two special-purpose instruments, to obtain a large volume of imaging and spectroscopic data. For more details see [15]. The 120 MP camera (for details see [14]) imaged the sky in five photometric bands (u , g , r , i , and z ; see appendix C for more details about astronomical flux measurements, and for a figure with the SDSS passbands). As a result of the first two phases of SDSS, Data Release 7 has publicly released photometry for 357 million unique sources detected in $\sim 12,000$ deg 2 of sky²³ (the full sky is equivalent to $\sim 40,000$ deg 2). For bright sources, the photometric precision is 0.01–0.02 mag (1–2% flux measurement errors), and the faint limit is $r \sim 22.5$. For more technical details about SDSS, see [1, 34, 42].

The SDSS imaging data were used to select a subset of sources for spectroscopic follow-up. A pair of spectrographs fed by optical fibers measured spectra for more than 600 galaxies, quasars and stars in each single observation. These spectra have wavelength coverage of 3800–9200 Å and a spectral resolving power of $R \sim 2000$. Data Release 7 includes about 1.6 million spectra, with about 900,000 galaxies, 120,000 quasars and 460,000 stars. The total volume of imaging and spectroscopic data products in the SDSS Data Release 7 is about 60 TB.

The second phase of the SDSS included many observations of the same patch of sky, dubbed “Stripe 82.” This opens up a new dimension of astronomical data: the time domain. The Stripe 82 data have led to advances in the understanding of many time-varying phenomena, from asteroid orbits to variable stars to quasars and supernovas. The multiple observations have also been combined to provide a catalog of nonvarying stars with excellent photometric precision.

In addition to providing an unprecedented data set, the SDSS has revolutionized the public dissemination of astronomical data by providing exquisite portals for easy data access, search, analysis, and download. For professional purposes, the Catalog Archive Server (CAS²⁴) and its SQL-based search engine is the most efficient way to get SDSS data. While detailed discussion of SQL is beyond the scope of this book,²⁵ we note that the SDSS site provides a very useful set of example queries²⁶ which can be quickly adapted to other problems.

Alongside the SDSS data, we also provide the Two Micron All Sky Survey (2MASS) photometry for stars from the SDSS Standard Star Catalog, described in [19]. 2MASS [32] used two 1.3 m telescopes to survey the entire sky in near-infrared light. The three 2MASS bands, spanning the wavelength range 1.2–2.2 μm (adjacent

²²<http://www.sdss.org>

²³<http://www.sdss.org/dr7/>

²⁴<http://cas.sdss.org/astrodr7/en/tools/search/sql.asp>

²⁵There are many available books about SQL since it is heavily used in industry and commerce. *Sams Teach Yourself SQL in 10 Minutes* by Forta (Sams Publishing) is a good start, although it took us more than 10 minutes to learn SQL; a more complete reference is *SQL in a Nutshell* by Kline, Kline, and Hunt (O'Reilly), and *The Art of SQL* by Faroult and Robson (O'Reilly) is a good choice for those already familiar with SQL.

²⁶<http://cas.sdss.org/astrodr7/en/help/docs/realquery.asp>

to the SDSS wavelength range on the red side), are called J , H , and K_s (“ s ” in K_s stands for “short”).

We provide several other data sets in addition to SDSS and 2MASS: the LINEAR database features time-domain observations of thousands of variable stars; the LIGO “Big Dog” data²⁷ is a *simulated* data set from a gravitational wave observatory; and the asteroid data file includes orbital data that come from a large variety of sources. For more details about these samples, see the detailed sections below.

We first describe tools and data sets for accessing SDSS imaging data for an arbitrary patch of sky, and for downloading an arbitrary SDSS spectrum. Several data sets specialized for the purposes of this book are described next and include galaxies with SDSS spectra, quasars with SDSS spectra, stars with SDSS spectra, a high-precision photometric catalog of SDSS standard stars, and a catalog of asteroids with known orbits and SDSS measurements.

Throughout the book, these data are supplemented by simulated data ranging from simple one-dimensional toy models to more accurate multidimensional representations of real data sets. The example code for each figure can be used to quickly reproduce these simulated data sets.

1.5.3. SDSS Imaging Data

The total volume of SDSS imaging data is measured in tens of terabytes and thus we will limit our example to a small (20 deg^2 , or 0.05% of the sky) patch of sky. Data for a different patch size, or a different direction on the sky, can be easily obtained by minor modifications of the SQL query listed below.

We used the following SQL query (fully reprinted here to illustrate SDSS SQL queries) to assemble a catalog of $\sim 330,000$ sources detected in SDSS images in the region bounded by $0^\circ < \alpha < 10^\circ$ and $-1^\circ < \delta < 1^\circ$ (α and δ are equatorial sky coordinates called the right ascension and declination).

```
SELECT
    round(p.ra,6) as ra, round(p.dec,6) as dec,
    p.run, --- comments are preceded by ---
    round(p.extinction_r,3) as rExtSFD, --- r band extinction from SFD
    round(p.modelMag_u,3) as uRaw,      --- ISM-uncorrected model mags
    round(p.modelMag_g,3) as gRaw,      --- rounding up model magnitudes
    round(p.modelMag_r,3) as rRaw,
    round(p.modelMag_i,3) as iRaw,
    round(p.modelMag_z,3) as zRaw,
    round(p.modelMagErr_u,3) as uErr,   --- errors are important!
    round(p.modelMagErr_g,3) as gErr,
    round(p.modelMagErr_r,3) as rErr,
    round(p.modelMagErr_i,3) as iErr,
    round(p.modelMagErr_z,3) as zErr,
    round(p.psfMag_u,3) as uRawPSF,    --- psf magnitudes
    round(p.psfMag_g,3) as gRawPSF,
    round(p.psfMag_r,3) as rRawPSF,
    round(p.psfMag_i,3) as iRawPSF,
    round(p.psfMag_z,3) as zRawPSF,
    round(p.psfMagErr_u,3) as upsfErr,
```

²⁷See <http://www.ligo.org/science/GW100916/>

```

round(p.psfMagErr_g,3) as gpsfErr,
round(p.psfMagErr_r,3) as rpsfErr,
round(p.psfMagErr_i,3) as ipsfErr,
round(p.psfMagErr_z,3) as zpsfErr,
p.type, --- tells if a source is resolved or not
(case when (p.flags & '16') = 0 then 1 else 0 end) as ISOLATED --- useful
INTO mydb.SDSSimagingSample
FROM PhotoTag p
WHERE
p.ra > 0.0 and p.ra < 10.0 and p.dec > -1 and p.dec < 1 --- 10x2 sq.deg.
and (p.type = 3 OR p.type = 6) and --- resolved and unresolved sources
(p.flags & '4295229440') = 0 and --- '4295229440' is magic code for no
--- DEBLENDED_AS_MOVING or SATURATED objects
p.mode = 1 and --- PRIMARY objects only, which implies
--- !BRIGHT && (!BLENDDED || NODEBLEND || nchild == 0)]
p.modelMag_r < 22.5 --- adopted faint limit (same as about SDSS limit)
-- the end of query

```

This query can be copied verbatim into the SQL window at the CASJobs site²⁸ (the CASJobs tool is designed for jobs that can require long execution time and requires registration). After running it, you should have your own database called SDSSimagingSample available for download.

The above query selects objects from the PhotoTag table (which includes a subset of the most popular data columns from the main table PhotoObjAll). Detailed descriptions of all listed parameters in all the available tables can be found at the CAS site.²⁹ The subset of PhotoTag parameters returned by the above query includes positions, interstellar dust extinction in the *r* band (from [28]), and the five SDSS magnitudes with errors in two flavors. There are several types of magnitudes measured by SDSS (using different aperture weighting schemes) and the so-called model magnitudes work well for both unresolved (*type*=6, mostly stars and quasars) and resolved (*type*=3, mostly galaxies) sources. Nevertheless, the query also downloads the so-called psf (point spread function) magnitudes. For unresolved sources, the model and psf magnitudes are calibrated to be on average equal, while for resolved sources, model magnitudes are brighter (because the weighting profile is fit to the observed profile of a source and thus can be much wider than the psf, resulting in more contribution to the total flux than in the case of psf-based weights from the outer parts of the source). Therefore, the difference between psf and model magnitudes can be used to recognize resolved sources (indeed, this is the gist of the standard SDSS “star/galaxy” separator whose classification is reported as *type* in the above query). More details about various magnitude types, as well as other algorithmic and processing details, can be found at the SDSS site.³⁰

The WHERE clause first limits the returned data to a 20 deg² patch of sky, and then uses several conditions to select unique stationary and well-measured sources above the chosen faint limit. The most mysterious part of this query is the use of processing flags. These 64-bit flags³¹ are set by the SDSS photometric processing

²⁸<http://casjobs.sdss.org/CasJobs/>

²⁹See *Schema Browser* at <http://skyserver.sdss3.org/dr8/en/help/browser/browser.asp>

³⁰<http://www.sdss.org/dr7/algorithms/index.html>

³¹[http://www.sdss.org/dr7/products/catalogs flags.html](http://www.sdss.org/dr7/products/catalogs	flags.html)

pipeline *photo* [24] and indicate the status of each object, warn of possible problems with the image itself, and warn of possible problems in the measurement of various quantities associated with the object. The use of these flags is unavoidable when selecting a data set with reliable measurements.

To facilitate use of this data set, we have provided code in `astroML.datasets` to download and parse this data. To do this, you must import the function `fetch_imaging_sample`:³²

```
In [1]: from astroML.datasets import \
          fetch_imaging_sample
In [2]: data = fetch_imaging_sample()
```

The first time this is called, the code will send an http request and download the data from the web. On subsequent calls, it will be loaded from local disk. The object returned is a *record array*, which is a data structure within NumPy designed for labeled data. Let us explore these data a bit:

```
In [3]: data.shape
Out[3]: (330753,)
```

We see that there are just over 330,000 objects in the data set. The names for each of the attributes of these objects are stored within the array data type, which can be accessed via the `dtype` attribute of `data`. The names of the columns can be accessed as follows:

```
In [4]: data.dtype.names[:5]
Out[4]: ('ra', 'dec', 'run', 'rExtSFD', 'uRaw')
```

We have printed only the first five names here using the array slice syntax `[:5]`. The data within each column can be accessed via the column name:

```
In [5]: data['ra'][:5]
Out[5]: array([ 0.358174,  0.358382,  0.357898,
               0.35791 ,  0.358881])

In [6]: data['dec'][:5]
Out[6]: array([-0.508718, -0.551157, -0.570892,
               -0.426526, -0.505625])
```

Here we have printed the right ascension and declination (i.e., angular position on the sky) of the first five objects in the catalog. Utilizing Python's plotting package Matplotlib, we show a simple scatter plot of the colors and magnitudes of the first 5000 galaxies and the first 5000 stars from this sample. The result can be seen in

³²Here and throughout we will assume the reader is using the IPython interface, which enables clean interactive plotting with Matplotlib. For more information, refer to appendix A.

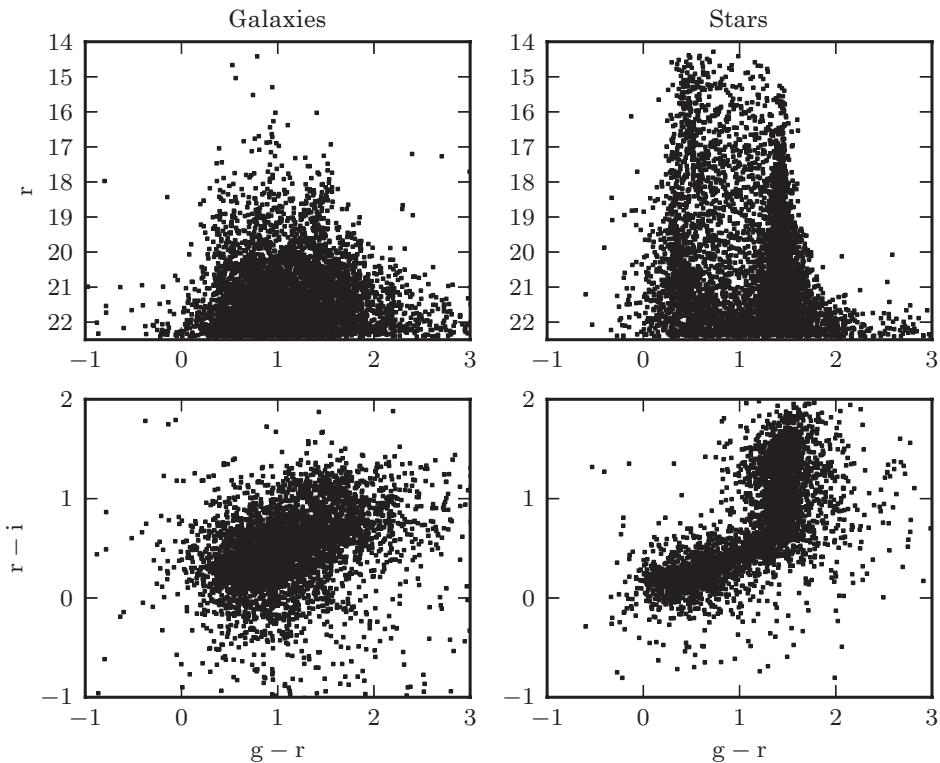


Figure 1.1. The r vs. $g - r$ color-magnitude diagrams and the $r - i$ vs. $g - r$ color-color diagrams for galaxies (left column) and stars (right column) from the SDSS imaging catalog. Only the first 5000 entries for each subset are shown in order to minimize the blending of points (various more sophisticated visualization methods are discussed in §1.6). This figure, and all the others in this book, can be easily reproduced using the `astroML` code freely downloadable from the supporting website.

figure 1.1. Note that as with all figures in this text, the Python code used to generate the figure can be viewed and downloaded on the book website.

Figure 1.1 suffers from a significant shortcoming: even with only 5000 points shown, the points blend together and obscure the details of the underlying structure. This blending becomes even worse when the full sample of 330,753 points is shown. Various visualization methods for alleviating this problem are discussed in §1.6. For the remainder of this section, we simply use relatively small samples to demonstrate how to access and plot data in the provided data sets.

1.5.4. Fetching and Displaying SDSS Spectra

While the above imaging data set has been downloaded in advance due to its size, it is also possible to access the SDSS database directly and in real time. In `astroML.datasets`, the function `fetch_sdss_spectrum` provides an interface to the FITS (Flexible Image Transport System; a standard file format in astronomy for manipulating images and tables³³) files located on the SDSS spectral server. This

³³See <http://fits.gsfc.nasa.gov/iaufwg/iaufwg.html>

operation is done in the background using the built-in Python module `urllib2`. For details on how this is accomplished, see the source code of `fetch_sdss_spectrum`.

The interface is very similar to those from other examples discussed in this chapter, except that in this case the function call must specify the parameters that uniquely identify an SDSS spectrum: the spectroscopic plate number, the fiber number on a given plate, and the date of observation (modified Julian date, abbreviated `mjd`). The returned object is a custom class which wraps the `pyfits` interface to the FITS data file.

```
In [1]: %pylab
Welcome to pylab, a matplotlib-based Python
environment [backend: TkAgg].
For more information, type 'help(pylab)'.

In [2]: from astroML.datasets import\
         fetch_sdss_spectrum
In [3]: plate = 1615 # plate number of the spectrum
In [4]: mjd = 53166 # modified Julian date
In [5]: fiber = 513 # fiber ID number on a given
         plate
In [6]: data = fetch_sdss_spectrum(plate, mjd, fiber)
In [7]: ax = plt.axes()
In [8]: ax.plot(data.wavelength(), data.spectrum,
         '-k')
In [9]: ax.set_xlabel(r'$\lambda (\text{\AA})$')
In [10]: ax.set_ylabel('Flux')
```

The resulting figure is shown in figure 1.2. Once the spectral data are loaded into Python, any desired postprocessing can be performed locally.

There is also a tool for determining the plate, `mjd`, and fiber numbers of spectra in a basic query. Here is an example, based on the spectroscopic galaxy data set described below.

```
In [1]: from astroML.datasets import tools
In [2]: target = tools.TARGET_GALAXY
# main galaxy sample
In [3]: plt, mjd, fib = tools.query_plate_mjd_fiber
        (5, primtarget=target)
In [4]: plt
Out[4]: array([266, 266, 266, 266, 266])

In [5]: mjd
Out[5]: array([51630, 51630, 51630, 51630, 51630])

In [6]: fib
Out[6]: array([27, 28, 30, 33, 35])
```

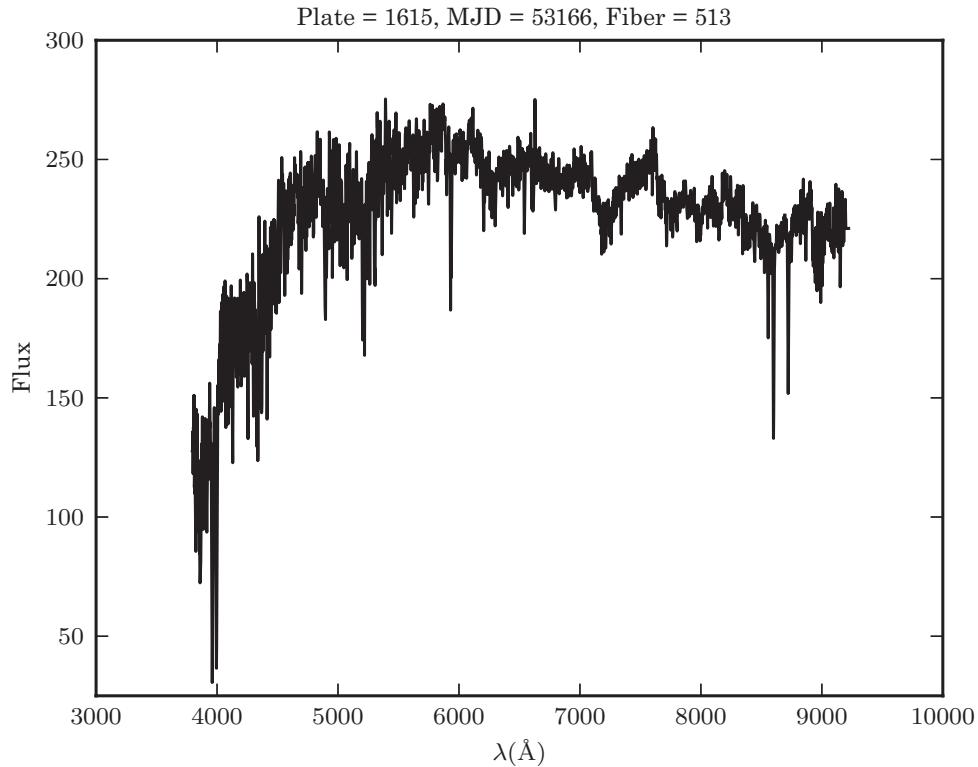


Figure 1.2. An example of an SDSS spectrum (the specific flux plotted as a function of wavelength) loaded from the SDSS SQL server in real time using Python tools provided here (this spectrum is uniquely described by SDSS parameters plate=1615, fiber=513, and mjd=53166).

Here we have asked for five objects, and received a list of five IDs. These could then be passed to the `fetch_sdss_spectrum` function to download and work with the spectral data directly. This function works by constructing a fairly simple SQL query and using `urllib` to send this query to the SDSS database, parsing the results into a NumPy array. It is provided as a simple example of the way SQL queries can be used with the SDSS database.

The plate and fiber numbers and mjd are listed in the next three data sets that are based on various SDSS spectroscopic samples. The corresponding spectra can be downloaded using `fetch_sdss_spectrum`, and processed as desired. An example of this can be found in the script `examples/datasets/compute_sdss_pca.py` within the `astroML` source code tree, which uses spectra to construct the spectral data set used in chapter 7.

1.5.5. Galaxies with SDSS Spectroscopic Data

During the main phase of the SDSS survey, the imaging data were used to select about a million galaxies for spectroscopic follow-up, including the main flux-limited sample (approximately $r < 18$; see the top-left panel in figure 1.1) and a smaller color-selected sample designed to include very luminous and distant galaxies (the

so-called giant elliptical galaxies). Details about the selection of the galaxies for the spectroscopic follow-up can be found in [36].

In addition to parameters computed by the SDSS processing pipeline, such as redshift and emission-line strengths, a number of groups have developed post-processing algorithms and produced so-called “value-added” catalogs with additional scientifically interesting parameters, such as star-formation rate and stellar mass estimates. We have downloaded a catalog with some of the most interesting parameters for $\sim 660,000$ galaxies using the query listed in appendix D submitted to the SDSS Data Release 8 database.

To facilitate use of this data set, in the AstroML package we have included a data set loading routine, which can be used as follows:

```
In [1]: from astroML.datasets import \
         fetch_sdss_specgals
In [2]: data = fetch_sdss_specgals()
In [3]: data.shape
Out[3]: (661598,)

In [4]: data.dtype.names[:5]
Out[4]: ('ra', 'dec', 'mjd', 'plate', 'fiberID')
```

As above, the resulting data is stored in a NumPy record array. We can use the data for the first 10,000 entries to create an example color-magnitude diagram, shown in figure 1.3.

```
In [5]: data = data[:10000] # truncate data
In [6]: u = data['modelMag_u']
In [7]: r = data['modelMag_r']
In [8]: rPetro = data['petroMag_r']
In [9]: %pylab
Welcome to pylab, a matplotlib-based Python
environment [backend: TkAgg].
For more information, type 'help(pylab)'.

In [10]: ax = plt.axes()
In [11]: ax.scatter(u-r, rPetro, s=4, lw=0, c='k')
In [12]: ax.set_xlim(1, 4.5)
In [13]: ax.set_ylim(18.1, 13.5)
In [14]: ax.set_xlabel('$u - r$')
In [15]: ax.set_ylabel('$r_{\mathrm{petrosian}}$')
```

Note that we used the Petrosian magnitudes for the magnitude axis and model magnitudes to construct the $u - r$ color; see [36] for details. Through squinted eyes, one can just make out a division at $u - r \approx 2.3$ between two classes of objects (see [2, 35] for an astrophysical discussion). Using the methods discussed in later chapters, we will be able to automate and quantify this sort of rough by-eye binary classification.

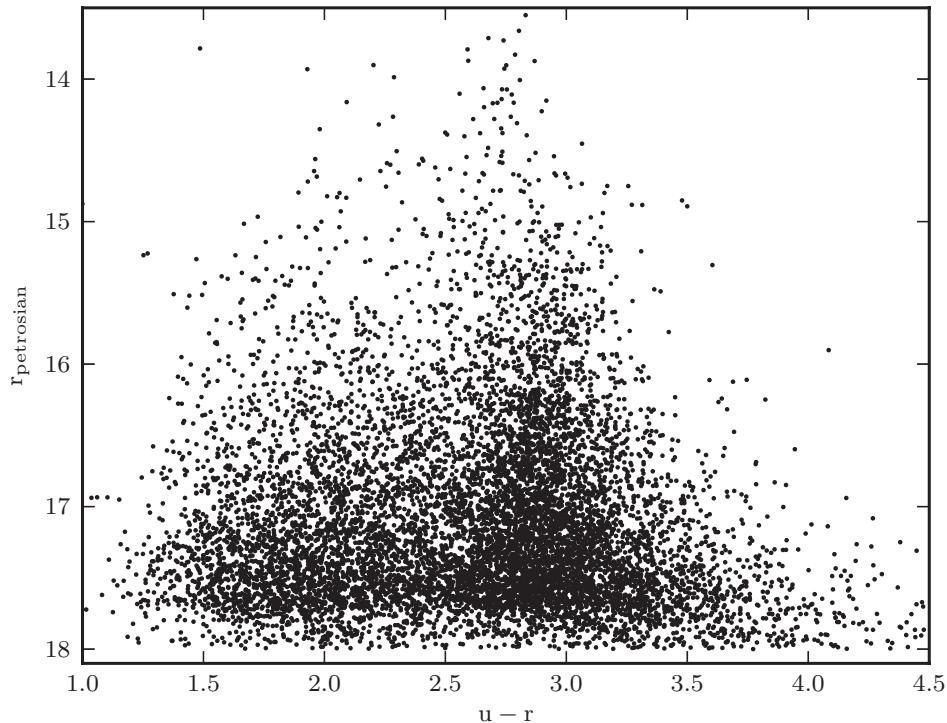


Figure 1.3. The r vs. $u - r$ color-magnitude diagram for the first 10,000 entries in the catalog of spectroscopically observed galaxies from the Sloan Digital Sky Survey (SDSS). Note two “clouds” of points with different morphologies separated by $u - r \approx 2.3$. The abrupt decrease of the point density for $r > 17.7$ (the bottom of the diagram) is due to the selection function for the spectroscopic galaxy sample from SDSS.

1.5.6. SDSS DR7 Quasar Catalog

The SDSS Data Release 7 (DR7) Quasar Catalog contains 105,783 spectroscopically confirmed quasars with highly reliable redshifts, and represents the largest available data set of its type. The construction and content of this catalog are described in detail in [29].

The function `astroML.datasets.fetch_dr7_quasar()` can be used to fetch these data as follows:

```
In [1]: from astroML.datasets import fetch_dr7_quasar
In [2]: data = fetch_dr7_quasar()
In [3]: data.shape
Out[3]: (105783,)

In [4]: data.dtype.names[:5]
Out[4]: ('sdssID', 'RA', 'dec', 'redshift', 'mag_u')
```

One interesting feature of quasars is the redshift dependence of their photometric colors. We can visualize this for the first 10,000 points in the data set as follows:

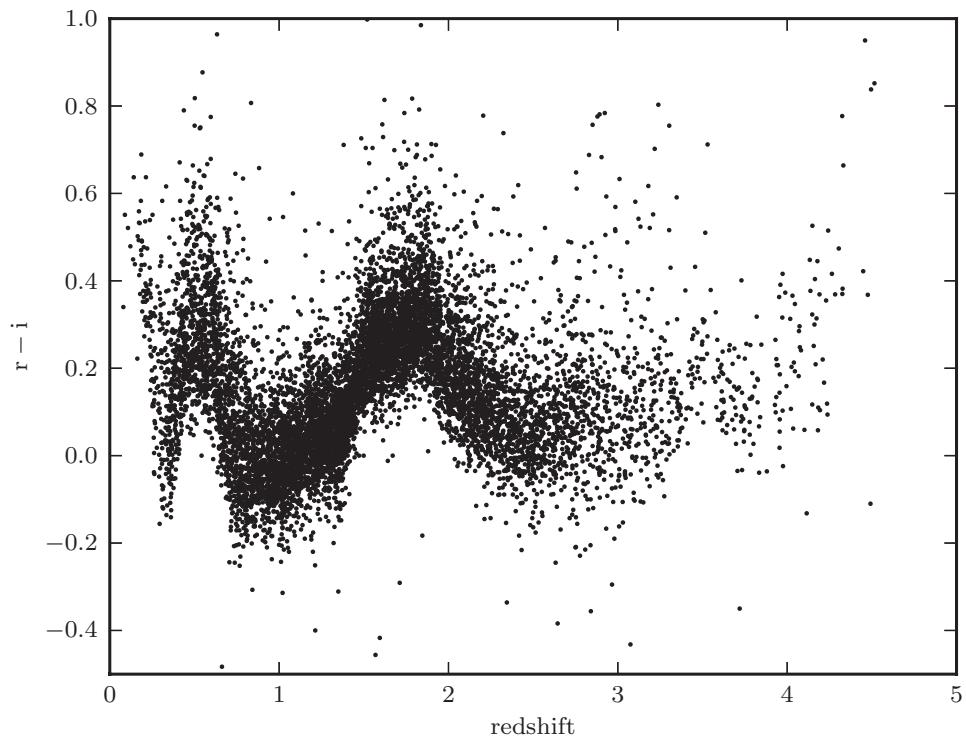


Figure 1.4. The $r - i$ color vs. redshift diagram for the first 10,000 entries from the SDSS Data Release 7 Quasar Catalog. The color variation is due to emission lines entering and exiting the r and i band wavelength windows.

```
In [5]: data = data[:10000]
In [6]: r = data['mag_r']
In [7]: i = data['mag_i']
In [8]: z = data['redshift']
In [9]: %pylab
Welcome to pylab, a matplotlib-based Python
environment [backend: TkAgg].
For more information, type 'help(pylab)'.

In [10]: ax = plt.axes()
In [11]: ax.scatter(z, r - i, s=4, c='black',
 linewidth=0)
In [12]: ax.set_xlim(0, 5)
In [13]: ax.set_ylim(-0.5, 1.0)
In [14]: ax.set_xlabel('redshift')
In [15]: ax.set_ylabel('r-i')
```

Figure 1.4 shows the resulting plot. The very clear structure in this diagram (and analogous diagrams for other colors) enables various algorithms for the photometric estimation of quasar redshifts, a type of problem discussed in detail in chapters 8–9.

1.5.7. SEGUE Stellar Parameters Pipeline Parameters

SDSS stellar spectra are of sufficient quality to provide robust and accurate values of the main stellar parameters, such as effective temperature, surface gravity, and metallicity (parametrized as [Fe/H]; this is the base 10 logarithm of the ratio of abundance of Fe atoms relative to H atoms, itself normalized by the corresponding ratio measured for the Sun, which is ~ 0.02 ; i.e., [Fe/H]=0 for the Sun). These parameters are estimated using a variety of methods implemented in an automated pipeline called SSPP (SEGUE Stellar Parameters Pipeline); a detailed discussion of these methods and their performance can be found in [5] and references therein.

We have selected a subset of stars for which, in addition to [Fe/H], another measure of chemical composition, [α /Fe] (for details see [21]), is also available from SDSS Data Release 9. Note that Data Release 9 is the first release with publicly available [α /Fe] data. These measurements meaningfully increase the dimensionality of the available parameter space; together with the three spatial coordinates and the three velocity components (the radial component is measured from spectra, and the two tangential components from angular displacements on the sky called proper motion), the resulting space has eight dimensions. To ensure a clean sample, we have selected $\sim 330,000$ stars from this catalog by applying various selection criteria that can be found in the documentation for function `fetch_sdss_sspp`.

The data set loader `fetch_sdss_sspp` for this catalog can be used as follows:

```
In [1]: from astroML.datasets import fetch_sdss_sspp
In [2]: data = fetch_sdss_sspp()
In [3]: data.shape
Out[3]: (327260,)
In [4]: data.dtype.names[:5]
Out[4]: ('ra', 'dec', 'Ar', 'rpsf', 'uErr')
```

As above, we use a simple example plot to show how to work with the data. Astronomers often look at a plot of surface gravity vs. effective temperature because it is related to the famous luminosity vs. temperature Hertzsprung–Russell diagram which summarizes well the theories of stellar structure. The surface gravity is typically expressed in the cgs system (in units of cm/s^2), and its logarithm is used in analysis (for orientation, $\log g$ for the Sun is ~ 4.44). As before, we plot only the first 10,000 entries, shown in figure 1.5.

```
In [5]: data = data[:10000]
In [6]: rpsf = data['rpsf'] # make some reasonable
       # cuts
In [7]: data = data[(rpsf > 15) & (rpsf < 19)]
In [8]: logg = data['logg']
In [9]: Teff = data['Teff']
In [10]: %pylab
Welcome to pylab, a matplotlib-based Python
environment [backend: TkAgg].
For more information, type 'help(pylab)'.
```

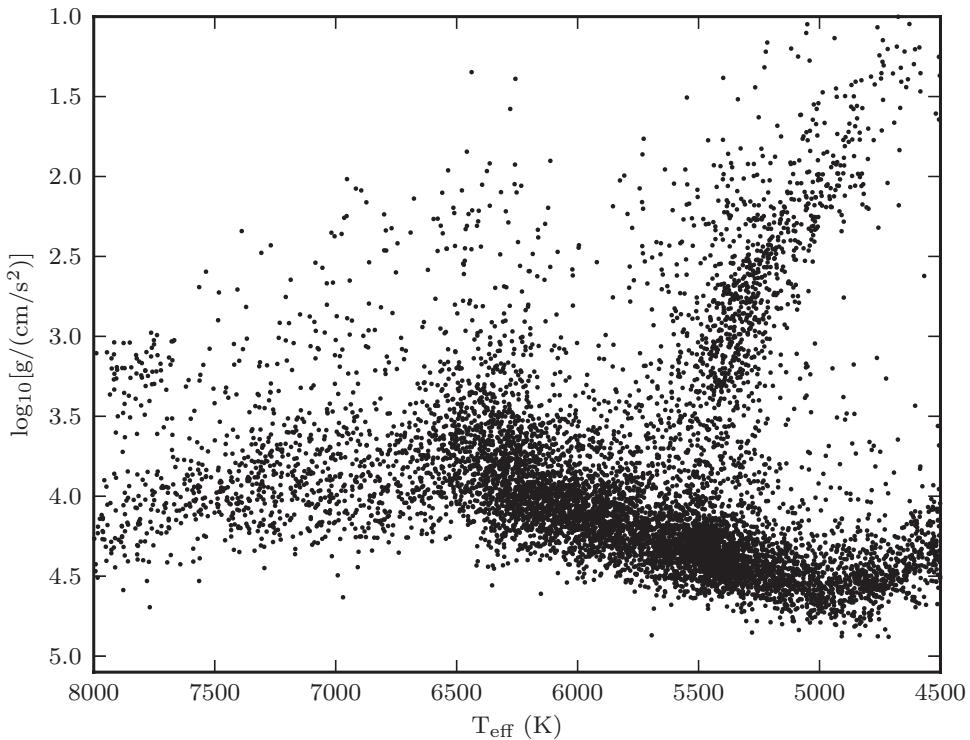


Figure 1.5. The surface gravity vs. effective temperature plot for the first 10,000 entries from the catalog of stars with SDSS spectra. The rich substructure reflects both stellar physics and the SDSS selection criteria for spectroscopic follow-up. The plume of points centered on $T_{\text{eff}} \sim 5300$ K and $\log g \sim 3$ is dominated by red giant stars, and the locus of points with $T_{\text{eff}} < 6500$ K and $\log g > 4.5$ is dominated by main sequence stars. Stars to the left from the main sequence locus are dominated by the so-called blue horizontal branch stars. The axes are plotted backward for ease of comparison with the classical Hertzsprung–Russell diagram: the luminosity of a star approximately increases upward in this diagram.

```
In [11]: ax = plt.axes()
In [12]: ax.scatter(Teff, logg, s=4, lw=0, c='k')
In [13]: ax.set_xlim(8000, 4500)
In [14]: ax.set_ylim(5.1, 1)
In [15]: ax.set_xlabel(r'$\mathrm{T_{eff}}\ (K)$')
In [16]: ax.set_ylabel(r'$\mathrm{\log_{10}[g / (cm/s^2)]}$')
```

1.5.8. SDSS Standard Star Catalog from Stripe 82

In a much smaller area of ~ 300 deg 2 , SDSS has obtained repeated imaging that enabled the construction of a more precise photometric catalog containing ~ 1 million stars (the precision comes from the averaging of typically over ten observations). These stars were selected as nonvariable point sources and have photometric precision better than 0.01 mag at the bright end (or about twice as good as single measurements). The size and photometric precision of this catalog make it a good choice for exploring various methods described in this book, such as stellar

locus parametrization in the four-dimensional color space, and search for outliers. Further details about the construction of this catalog and its contents can be found in [19].

There are two versions of this catalog available from `astroML.datasets`. Both are accessed with the function `fetch_sdss_S82standards`. The first contains just the attributes measured by SDSS, while the second version includes a subset of stars cross-matched to 2MASS. This second version can be obtained by calling

```
fetch_sdss_S82standards(crossmatch_2mass = True).
```

The following shows how to fetch and plot the data:

```
In [1]: from astroML.datasets import\
         fetch_sdss_S82standards
In [2]: data = fetch_sdss_S82standards()
In [3]: data.shape
Out[3]: (1006849,)
In [4]: data.dtype.names[:5]
Out[4]: ('RA', 'DEC', 'RArms', 'DECrms', 'Ntot')
```

Again, we will create a simple color–color scatter plot of the first 10,000 entries, shown in figure 1.6.

```
In [5]: data = data[:10000]
In [6]: g = data['mmu_g'] # g-band mean magnitude
In [7]: r = data['mmu_r'] # r-band mean magnitude
In [8]: i = data['mmu_i'] # i-band mean magnitude
In [9]: %pylab
Welcome to pylab, a matplotlib-based Python
environment [backend: TkAgg].
For more information, type 'help(pylab)'.

In [10]: ax = plt.axes()
In [11]: ax.scatter(g - r, r - i, s=4, c='black',
                  linewidth=0)
In [12]: ax.set_xlabel('g - r')
In [13]: ax.set_ylabel('r - i')
```

1.5.9. LINEAR Stellar Light Curves

The LINEAR project has been operated by the MIT Lincoln Laboratory since 1998 to discover and track near-Earth asteroids (the so-called “killer asteroids”). Its archive now contains approximately 6 million images of the sky, most of which are 5 MP images covering 2 deg². The LINEAR image archive contains a unique combination of sensitivity, sky coverage, and observational cadence (several hundred observations per object). A shortcoming of original reductions of LINEAR data is that its photometric calibration is fairly inaccurate because the effort was focused on

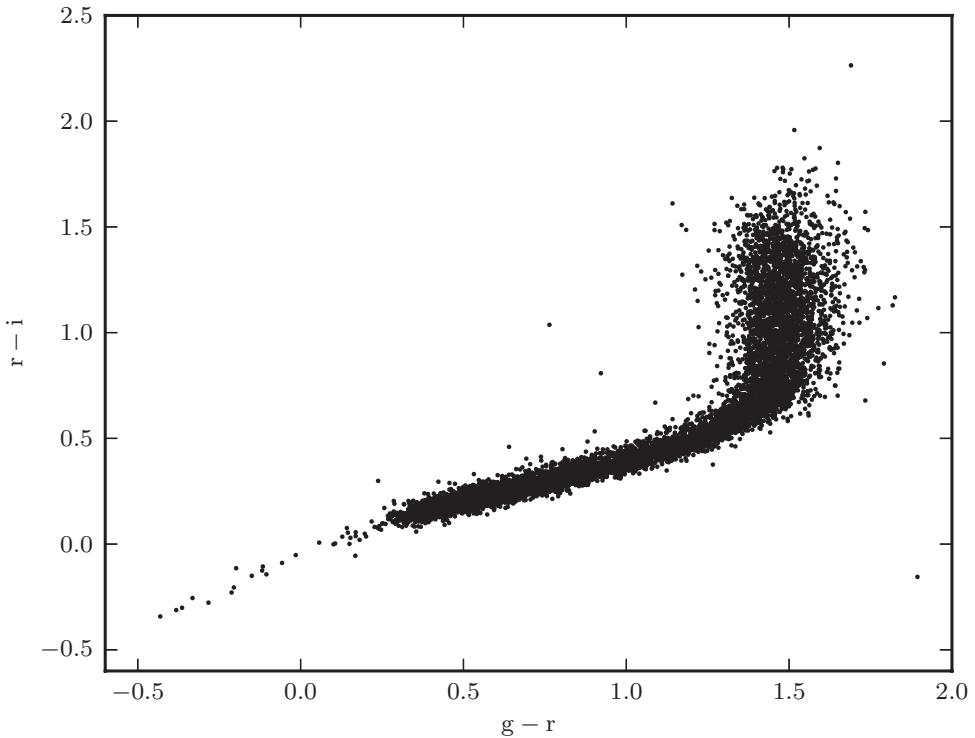


Figure 1.6. The $g - r$ vs. $r - i$ color–color diagram for the first 10,000 entries in the Stripe 82 Standard Star Catalog. The region with the highest point density is dominated by main sequence stars. The thin extension toward the lower-left corner is dominated by the so-called blue horizontal branch stars and white dwarf stars.

astrometric observations of asteroids. Here we use recalibrated LINEAR data from the sky region covered by SDSS which aided recalibration [30]. We focus on 7000 likely periodic variable stars. The full data set with 20 million light curves is publicly available.³⁴

The loader for the LINEAR data set is `fetch_LINEAR_sample`. This data set contains light curves and associated catalog data for over 7000 objects:

```
In [1]: from astroML.datasets import \
         fetch_LINEAR_sample
In [2]: data = fetch_LINEAR_sample()
In [3]: gr = data.targets['gr']    # g-r color
In [4]: ri = data.targets['ri']    # r-i color
In [5]: logP = data.targets['LP1']
# log_10(period) in days
In [6]: gr.shape
Out[6]: (7010,)

In [7]: id = data.ids[2756]      # choose one id from the
       # sample
```

³⁴The LINEAR Survey Photometric Database is available from <https://astroweb.lanl.gov/lineardb/>

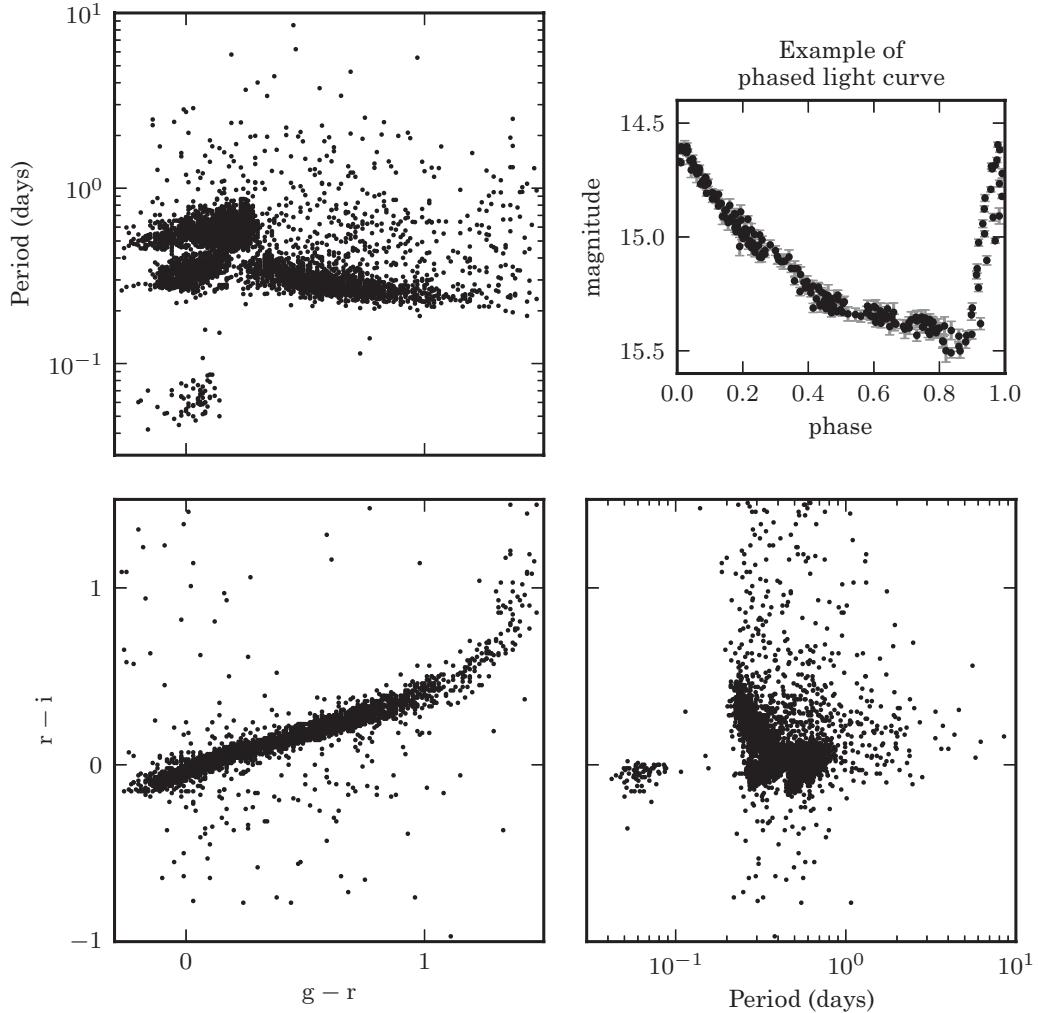


Figure 1.7. An example of the type of data available in the LINEAR data set. The scatter plots show the $g - r$ and $r - i$ colors, and the variability period determined using a Lomb–Scargle periodogram (for details see chapter 10). The upper-right panel shows a phased light curve for one of the over 7000 objects.

```
In [8]: id
Out[8]: 18527462

In [9]: t, mag, dmag = data[id].T
# access light curve data
In [10]: logP = data.get_target_parameter(id, 'LP1')
```

The somewhat cumbersome interface is due to the size of the data set: to avoid the overhead of loading all of the data when only a portion will be needed in any given script, the data are accessed through a class interface which loads the needed data on demand. Figure 1.7 shows a visualization of the data loaded in the example above.

1.5.10. SDSS Moving Object Catalog

SDSS, although primarily designed for observations of extragalactic objects, contributed significantly to studies of Solar system objects. It increased the number of asteroids with accurate five-color photometry by more than a factor of one hundred, and to a flux limit about one hundred times fainter than previous multicolor surveys. SDSS data for asteroids is collated and available as the Moving Object Catalog³⁵ (MOC). The 4th MOC lists astrometric and photometric data for ~472,000 Solar system objects. Of those, ~100,000 are unique objects with known orbital elements obtained by other surveys.

We can use the provided Python utilities to access the MOC data. The loader is called `fetch_moving_objects`.

```
In [1]: from astroML.datasets import \
    fetch_moving_objects
In [2]: data = fetch_moving_objects(Parker2008_cuts=
    True)
In [3]: data.shape
Out[3]: (33160,)

In [4]: data.dtype.names[:5]
Out[4]: ('moID', 'sdss_run', 'sdss_col', 'sdss_field',
    'sdss_obj')
```

As an example, we make a scatter plot of the orbital semimajor axis vs. the orbital inclination angle for the first 10,000 catalog entries (figure 1.8). Note that we have set a flag to make the data quality cuts used in [26] to increase the measurement quality for the resulting subsample. Additional details about this plot can be found in the same reference, and references therein.

```
In [5]: data = data[:10000]
In [6]: a = data['aprime']
In [7]: sini = data['sin_iprime']
In [8]: %pylab
Welcome to pylab, a matplotlib-based Python
environment [backend: TkAgg].
For more information, type 'help(pylab)'.

In [9]: ax = plt.axes()
In [10]: ax.scatter(a, sini, s=4, c='black',
    linewidth=0)
In [11]: ax.set_xlabel('Semi-major Axis (AU)')
In [12]: ax.set_ylabel('Sine of Inclination Angle')
```

³⁵<http://www.astro.washington.edu/users/ivezic/sdssmoc/sdssmoc.html>

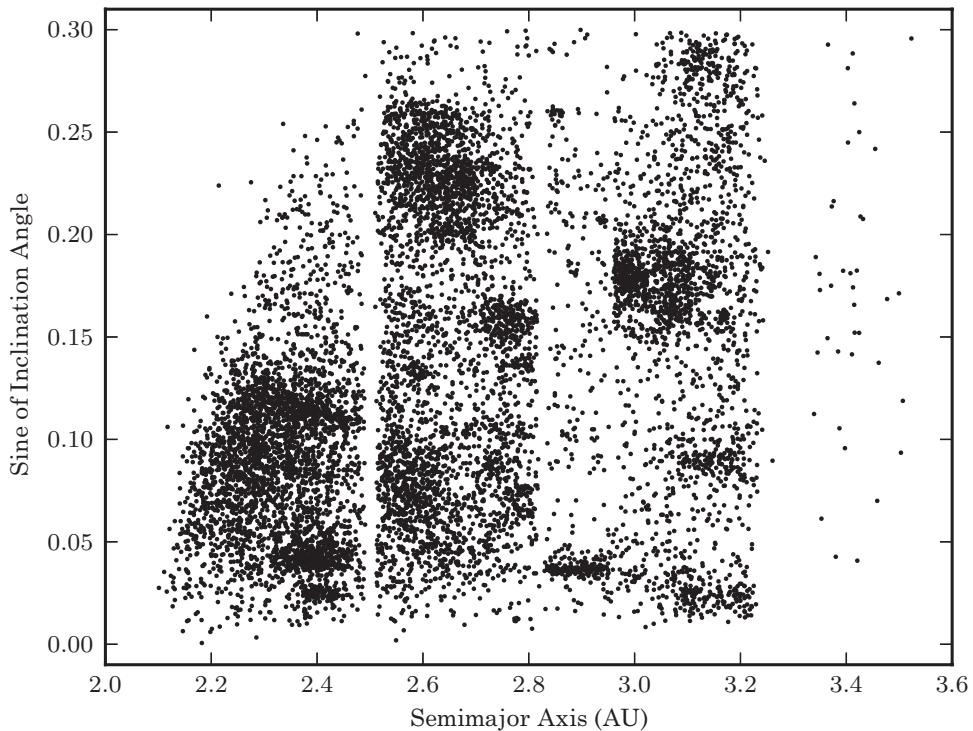


Figure 1.8. The orbital semimajor axis vs. the orbital inclination angle diagram for the first 10,000 catalog entries from the SDSS Moving Object Catalog (after applying several quality cuts). The gaps at approximately 2.5, 2.8, and 3.3 AU are called the Kirkwood gaps and are due to orbital resonances with Jupiter. The several distinct clumps are called asteroid families and represent remnants from collisions of larger asteroids.

1.6. Plotting and Visualizing the Data in This Book

Data visualization is an important part of scientific data analysis, both during exploratory analysis (e.g., to look for problems in data, searching for patterns, and informing quantitative hypothesis) and for the presentation of results. There are a number of books of varying quality written on this topic. An exceptional book is *The Visual Display of Quantitative Information* by Tufte [37], with excellent examples of both good and bad graphics, as well as clearly exposed design principles. Four of his principles that directly pertain to large data sets are (i) present many numbers in a small space, (ii) make large data sets coherent, (iii) reveal the data at several levels of detail, and (iv) encourage the eye to compare different pieces of data. For a recent review of high-dimensional data visualization in astronomy see [11].

1.6.1. Plotting Two-Dimensional Representations of Large Data Sets

The most fundamental quantity we typically want to visualize and understand is the distribution or density of the data. The simplest way to do this is via a scatter plot. When there are too many points to plot, individual points tend to blend together in dense regions of the plot. We must find an effective way to model the density. Note that, as we will see in the case of the histogram (§5.7.2), visualization of the density cannot be done ad hoc, that is, estimating the density is a statistical problem in

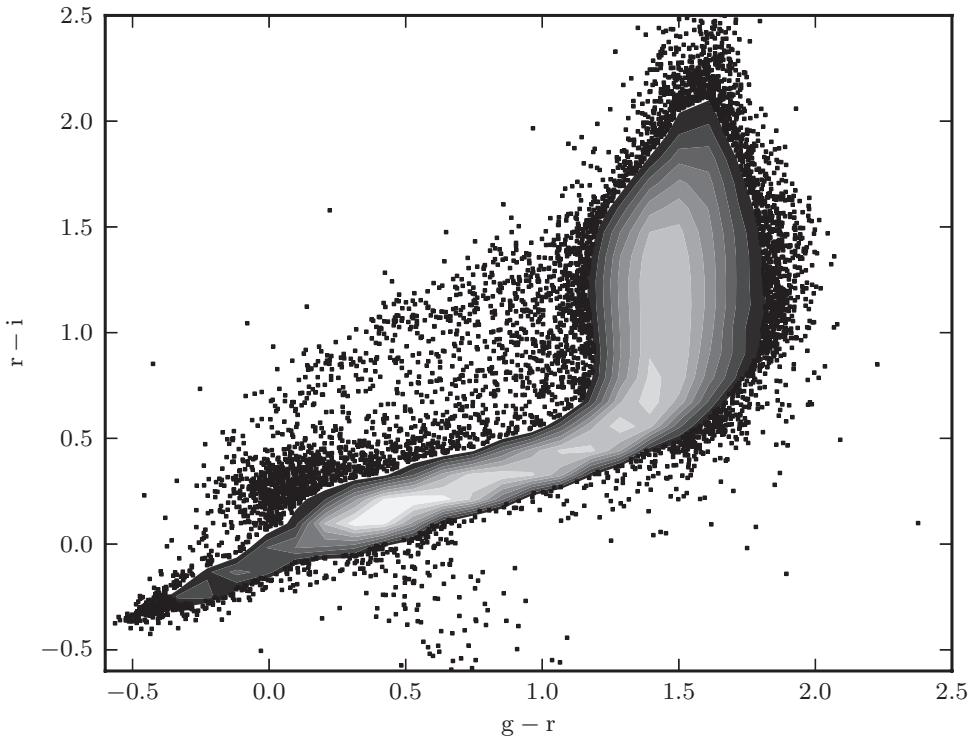


Figure 1.9. Scatter plot with contours over dense regions. This is a color–color diagram of the entire set of SDSS Stripe 82 standard stars; cf. figure 1.6.

itself—choices in simple visualizations of the density may undersmooth or oversmooth the data, misleading the analyst about its properties (density estimation methods are discussed in chapter 6).

A visualization method which addresses this blending limitation is the *contour plot*. Here the contours successfully show the distribution of dense regions, but at the cost of losing information in regions with only a few points. An elegant solution is to use contours for the high-density regions, and show individual points in low-density regions (due to Michael Strauss from Princeton University, who pioneered this approach with SDSS data). An example is shown in figure 1.9 (compare to the scatter plot of a subset of this data in figure 1.6).

Another method is to pixelize the plotted diagram and display the counts of points in each pixel (this “two-dimensional histogram” is known as a *Hess diagram* in astronomy, though this term is often used to refer specifically to color–magnitude plots visualized in this way). The counts can be displayed with different “stretch” (or mapping functions) in order to improve dynamic range (e.g., a logarithmic stretch). A Hess diagram for the color–color plot of the SDSS Stripe 82 standard stars is shown in figure 1.10.

Hess diagrams can be useful in other ways as well. Rather than simply displaying the count or density of points as a function of two parameters, one often desires to show the variation of a separate statistic or measurement. An example of this is shown in figure 1.11. The left panel shows the Hess diagram of the density of points as a function of temperature and surface gravity. The center panel shows a Hess diagram, except here the value in each pixel is the mean metallicity ([Fe/H]).

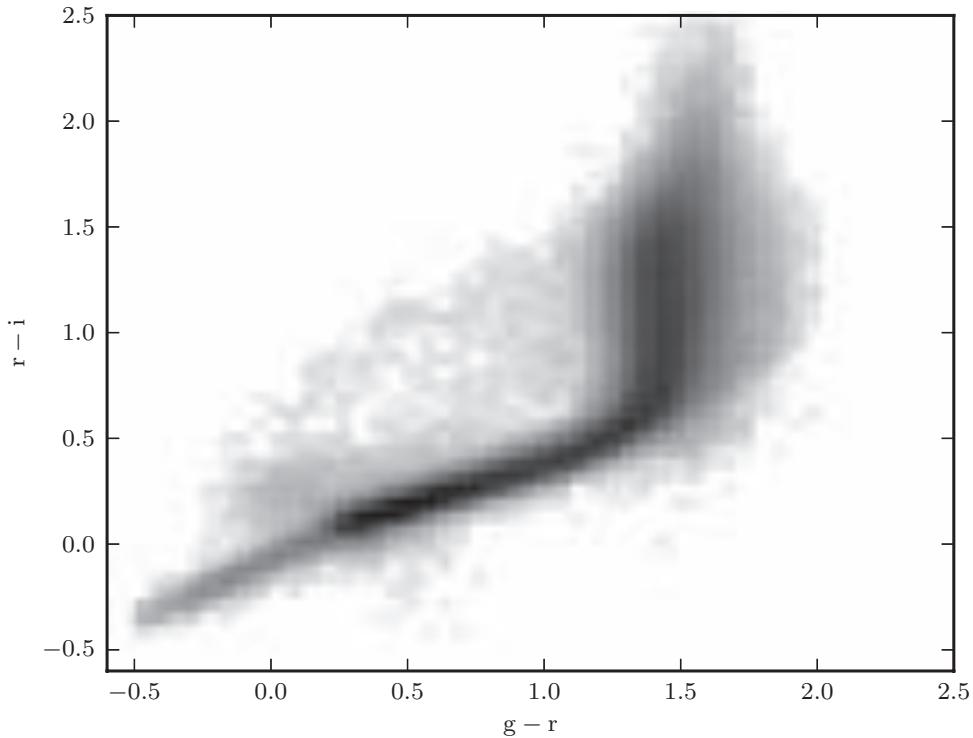


Figure 1.10. A Hess diagram of the $r - i$ vs. $g - r$ colors for the entire set of SDSS Stripe 82 standard stars. The pixels are colored with a logarithmic scaling; cf. figures 1.6 and 1.9.

The number density contours are overplotted for comparison. The grayscale color scheme in the middle panel can lead to the viewer missing fine changes in scale: for this reason, the right panel shows the same plot with a multicolor scale. This is one situation in which a multicolored scale allows better representation of information than a simple grayscale. Combining the counts and mean metallicity into a single plot provides much more information than the individual plots themselves.

Sometimes the quantity of interest is the density variation traced by a sample of points. If the number of points per required resolution element is very large, the simplest method is to use a Hess diagram. However, when points are sparsely sampled, or the density variation is large, it can happen that many pixels have low or vanishing counts. In such cases there are better methods than the Hess diagram where, in low-density regions, we might display a model for the density distribution as discussed, for example, in §6.1.1.

1.6.2. Plotting in Higher Dimensions

In the case of three-dimensional data sets (i.e., three vectors of length N , where N is the number of points), we have already seen examples of using color to encode a third component in a two-dimensional diagram. Sometimes we have four data vectors and would like to find out whether the position in one two-dimensional diagram is correlated with the position in another two-dimensional diagram. For example, we can ask whether two-dimensional color information for asteroids is correlated with their orbital semimajor axis and inclination [18], or whether the color and luminosity of galaxies are correlated with their position in a spectral emission-line diagram [33].

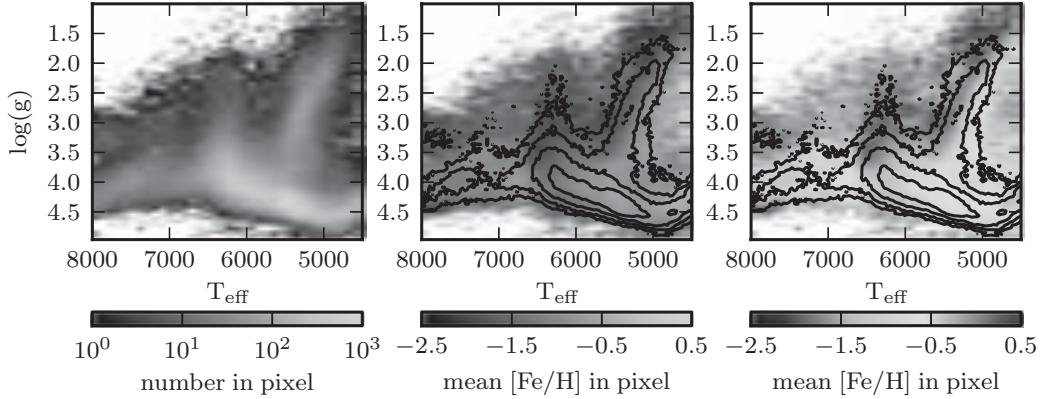


Figure 1.11. A Hess diagram of the number per pixel (left) and [Fe/H] metallicity (center, right) of SEGUE Stellar Parameters Pipeline stars. In the center and right panels, contours representing the number density are overplotted for comparison. These two panels show identical data, but compare a grayscale and multicolor plotting scheme. This is an example of a situation in which multiple colors are very helpful in distinguishing close metallicity levels. This is the same data as shown in figure 1.5. See color plate 1.

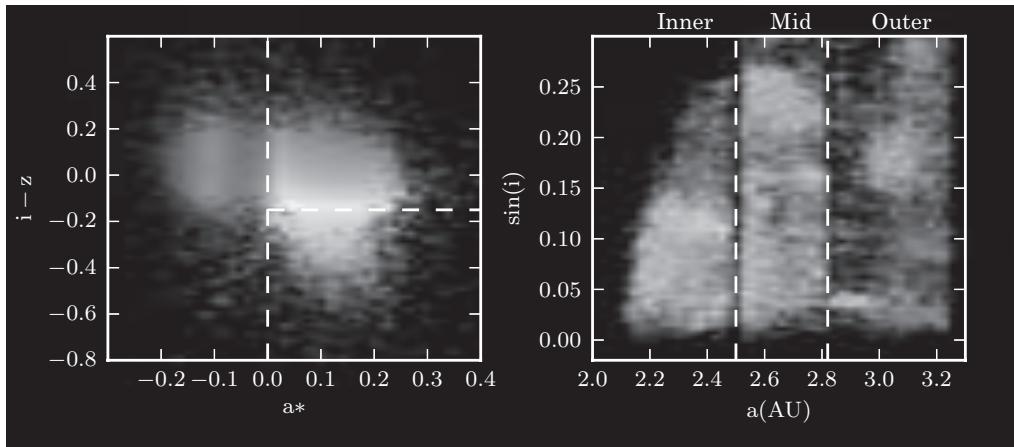


Figure 1.12. A multicolor scatter plot of the properties of asteroids from the SDSS Moving Object Catalog (cf. figure 1.8). The left panel shows observational markers of the chemical properties of the asteroids: two colors a^* and $i - z$. The right panel shows the orbital parameters: semimajor axis a vs. the sine of the inclination. The color of points in the right panel reflects their position in the left panel. See color plate 2.

Let us assume that the four data vectors are called (x, y, z, w) . It is possible to define a continuous two-dimensional color palette that assigns a unique color to each data pair from, say, (z, w) . Then we can plot the $x - y$ diagram with each symbol, or pixel, color coded according to this palette (of course, one would want to show the $z - w$ diagram, too). An example of this visualization method, based on [18], is shown in figure 1.12.

For higher-dimensional data, visualization can be very challenging. One possibility is to seek various low-dimensional projections which preserve certain “interesting” aspects of the data set. Several of these *dimensionality reduction* techniques are discussed in chapter 7.

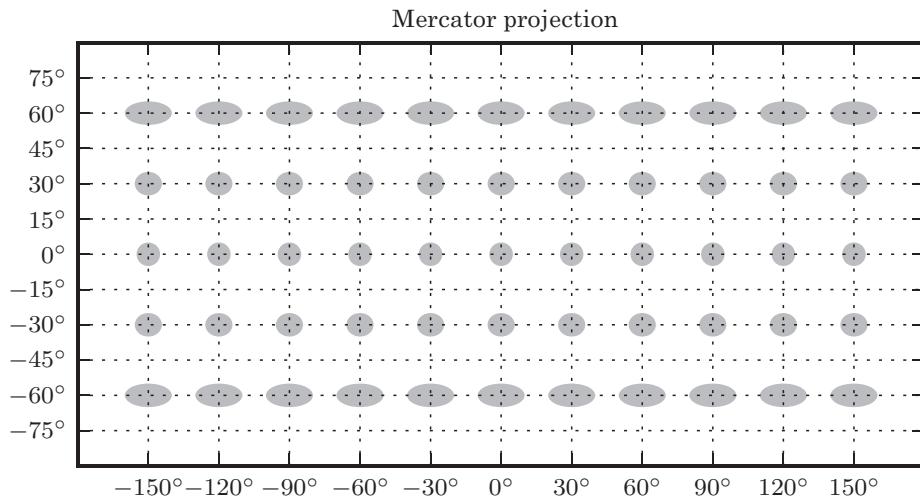


Figure 1.13. The Mercator projection. Shown are the projections of circles of constant radius 10° across the sky. Note that the area is not preserved by the Mercator projection: the projection increases the size of finite regions on the sphere, with a magnitude which increases at high latitudes.

1.6.3. Plotting Representations of Data on the Sky

Plotting the distributions or densities of sources as they would appear on the sky is an integral part of many large-scale analyses (including the analysis of the cosmic microwave background or the angular clustering of galaxies). The projection of various spherical coordinate systems (equatorial, ecliptic, galactic) to a plane is often used in astronomy, geography, and other sciences. There are a few dozen different projections that can be found in the literature, but only a few are widely used. There are always distortions associated with projecting a curved surface onto a plane, and various projections are constructed to preserve different properties (e.g., distance, angle, shape, area).

The Mercator projection is probably the most well known since it was used for several centuries for nautical purposes. The lines of constant true compass bearing (called loxodromes or rhumb lines) are straight line segments in this projection, hence its use in navigation. Unfortunately, it distorts the size of map features. For example, world maps in this projection can be easily recognized by the size of Greenland being about the same as the size of Africa (with the latter being much larger in reality). This can be seen from the sizes of the projected circles (called Tissot's indicatrix) in figure 1.13. Projections that preserve the feature size, known as equal-area projections, are more appropriate for use in astronomy, and here we review and illustrate a few of the most popular choices.

The Hammer and Aitoff projections are visually very similar. The former is an equal-area projection and the latter is an equal-distance projection. Sometimes, the Hammer projection is also referred to as the Hammer–Aitoff projection. They show an entire sphere centered on the equator and rescaled to cover twice as much equatorial distance as polar distance (see figure 1.14). For example, these projections were used for the all-sky maps produced by IRAS (the InfraRed Astronomy Satellite).

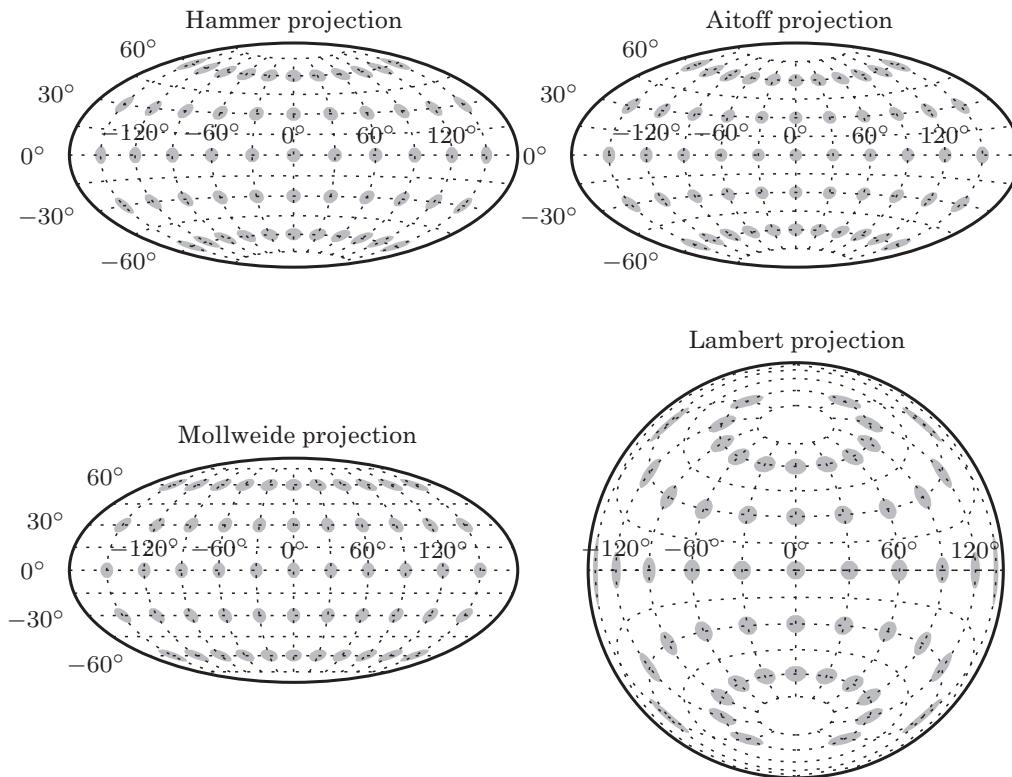


Figure 1.14. Four common full-sky projections. The shaded ellipses represent the distortion across the sky: each is projected from a circle of radius 10° on the sphere. The extent to which these are distorted and/or magnified shows the distortion inherent to the mapping.

The Mollweide projection is another equal-area projection, similar to the Hammer projection, except for straight parallels of latitude instead of the Hammer's curved parallels (developed by an astronomer). It is also known as the Babinet projection, elliptical projection, and homographic (or homalographic) projection. This projection was used to visualize the WMAP (Wilkinson Microwave Anisotropy Probe) maps.

The Lambert azimuthal equal-area projection maps spherical coordinates to a disk. It is especially useful for projecting the two sky hemispheres into two disks.

In general, given spherical coordinates, (α, δ) , the projected planar coordinates, (x, y) , are computed using formulas for a particular projection. For example, for the Hammer projection planar coordinates can be computed from

$$x = \frac{2\sqrt{2} \cos(\delta) \sin(\alpha/2)}{\sqrt{1 + \cos(\delta) \cos(\alpha/2)}} \quad (1.5)$$

and

$$y = \frac{\sqrt{2} \sin(\delta)}{\sqrt{1 + \cos(\delta) \cos(\alpha/2)}}. \quad (1.6)$$

The inverse transformation can be computed as

$$\alpha = 2 \arctan \left[\frac{zx}{2(z^2 - 1)} \right] \quad (1.7)$$

and

$$\delta = \arcsin(zy), \quad (1.8)$$

where

$$z = \sqrt{1 - (x/4)^2 - (y/2)^2}. \quad (1.9)$$

These projections are available in Matplotlib by specifying the `projection` keyword when building the axis. See the source code associated with figure 1.14 for an example of how this can be accomplished in practice. For more obscure projections, the basemap toolkit,³⁶ an add-on to Matplotlib, has a more complete set of utilities. These are primarily geared toward visualization for earth sciences, but can be very useful for astronomical sky projections as well.

A concept related to spherical projection is pixelization of a spherical surface. One of the most useful tools is HEALPix (Hierarchical Equal Area isoLatitude Pixelization). HEALPix subdivides a sphere into equal-area pixels (which are not squares but rather curvilinear quadrilaterals). This tessellation is done hierarchically, with higher levels corresponding to smaller pixels. The lowest resolution partition includes 12 pixels, and each new level divides each pixel into four new ones (see figure 1.15). For example, to reach ~ 3 arcmin resolution, it takes about 12 million pixels. Pixels are distributed on lines of constant latitude, which simplifies and speeds up analysis based on spherical harmonics [12]. The HEALPix code (in IDL and Fortran 90) is publicly available from NASA.³⁷ A Python version, called HealPy is also available.³⁸ The lower panel of figure 1.15 shows an example of raw WMAP data, in a Mollweide projection using data in a HEALPix format.

1.7. How to Efficiently Use This Book

We hope that this book will be found useful both as formal course material, and as a self-study guide and reference book. Sufficient statistical background is provided in chapters 2–5 to enable a semester-long course on astronomical statistics (perhaps with one additional chapter from chapters 6–10 to make the course more data-analysis oriented). On the other hand, chapters 6–10 (together with supporting chapter 1) can enable a semester-long course on data mining and machine learning in astronomy. Unlike most textbooks, we do not provide specific exercises with answers. The main reason is that modern scientific data analysis is intimately intertwined with the writing and execution of efficient computer code, and we have designed this book as a practical text with that fact in mind. If a lecturer prefers problem assignments, we highly recommend exercises from Lup93 for a course on astronomical statistics.

A unique feature of this text is the free availability of example code to fetch relevant data sets and recreate every figure in each chapter of the book. This code

³⁶<http://matplotlib.github.com/basemap/>

³⁷Details about HEALPix are available from <http://healpix.jpl.nasa.gov/>

³⁸<http://healpy.readthedocs.org/en/latest/index.html>

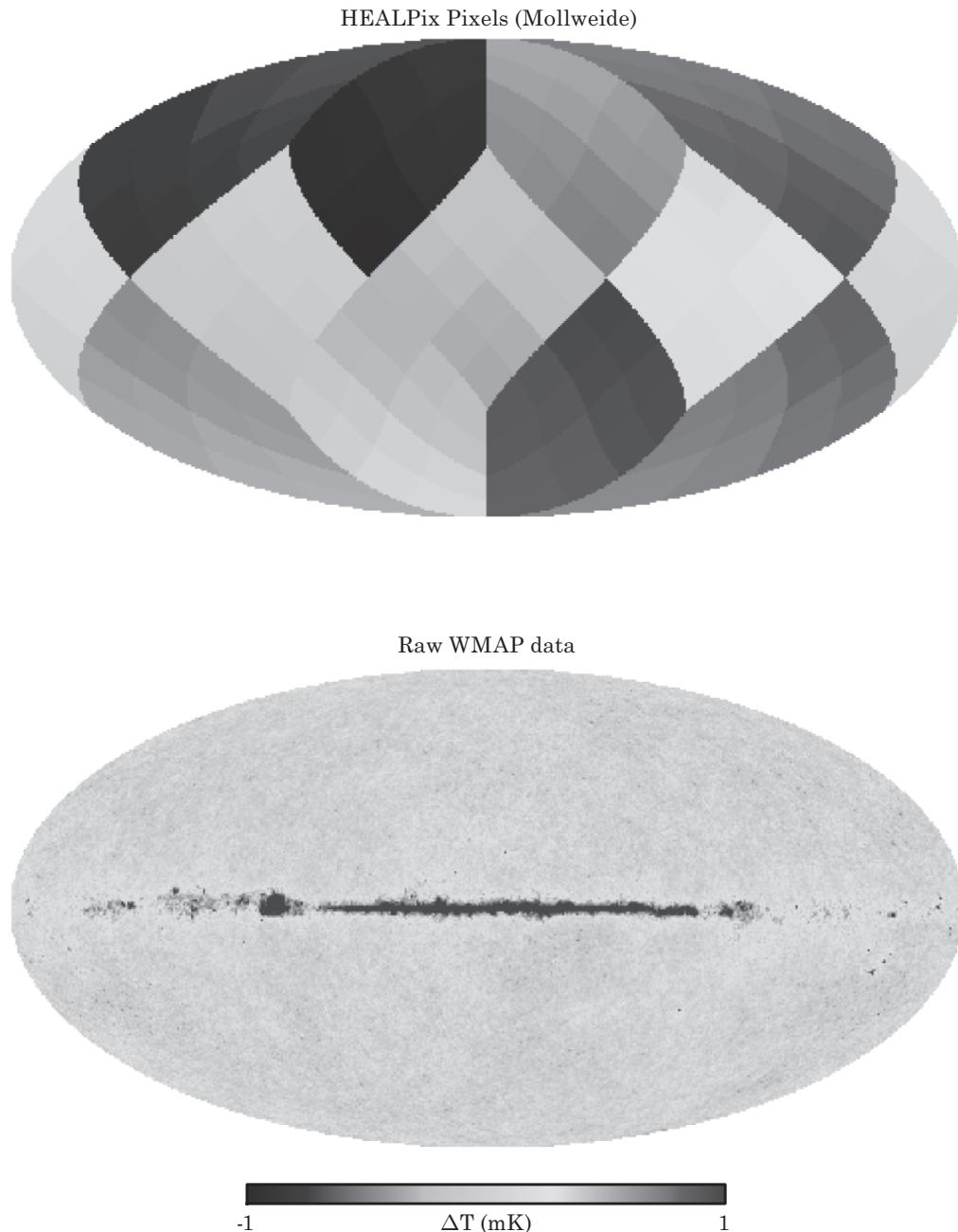


Figure 1.15. The top panel shows HEALPix pixels in nested order. The 12 fundamental sky divisions can be seen, as well as the hierarchical nature of the smaller pixels. This shows a pixelization with $n_{\text{side}} = 4$, that is, each of the 12 large regions has 4×4 pixels, for a total of 192 pixels. The lower panel shows a seven-year co-add of raw WMAP data, plotted using the HEALPix projection using the HealPy package. This particular realization has $n_{\text{side}} = 512$, for a total of 3,145,728 pixels. The pixels are roughly 6.8 arcminutes on a side. See color plate 3.

is available online at <http://www.astroML.org>, where the examples are organized by figure number. Additionally, throughout this text we include minimal code snippets which are meant to give a flavor of how various tools can be used. These snippets are not generally meant to be complete examples; this is the purpose of the online resources.

All code snippets in the book are set aside and appear like this. They will show some minimal code for purposes of illustration. For example, this is how to compute the cosine of a sequence of numbers:

```
import numpy as np
x = np.random.random(100)    # 100 numbers between 0
                             # and 1
cos_x = np.cos(x)          # cosine of each element.
```

For more details on the essential modules for scientific computing in Python, see appendix A.

To take advantage of this book layout, we suggest downloading, examining, modifying, and experimenting with the source code used to create each figure in this text. In order to run these examples on your own machine, you need to install AstroML and its dependencies. A discussion of installation requirements can be found in appendix B, and on the AstroML website.

You can test the success of the installation by plotting one of the example figures from this chapter. For example, to plot figure 1.1, download the source code from http://www.astroml.org/book_figures/chapter1/ and run the code. The data set will be downloaded and the code should generate a plot identical to figure 1.1. You can then modify the code: for example, rather than $g - r$ and $r - i$ colors, you may wish to see the diagram for $u - g$ and $i - z$ colors.

To get the most out of reading this book, we suggest the following interactive approach: When you come across a section which describes a technique or method which interests you, first find the associated figure on the website and copy the source code into a file which you can modify. Experiment with the code: run the code several times, modifying it to explore how variations of the input (e.g., number of points, number of features used or visualized, type of data) affect the results. See if you can find combinations of parameters that improve on the results shown in the book, or highlight the strengths and weaknesses of the method in question. Finally, you can use the code as a template for running a similar method on your own research data.

We hope that this interactive way of reading the text, working with the data firsthand, will give you the experience and insight needed to successfully apply data mining and statistical learning approaches to your own research, whether it is in astronomy or another data-intensive science.

References

- [1] Abazajian, K. N., J. K. Adelman-McCarthy, M. A. Agüeros, and others (2009). The Seventh Data Release of the Sloan Digital Sky Survey. *ApJS* 182, 543–558.
- [2] Baldry, I. K., K. Glazebrook, J. Brinkmann, and others (2004). Quantifying the bimodal color-magnitude distribution of galaxies. *ApJ* 600, 681–694.
- [3] Ball, N. M. and R. J. Brunner (2010). Data mining and machine learning in astronomy. *International Journal of Modern Physics D* 19, 1049–1106.
- [4] Barlow, R. (1989). *Statistics. A Guide to the Use of Statistical Methods in the Physical Sciences*. The Manchester Physics Series, New York: Wiley, 1989.

- [5] Beers, T. C., Y. Lee, T. Sivarani, and others (2006). The SDSS-I Value-Added Catalog of stellar parameters and the SEGUE pipeline. *Mem.S.A.I.* 77, 1171.
- [6] Bishop, C. M. (2006). *Pattern Recognition and Machine Learning*. Springer.
- [7] Borne, K. (2009). Scientific data mining in astronomy. *ArXiv:astro-ph/0911.0505*.
- [8] Borne, K., A. Accomazzi, J. Bloom, and others (2009). Astroinformatics: A 21st century approach to astronomy. In *Astro2010: The Astronomy and Astrophysics Decadal Survey*. *ArXiv:astro-ph/0909.3892*.
- [9] Feigelson, E. D. and G. J. Babu (2012). *Modern Statistical Methods for Astronomy With R Applications*. Cambridge University Press.
- [10] Feigelson, E. D. and G. Jogesh Babu (2012). Statistical methods for astronomy. *ArXiv:astro-ph/1205.2064*.
- [11] Goodman, A. A. (2012). Principles of high-dimensional data visualization in astronomy. *Astronomische Nachrichten* 333, 505.
- [12] Górski, K. M., E. Hivon, A. J. Banday, and others (2005). HEALPix: A framework for high-resolution discretization and fast analysis of data distributed on the sphere. *ApJ* 622, 759–771.
- [13] Gregory, P. C. (2005). *Bayesian Logical Data Analysis for the Physical Sciences: A Comparative Approach with ‘Mathematica’ Support*. Cambridge University Press.
- [14] Gunn, J. E., M. Carr, C. Rockosi, and others (1998). The Sloan Digital Sky Survey photometric camera. *AJ* 116, 3040–3081.
- [15] Gunn, J. E., W. A. Siegmund, E. J. Mannery, and others (2006). The 2.5 m telescope of the Sloan Digital Sky Survey. *AJ* 131, 2332–2359.
- [16] Hastie, T., R. Tibshirani, and J. Friedman (2009). *The Elements of Statistical Learning: Data Mining, Inference, and Prediction*. Springer.
- [17] Hobson, M. P., A. H. Jaffe, A. R. Liddle, P. Mukherjee, and D. Parkinson (2010). *Bayesian Methods in Cosmology*. Cambridge: University Press.
- [18] Ivezić, Ž., R. H. Lupton, M. Jurić, and others (2002). Color confirmation of asteroid families. *AJ* 124, 2943–2948.
- [19] Ivezić, Ž., J. A. Smith, G. Miknaitis, and others (2007). Sloan Digital Sky Survey Standard Star Catalog for Stripe 82: The dawn of industrial 1% optical photometry. *AJ* 134, 973–998.
- [20] Jaynes, E. T. (2003). *Probability Theory: The Logic of Science*. Cambridge University Press.
- [21] Lee, Y. S., T. C. Beers, C. Allende Prieto, and others (2011). The SEGUE Stellar Parameter Pipeline. V. estimation of alpha-element abundance ratios from low-resolution SDSS/SEGUE stellar spectra. *AJ* 141, 90.
- [22] Loredo, T. and the Astro/Info Working Group (2009). The astronomical information sciences: A keystone for 21st-century astronomy. Position paper for the Astro2010 Decadal Survey, # 34.
- [23] Lupton, R. (1993). *Statistics in Theory and Practice*. Princeton University Press.
- [24] Lupton, R. H., J. E. Gunn, Ž. Ivezić, and others (2001). The SDSS imaging pipelines. In F. R. Harnden Jr., F. A. Primini, and H. E. Payne (Ed.), *Astronomical Data Analysis Software and Systems X*, Volume 238 of *Astronomical Society of the Pacific Conference Series*, pp. 269.
- [25] MacKay, D. J. C. (2010). *Information Theory, Inference, and Learning Algorithms*. Cambridge: University Press.
- [26] Parker, A., Ž. Ivezić, M. Jurić, and others (2008). The size distributions of asteroid families in the SDSS Moving Object Catalog 4. *Icarus* 198, 138–155.

- [27] Press, W. H., S. A. Teukolsky, W. T. Vetterling, and B. P. Flannery (1992). *Numerical Recipes in FORTRAN. The Art of Scientific Computing*. Cambridge: University Press.
- [28] Schlegel, D. J., D. P. Finkbeiner, and M. Davis (1998). Maps of dust infrared emission for use in estimation of reddening and cosmic microwave background radiation foregrounds. *ApJ* 500, 525.
- [29] Schneider, D. P., G. T. Richards, P. B. Hall, and others (2010). The Sloan Digital Sky Survey Quasar Catalog. V. Seventh Data Release. *AJ* 139, 2360–2373.
- [30] Sesar, B., J. S. Stuart, Ž. Ivezić, and others (2011). Exploring the variable sky with LINEAR. I. photometric recalibration with the Sloan Digital Sky Survey. *AJ* 142, 190.
- [31] Sivia, D. S. (2006). *Data Analysis: A Bayesian Tutorial*. Oxford University Press.
- [32] Skrutskie, M. F., R. M. Cutri, R. Stiening, and others (2006). The Two Micron All Sky Survey (2MASS). *AJ* 131, 1163–1183.
- [33] Smolčić, V., Ž. Ivezić, M. Gaćeša, and others (2006). The rest-frame optical colours of 99000 Sloan Digital Sky Survey galaxies. *MNRAS* 371, 121–137.
- [34] Stoughton, C., R. H. Lupton, M. Bernardi, and others (2002). Sloan Digital Sky Survey: Early Data Release. *AJ* 123, 485–548.
- [35] Strateva, I., Ž. Ivezić, G. R. Knapp, and others (2001). Color separation of galaxy types in the Sloan Digital Sky Survey imaging data. *AJ* 122, 1861–1874.
- [36] Strauss, M. A., D. H. Weinberg, R. H. Lupton, and others (2002). Spectroscopic target selection in the Sloan Digital Sky Survey: The main galaxy sample. *AJ* 124, 1810–1824.
- [37] Tufte, E. R. (2009). *The Visual Display of Quantitative Information*. Cheshire, Connecticut: Graphics Press.
- [38] Wall, J. V. and C. R. Jenkins (2003). *Practical Statistics for Astronomers*. Cambridge University Press, 2003.
- [39] Wasserman, L. (2010a). *All of Nonparametric Statistics*. Springer.
- [40] Wasserman, L. (2010b). *All of Statistics: A Concise Course in Statistical Inference*. Springer.
- [41] Way, M., J. Scargle, K. Ali, and A. Srivastava (2012). *Advances in Machine Learning and Data Mining for Astronomy*. Chapman and Hall/CRC Data Mining and Knowledge Discovery Series. Taylor and Francis.
- [42] York, D. G., J. Adelman, J. E. Anderson, Jr., and others (2000). The Sloan Digital Sky Survey: Technical summary. *AJ* 120, 1579–1587.