

THE ANALYSIS OF DATA

VOLUME 1

Probability

Guy Lebanon

First Edition

2012

Probability. The Analysis of Data, Volume 1.

First Edition, First Printing, 2013

<http://theanalysisofdata.com>

Copyright ©2013 by Guy Lebanon. All rights reserved. This book or any portion thereof may not be reproduced or used in any manner whatsoever without the express written permission of the author.

In memory of Alex Lebanon

Contents

| | |
|--|-----------|
| Contents | 5 |
| Preface | 9 |
| Mathematical Notations | 14 |
| 1 Basic Definitions | 19 |
| 1.1 Sample Space and Events | 19 |
| 1.2 The Probability Function | 21 |
| 1.3 The Classical Probability Model on Finite Spaces | 24 |
| 1.4 The Classical Model on Continuous Spaces | 27 |
| 1.5 Conditional Probability and Independence | 28 |
| 1.6 Basic Combinatorics for Probability | 32 |
| 1.7 Probability and Measure Theory* | 38 |
| 1.8 Notes | 39 |
| 1.9 Exercises | 39 |
| 2 Random Variables | 41 |
| 2.1 Basic Definitions | 41 |
| 2.2 Functions of a Random Variable | 49 |
| 2.2.1 Discrete $g(X)$ | 50 |
| 2.2.2 Continuous $g(X)$ | 51 |
| 2.3 Expectation and Variance | 53 |
| 2.4 Moments and the Moment Generating Function | 56 |
| 2.5 Random Variables and Measurable Functions* | 57 |
| 2.6 Notes | 58 |
| 2.7 Exercises | 59 |
| 3 Important Random Variables | 61 |
| 3.1 The Bernoulli Trial Random Variable | 61 |
| 3.2 The Binomial Random Variable | 62 |
| 3.3 The Geometric Random Variable | 64 |
| 3.4 The Hypergeometric Random Variable | 66 |
| 3.5 The Negative Binomial Random Variable | 68 |

| | | |
|----------|--|------------|
| 3.6 | The Poisson Random Variable | 70 |
| 3.7 | The Uniform Random Variable | 73 |
| 3.8 | The Exponential Random Variable | 75 |
| 3.9 | The Gaussian Random Variable | 77 |
| 3.10 | The Gamma and χ^2 Distributions | 80 |
| 3.11 | The t Random Variable | 84 |
| 3.12 | The Beta Random Variable | 85 |
| 3.13 | Mixture Random Variable | 87 |
| 3.14 | The Empirical Random Variable | 90 |
| 3.15 | The Smoothed Empirical Random Variable | 91 |
| 3.16 | Notes | 95 |
| 3.17 | Exercises | 96 |
| 4 | Random Vectors | 97 |
| 4.1 | Basic Definitions | 97 |
| 4.2 | Joint Pmf, Pdf, and Cdf Functions | 101 |
| 4.3 | Marginal Random Vectors | 102 |
| 4.4 | Functions of a Random Vector | 104 |
| 4.5 | Conditional Probabilities and Random Vectors | 108 |
| 4.6 | Moments | 110 |
| 4.7 | Conditional Expectations | 116 |
| 4.8 | Moment Generating Function | 117 |
| 4.9 | Random Vectors and Independent σ -Algebras* | 118 |
| 4.10 | Notes | 119 |
| 4.11 | Exercises | 119 |
| 5 | Important Random Vectors | 121 |
| 5.1 | The Multinomial Random Vector | 121 |
| 5.2 | The Multivariate Normal Random Vector | 125 |
| 5.3 | The Dirichlet Random Vector | 131 |
| 5.4 | Mixture Random Vectors | 137 |
| 5.5 | The Exponential Family Random Vector | 139 |
| 5.6 | Notes | 141 |
| 5.7 | Exercises | 141 |
| 6 | Random Processes | 143 |
| 6.1 | Basic Definitions | 143 |
| 6.2 | Random Processes and Marginal Distributions | 146 |
| 6.3 | Moments | 149 |
| 6.4 | One Dimensional Random Walk | 151 |
| 6.5 | Random Processes and Measure Theory* | 156 |
| 6.6 | The Borel-Cantelli Lemmas and the Zero-One Law* | 157 |
| 6.7 | Notes | 160 |
| 6.8 | Exercises | 161 |

| | | |
|----------|---|------------|
| 7 | Important Random Processes | 163 |
| 7.1 | Markov Chains | 163 |
| 7.1.1 | Basic Definitions | 163 |
| 7.1.2 | Examples | 166 |
| 7.1.3 | Transience, Persistence, and Irreducibility | 172 |
| 7.1.4 | Persistence and Transience of the Random Walk Process | 175 |
| 7.1.5 | Periodicity in Markov Chains | 177 |
| 7.1.6 | The Stationary Distribution | 178 |
| 7.2 | Poisson Processes | 182 |
| 7.2.1 | Postulates and Differential Equation | 182 |
| 7.2.2 | Relationship to Poisson Distribution | 184 |
| 7.2.3 | Relationship to Exponential Distribution | 184 |
| 7.2.4 | Relationship to Binomial Distribution | 186 |
| 7.2.5 | Relationship to Uniform Distribution | 186 |
| 7.3 | Gaussian Processes | 186 |
| 7.3.1 | The Wiener Process | 191 |
| 7.4 | Notes | 192 |
| 7.5 | Exercises | 193 |
| 8 | Limit Theorems | 195 |
| 8.1 | Modes of Stochastic Convergence | 195 |
| 8.2 | Relationships Between Modes of Convergences | 196 |
| 8.3 | Dominated Convergence Theorem for Vectors* | 199 |
| 8.4 | Scheffe's Theorem | 200 |
| 8.5 | The Portmanteau Theorem | 201 |
| 8.6 | The Law of Large Numbers | 205 |
| 8.7 | The Characteristic Function* | 209 |
| 8.8 | Levy's Continuity Theorem | 211 |
| 8.9 | The Central Limit Theorem | 213 |
| 8.10 | Continuous Mapping Theorem | 217 |
| 8.11 | Slustky's Theorems | 218 |
| 8.12 | Notes | 219 |
| 8.13 | Exercises | 220 |
| A | Set Theory | 223 |
| A.1 | Basic Definitions | 223 |
| A.2 | Functions | 229 |
| A.3 | Cardinality | 230 |
| A.4 | Limits of Sets | 233 |
| A.5 | Notes | 235 |
| A.6 | Exercises | 235 |

| | | |
|----------|--|------------|
| B | Metric Spaces | 237 |
| | B.1 Basic Definitions | 237 |
| | B.2 Limits | 239 |
| | B.3 Continuity | 244 |
| | B.4 The Euclidean Space | 245 |
| | B.5 Growth of Functions | 254 |
| | B.6 Notes | 254 |
| | B.7 Exercises | 254 |
| C | Linear Algebra | 257 |
| | C.1 Basic Definitions | 257 |
| | C.2 Rank | 265 |
| | C.3 Eigenvalues, Determinant, and Trace | 266 |
| | C.4 Positive Semi-Definite Matrices | 272 |
| | C.5 Singular Value Decomposition | 274 |
| | C.6 Notes | 276 |
| | C.7 Exercises | 276 |
| D | Differentiation | 279 |
| | D.1 Univariate Differentiation | 279 |
| | D.2 Taylor Expansion and Power Series | 284 |
| | D.3 Notes | 287 |
| | D.4 Exercises | 288 |
| E | Measure Theory* | 289 |
| | E.1 σ -Algebras* | 289 |
| | E.2 The Measure Function* | 292 |
| | E.3 Caratheodory's Extension Theorem* | 293 |
| | E.3.1 Dynkin's Theorem* | 293 |
| | E.3.2 Outer Measure* | 295 |
| | E.3.3 The Extension Theorem* | 298 |
| | E.4 Independent σ -Algebras* | 299 |
| | E.5 Important Measure Functions* | 300 |
| | E.5.1 Discrete Measure Functions* | 301 |
| | E.5.2 The Lebesgue Measure* | 301 |
| | E.5.3 The Lebesgue-Stieltjes Measure* | 303 |
| | E.6 Measurability of Functions* | 304 |
| | E.7 Notes | 305 |
| F | Integration | 307 |
| | F.1 Riemann Integral | 307 |
| | F.2 Integration and Differentiation | 311 |
| | F.3 The Lebesgue Integral* | 315 |
| | F.3.1 Relation between the Riemann and the Lebesgue Integrals* | 325 |
| | F.3.2 Transformed Measures* | 326 |
| | F.4 Product Measures* | 327 |

| | |
|---|------------|
| <i>CONTENTS</i> | 9 |
| F.5 Integration over Product Spaces* | 330 |
| F.5.1 The Lebesgue Measure over \mathbb{R}^{d*} | 332 |
| F.6 Multivariate Differentiation and Integration | 333 |
| F.7 Notes | 336 |
| Bibliography | 337 |
| Index | 341 |

Preface

The Analysis of Data Project

The Analysis of Data (TAOD) project provides educational material in the area of data analysis.

- The project features comprehensive coverage of all relevant disciplines including probability, statistics, computing, and machine learning.
- The content is almost self-contained and includes mathematical prerequisites and basic computing concepts.
- The R programming language is used to demonstrate the contents. Full code is available, facilitating reproducibility of experiments and letting readers experiment with variations of the code.
- The presentation is mathematically rigorous, and includes derivations and proofs in most cases.
- HTML versions are freely available on the website <http://theanalysisofdata.com>. Hardcopies are available at affordable prices.

Volume 1: Probability

This volume focuses on probability theory. There are many excellent textbooks on probability, and yet this book differs from others in several ways.

- Probability theory is a wide field. This book focuses on the parts of probability that are most relevant for statistics and machine learning.
- The book contains almost all of the mathematical prerequisites, including set theory, metric spaces, linear algebra, differentiation, integration, and measure theory.
- Almost all results in the book appear with a proof.

- Probability textbooks are typically either elementary or advanced. This book strikes a balance by attempting to avoid measure theory where possible, but resorting to measure theory and other advanced material in a few places where they are essential.
- The book uses R to illustrate concepts. Full code is available in the book, facilitating reproducibility of experiments and letting readers experiment with variations of the code.

I am not aware of a single textbook that covers the material from probability theory that is necessary and sufficient for an in-depth understanding of statistics and machine learning. This book represents my best effort in that direction.

Since this book is part of a series of books on data analysis, it does not include any statistics or machine learning. Such content is postponed to future volumes.

Website

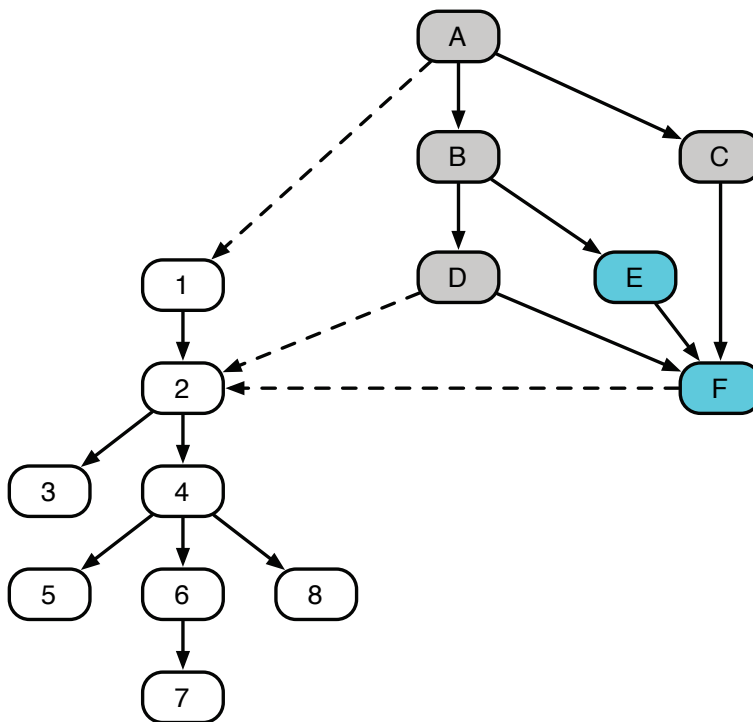
A companion website (<http://theanalysisofdata.com>) contains an HTML version of this book, errata, and additional multimedia material. The website will also link to additional TAOD volumes as they become available.

Mathematical Appendices

A large part of the book contains six appendices on mathematical prerequisites. Probability requires knowledge of many branches of mathematics, including calculus, linear algebra, set theory, metric spaces, measure, and Lebesgue integration. Instead of referring the reader to a large collection of math textbooks we include here all of the necessary prerequisites. References are provided in the notes sections at the end of each chapter for additional resources.

Dependencies

The diagram below indicates the dependencies between the different chapters of the book. Appendix chapters are shaded and dependencies between appendix chapters and regular chapters are marked by dashed arrows. It is not essential to strictly adhere to this dependency as many chapters require only a brief familiarity with some issues in the chapters that they depend on.



Starred sections correspond to material that requires measure theory, or that can be better appreciated with knowledge of measure theory.

R Code

The book contains many fragments of R code, aimed to illustrate probability theory and its applications. The code is included so that the reader can reproduce the results as well as modify the code and run variations of it. In order to appreciate the code and modify it, the reader will need a basic understanding of the R programming language and R graphics. A good introduction to R is available from the CRAN website at <http://cran.r-project.org/doc/manuals/R-intro.pdf>. Alternatively, two chapters from volume 2 of the analysis of data series (R Programming and R Graphics) are freely available online at <http://theanalysisofdata.com>.

The book uses a variety of R graphics packages, including `base`, `lattice`, and `ggplot2`. The most frequently used package is `ggplot2`, which is described in detail in [49] (see also <http://had.co.nz/ggplot2/> and the R Graphics chapter available at <http://theanalysisofdata.com>.)

To ensure that the code fragments run locally as they appear in the book, install and bring into scope the required packages using the `install.packages` and `library` R functions. For example, to install and load the `ggplot2` package type the following commands in the R prompt.

```
# this is a comment
install.packages("ggplot2") # installs the package ggplot2
library("ggplot2") # brings ggplot2 into scope
```

The code fragments throughout the book are annotated with output displayed by the R interpreter. This output is displayed following two hash symbols which are interpreted as comments by R (see below).

```
a = pi
print(a) # note the ## symbols preceding the output below

## [1] 3.142

print(a + 1)

## [1] 4.142
```

This format makes it easy to copy a code fragment from an HTML page and paste it directly into R (the output following the hash symbols will be interpreted as comments and not produce a syntax error).

Acknowledgements

The following people made technical suggestions that helped improve the contents: Krishnakumar Balasubramanian, Rohit Banga, Joshua Dillon, Sanjeet Hajarnis, Oded Green, Yi Mao, Seungyeon Kim, Joonseok Lee, Fuxin Li, Nishant Mehta, Yaron Rachlin, Parikshit Ram, Kaushik Rangadurai, Neil Slagle, Mingxuan Sun, Brian Steber, Gena Tang, Long Tran. In addition, many useful comments were received through a discussion board during my fall 2011 class computational data analysis at Georgia Tech. These comments were mostly anonymized, but some commentators who identified themselves appear above.

Katharina Probst, Neil Slagle and Laura Usselman edited portions of this book, and made many useful suggestions. The book features a combination of text, equations, graphs, and R code, made possible by the `knitr` package. I thank Yihui Xie for implementing `knitr`, and for his help through the `knitr` Google discussion group. Katharina Probst helped with web development and design.

Mathematical Notations

Logic

| | |
|----------------------------|----------------|
| \forall | for all |
| \exists | there exists |
| \Rightarrow | implies |
| \Leftrightarrow | if and only if |
| $\stackrel{\text{def}}{=}$ | defined as |

For example, $\forall a > 0, 1/a > 0$ reads “for every $a > 0$ we have $1/a > 0$ ” and $\forall a > 0, \exists b > 0, a/b = 1 \Rightarrow b/a = 1$ reads “for all $a > 0$ there exists $b > 0$ such that, whenever $a/b = 1$, we also have $b/a = 1$ ”.

Sets and Functions (Chapter A)

| | |
|------------------|---|
| Ω | sample space |
| ω | an element of the sample space Ω |
| I_A | indicator function $I_A(x) = 1$ if $x \in A$ and 0 otherwise |
| δ_{ij} | Kronercker’s delta: $\delta_{ij} = 1$ if $i = j$ and 0 otherwise |
| \mathbb{N} | natural numbers $\{1, 2, 3, \dots\}$ |
| \mathbb{Z} | integers $\{\dots, -2, -1, 0, 1, 2, \dots\}$ |
| \mathbb{Q} | rational numbers |
| $A \times B$ | Cartesian product of two sets |
| A^k | repeated Cartesian product (k times) of a set: $A \times \dots \times A$ |
| A^c | complement of a set: $\Omega \setminus A$ |
| 2^A | power set of the set A |
| $ A $ | number of elements in a finite set |
| $f^{-1}(A)$ | pre-image of the function f : $\{x : f(x) \in A\}$ |
| $f \equiv g$ | $f(x) = g(x)$ for all x |
| $x^{(n)}$ | sequence of elements in a set |
| $A_n \nearrow A$ | convergence of a sequence of increasing sets |
| $A_n \searrow A$ | convergence of a sequence of decreasing sets |

Combinatorics (Section 1.6)

| | |
|----------------|--|
| $n!$ | factorial function $n(n-1) \cdots 2 \cdot 1$ |
| $\binom{n}{r}$ | $n!/(n-r)!$ |
| $\binom{n}{r}$ | $n!/(r!(n-r)!)$ |

Metric spaces (Chapter B)

| | |
|---|---|
| \mathbb{R} | real numbers |
| \mathbb{R}^d | set of d -dimensional vectors of real numbers |
| $B_a(\mathbf{x})$ | open ball of radius a centered at \mathbf{x} |
| \mathbf{x} | vector in \mathbb{R}^d |
| x_i | the i -component of the vector \mathbf{x} |
| $\mathbf{x}^{(n)}$ | sequence of vectors in \mathbb{R}^d |
| $x_i^{(j)}$ | the i -component of the vector $\mathbf{x}^{(j)}$ |
| $\langle \mathbf{x}, \mathbf{y} \rangle$ | inner product $\sum_i x_i y_i$ |
| $\ \mathbf{x}\ $ | Euclidean norm $\sqrt{\sum_i x_i^2}$ |
| $d(\mathbf{x}, \mathbf{y})$ | Euclidean distance $\sqrt{\sum_i (x_i - y_i)^2}$ |
| $\mathbf{x}^{(n)} \rightarrow \mathbf{x}$ | convergence of a sequence |
| $f_n \rightarrow f$ | pointwise convergence of a sequence of functions |
| $f_n \nearrow f$ | pointwise convergence of an increasing sequence of functions |
| $f_n \searrow f$ | pointwise convergence of a decreasing sequence of functions |
| $O(f)$ | big O growth notation |
| $o(f)$ | little o growth notation |

Note in particular that we denote vectors in bold face, for example \mathbf{x} , and the scalar components of such vectors using subscripts (non-bold face) $\mathbf{x} = (x_1, \dots, x_d)$. We refer to sequence of vectors using superscripts, for example $\mathbf{x}^{(n)}$ and their scalar components as $\mathbf{x}^{(n)} = (x_1^{(n)}, \dots, x_d^{(n)})$.

Probability (Chapters 1, 2, 4, 8)

| | |
|-------------------------------|--|
| P | probability function |
| P_E | probability conditioned on the event E |
| F_X | cumulative distribution function (cdf) corresponding to the RV X |
| f_X | probability density function (pdf) corresponding to the RV X |
| p_X | probability mass function (pmf) corresponding to the RV X |
| E | expectation |
| Var | variance |
| Cov | covariance |
| std | standard deviation |
| \sim | distributed according to |
| $\stackrel{\text{iid}}{\sim}$ | independent identically distributed (iid) sampling |
| \mathbb{P}_k | simplex of probability functions over $\Omega = \{1, \dots, k\}$ |
| i.o. | infinitely often |
| $\xrightarrow{\text{as}}$ | convergence with probability 1 |
| \xrightarrow{p} | convergence in probability |
| \rightsquigarrow | convergence in distribution |

Matrices (Chapter C)

| | |
|---------------------------|--|
| A^\top | matrix transpose |
| A^n | the matrix A raised to the n -power |
| A_{ij}^n | the i, j element of the matrix A^n |
| A^{-1} | inverse of matrix A |
| I | identity matrix |
| $\text{tr } A$ | trace of the matrix A |
| $\det A$ | determinant of the matrix A |
| $\text{row } A$ | row space of the matrix A |
| $\text{col } A$ | column space of the matrix A |
| $\dim A$ | dimension of a linear space |
| $\text{rank } A$ | rank of the matrix A |
| $\text{diag}(\mathbf{v})$ | diagonal matrix whose diagonal is given by the vector \mathbf{v} |
| $A \otimes B$ | Kronecker product of two matrices A, B |

All vectors are assumed to be column vectors unless stated explicitly otherwise. For example, if \mathbf{x} is an arbitrary vector and A a matrix the expression $\mathbf{x}^\top A \mathbf{x}$ represents a scalar.

Differentiation (Chapter D)

| | |
|--|---|
| df/dx | derivative |
| d^2f/dx^2 | second order derivative |
| $\partial f/\partial x_i$ | partial derivative |
| ∇f | gradient vector of partial derivatives of $f : \mathbb{R}^k \rightarrow \mathbb{R}$ |
| $\nabla \mathbf{f}$ | Jacobian matrix of partial derivatives of $\mathbf{f} : \mathbb{R}^k \rightarrow \mathbb{R}^m$ |
| $\partial^2 f/\partial x_i \partial x_j$ | second order partial derivative |
| $\nabla^2 f$ | Hessian matrix of second order partial derivatives of $f : \mathbb{R}^k \rightarrow \mathbb{R}$ |

We consider the gradient vector ∇f as a row vector. This ensures that the notations ∇f and $\nabla \mathbf{f}$ are consistent.

Measure and Integration (Chapters E, F)

| | |
|---------------------------------------|--|
| $\sigma(\mathcal{C})$ | σ algebra generated by the set of sets \mathcal{C} |
| μ | measure |
| $\int d\mu$ | Lebesgue integral with respect to measure μ |
| $\int d\mathbf{P}$ | Lebesgue integral with respect to probability measure \mathbf{P} |
| $\int dF_X$ | Lebesgue integral with respect to probability measure corresponding to cdf F_X |
| $\int dx$ | Lebesgue integral with respect to Lebesgue measure or Riemann integral |
| $\mathcal{B}(\mathbb{R}^d)$ | Borel σ -algebra over metric space \mathbb{R}^d |
| $\mathcal{F}_1 \otimes \mathcal{F}_1$ | product σ -algebra |
| $\mu_1 \times \mu_2$ | product measure |
| a.e. | almost everywhere (except on a set of measure zero) |

